

TABLE OF CONTENTS

Contributed Articles

- 1 Attribution and Obfuscation of Neural Text Authorship: A Data Mining Perspective.
Adaku Uchendu, Thai Le, and Dongwon Lee
- 19 The Need for Unsupervised Outlier Model Selection: A Review and Evaluation of Internal Evaluation Strategies.
Martin Q. Ma, Yue Zhao, Xiaorong Zhang, and Leman Akoglu
- 36 Stop Using the Elbow Criterion for k-means, and How to Choose the Number of Clusters Instead.
Erich Schubert
- 43 Did You Train on My Dataset? Towards Public Dataset Protection with Clean-Label Backdoor Watermarking.
Ruixiang Tang, Qizhang Feng, Ninghao Liu, Fan Yang, and Xia Hu
- 54 Privacy-Preserving Graph Machine Learning from Data to Computation: A Survey.
Dongqi Fu, Wenxuan Bao, Ross Maciejewski, Hanghang Tong, and Jingrui He
- 73 Adaptive Risk-Aware Bidding with Budget Constraint in Display Advertising.
Zhimeng Jiang, Kaixiong Zhou, Mi Zhang, Rui Chen, Xia Hu, and Soo-Hyun Choi

Editor-in-Chief:
Xiangliang Zhang

Associate Editors:
Brian Davison
Jiayu Zhou
Srijan Kumar
<http://www.kdd.org/explorations/>



**Association for
Computing Machinery**

Advancing Computing as a Science & Profession

Attribution and Obfuscation of Neural Text Authorship: A Data Mining Perspective

Adaku Uchendu
Penn State University
PA, USA
azu5030@psu.edu

Thai Le
University of Mississippi
MS, USA
thaile@olemiss.edu

Dongwon Lee
Penn State University
PA, USA
dongwon@psu.edu

ABSTRACT

Two interlocking research questions of growing interest and importance in privacy research are *Authorship Attribution* (AA) and *Authorship Obfuscation* (AO). Given an artifact, especially a text t in question, an AA solution aims to accurately attribute t to its true author out of many candidate authors while an AO solution aims to modify t to hide its true authorship. Traditionally, the notion of authorship and its accompanying privacy concern is only toward *human* authors. However, in recent years, due to the explosive advancements in Neural Text Generation (NTG) techniques in NLP, capable of synthesizing human-quality open-ended texts (so-called “neural texts”), one has to now consider authorships by humans, machines, or their combination. Due to the implications and potential threats of neural texts when used maliciously, it has become critical to understand the limitations of traditional AA/AO solutions and develop novel AA/AO solutions in dealing with neural texts. In this survey, therefore, we make a comprehensive review of recent literature on the attribution and obfuscation of neural text authorship from a Data Mining perspective, and share our view on their limitations and promising research directions.

1. INTRODUCTION

Natural Language Generation (NLG) is a broad term for AI techniques to produce high-quality human-understandable texts in some human languages, and often encompasses terms such as machine translation, dialogue generation, text summarization, data-to-text generation, Question-Answer generation, and open-ended or story generation [72, 119]. Among these, in particular, this survey focuses on the open-ended text generation aspect of NLG. Since the advent of the Transformers architecture in 2018, the field of NLG has experienced exponential improvement. Before 2018, leading NLG models were only able to generate a few sentences coherently. However, after adopting the Transformer architecture into deep learning-based Language models (LMs), NLG models could generate more than a few sentences (i.e., ≥ 200 words) coherently. GPT-1 [92] by OpenAI is one of the first such NLG models. Since then, many other Transformer-based LMs with the capacity to generate long coherent texts have been

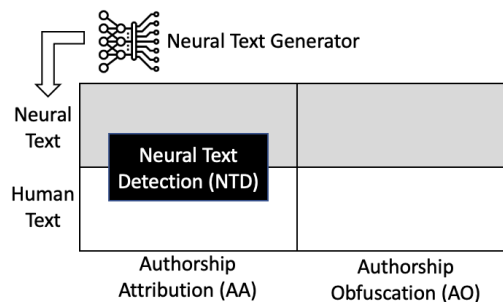


Figure 1: The figure illustrates the quadrant of research problems where (1) the **GRAY** quadrants are the focus of this survey, and (2) The **BLACK** box indicates the specialized binary AA problem to distinguish neural texts from human texts.

released (e.g., FAIR [16, 82], CTRL [59], PPLM [25], T5 [94], Wu-Dao ¹). In fact, as of February 2023, huggingface’s [113] model repo houses about 8,300 variants of text-generative LMs². In this survey, we refer to these LMs as **Neural Text Generator (NTG)** since they are neural network-based LMs with text-generative abilities. Further, we refer to the texts generated by NTG as “**neural**” **texts**³, as opposed to normal texts written by humans as **human texts**.

As the qualities of NTGs improve, neural texts become more easily misconstrued as human-written [18, 45, 52, 108, 117], exacerbating the difficulty of distinguishing neural texts from human texts. For instance, therefore, such a text generation capability can be misused to generate misinformation [13, 14, 117], fake reviews [3] and political propaganda [109] at scale with little cost. These problems lead to the need to effectively distinguish neural texts from human texts, the so-called **Neural Text Detection (NTD)** problem, which is a sub-problem of a widely studied problem in the privacy community—i.e., authorship attribution. In fact, two interlocking research questions in privacy research, heavily studied but of growing interest, are **Authorship Attribution (AA)** and **Authorship Obfuscation (AO)**. Given an artifact, especially a text t in question, an AA solution aims to accurately at-

¹<https://github.com/BAAI-WuDao>

²https://huggingface.co/models?pipeline_tag=text-generation

³Other names for neural text include *AI(-generated) text* [66], *Machine(-generated/written) text* [1, 5, 36, 38, 45, 89, 96, 104, 107, 108, 117], *Artificial text* [67], *Computer-generated text* [102], *Deepfake text* [91], *Auto-generated text* [85], and *Synthetic text* [13, 27, 46, 81].

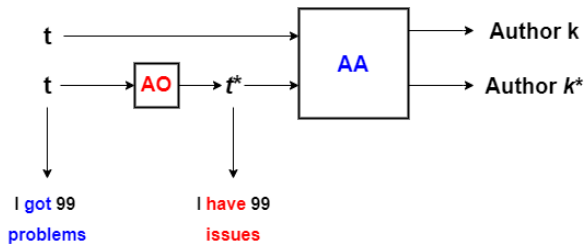


Figure 2: Illustration of both AA and AO problems on neural texts

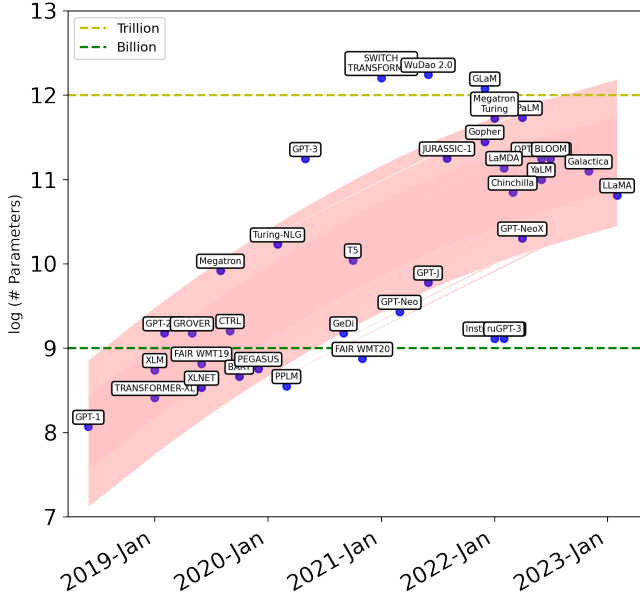


Figure 3: Evolution of Neural Text Generators (NTGs) from 2018 to 2023 (Y-axis is a log plot of # of parameters).

tribute t to its true author out of k candidate authors while an AO solution aims to modify t to hide its true authorship. Therefore, NTD is a specialized case, a Turing Test, of AA with $k = 2$ authors (i.e., human vs. machine). Figure 1 illustrates the quadrant of research problems, while Figure 2 illustrates how both AA and AO problems work hand-in-hand.

Traditionally, the notion of authorship and its accompanying privacy concern in both AA and AO problems are only toward *human* authors. However, with the arrival of generative AI technologies and due to the potential threats of misused neural texts, one now has to consider authorships by humans, machines, or their combination, and re-thinks about effective solutions for both AA and AO problems for neural texts. Hence, to guide these developments, in this survey, we provide a detailed analysis of both AA and AO problems, their existing solutions, and our perspective on the open challenges. As both AA and AO problems are essentially computational learning problems, we discuss the landscape from *A Data Mining Perspective* and call attention to the security challenges that need to be solved. We believe that the issues of these novel AA/AO problems for neural texts are “nuanced” and therefore require nuanced solutions from the Data Mining and Machine Learning community.

2. NEURAL TEXT GENERATION

We first select a handful of Neural Text Generator (NTG) that we focus on in this survey, and introduce a list of popular datasets with neural texts.

2.1 Neural Text Generators (NTGs)

Those NTGs studied in this survey are large-scale probabilistic LMs that are capable of generating long-coherent texts (e.g., ≥ 200 words). These LMs are trained on massive amounts of unstructured texts. Based on its architecture structure (e.g., encoder-decoder or decoder only), these LMs use a prompt, a snippet of human-written text, to guide the generation of texts, emulating the most similar style from the training set and predicting one token at a time. Recent works such as [72, 119] survey these NTGs in detail. The progress of NTGs in recent years has been expeditious. As shown in Figure 3, for instance, the sizes of NTGs with respect to their parameters are growing at an exponential rate, yielding the rapid improvement in the quality of neural texts, thus exacerbating the AA/AO problems. Table 1 shows a summary of state-of-the-art NTGs, where many entries are drawn from [107, 108].

Within LMs, in particular, hyperparameters matter a great deal when generating texts. The choice of these hyperparameters, referred to as *decoding strategies*, greatly affects the quality of generated neural texts. According to [50], there are 6 *decoding strategies*: (1) *Greedy sampling* selects the best probable word, (2) *Random sampling* does a stochastic search for a sufficient word, (3) *Top-k sampling* samples from top-k most probable words, (4) *Beam search* searches for most probable candidate sequences, (5) *Nucleus (Top-p) sampling* samples similar to top-k, but its focus is on the smallest possible set of top words, such that the sum of their probabilities is $\geq p$, and (6) *Temperature* scales logits to either increase or decrease the entropy of sampling. NTG models often use *Top-k sampling*, *Beam search*, *Nucleus (Top-p) sampling*, and *Temperature* decoding strategies as they produce higher quality texts than other decoding strategies. In fact, [50] reports that neural texts generated with the *Nucleus (Top-p) sampling* strategy are more challenging to attribute authorships.

2.2 Neural Text Datasets

To investigate both AA/AO problems for neural texts, one needs benchmark datasets of neural texts. Table 2 describes a list of publicly available datasets that contain neural texts. Labels of the texts in the datasets are either binary (neural vs. human text) or multi-class (having multiple neural texts generated by different NTGs vs. 1 human label). The majority of recent studies have focused on the binary case of the Turing Test to check if a given text is written by a human author or machine (i.e., one of NTGs). Researchers utilized clever labeling and generation methods to build datasets: (1) *Binary dataset*: Researchers first collect human-written texts (e.g., news, blogs, stories, or recipes) and use snippets of these texts as prompts to the chosen NTG to generate a machine-written (neural) text. (2) *Multi-class dataset*: Starting with human-written texts as prompts, this generates multiple neural texts by using different NTG architectures (i.e., human vs. GPT-1, GROVER, PPLM, etc.), different pre-trained sizes of the same NTG architecture (i.e., human vs. GPT-2 small vs. GPT-2 medium vs. GPT-2 large vs. GPT-2 XL, etc.), different decoding strategies (i.e., human vs. GPT-2 top-k vs. GPT-2 top-p, etc.).

NTG	Author	Description
GPT-1 [92]	OpenAI	It used Transformers to model a simple concept - to predict the next token, given the previous token.
GPT-2 [93]	OpenAI	GPT-1 scaled up. There are 4 GPT-2 pre-trained models - <i>small</i> (124 million parameters), <i>medium</i> (355 million parameters), <i>large</i> (774 million parameters), and <i>x-large</i> (1558 million parameters)
GPT-3 [13]	OpenAI	GPT-2, scaled up - increasing parameter and train data size.
GROVER [117]	AllenAI	Similar to GPT-2 architecture and trained to generate political news. There are 3 pre-trained models: <i>GROVER-base</i> , <i>GROVER-large</i> , <i>GROVER-mega</i>
CTRL [59]	Salesforce	Conditional Transformer LM For controllable generation uses control codes to guide generation
XLM [20]	Facebook	A Cross-lingual Language Model trained on various languages. Only the English model is used for AA
XLNET [116]	Google	A generalized auto-regressive pre-training method that adopts the Transformer-XL framework
FAIR_wmt [16, 82]	Facebook	FAIR_wmt has 3 language models - English, Russian, and German. Only the English model ⁴ is used, which has 2 models - <i>WMT19</i> [82] and <i>WMT20</i> [16].
TRANSFORMER_XL [24]	Google	Another Transformer model that learns long-term dependency to improve long coherent text generation
PPLM [25]	Uber	The Plug and Play Language Models (PPLM) model upon GPT-2 by fusing the GPT-2 medium with a bag of words (BoW) models. These BoW models are <i>legal</i> , <i>military</i> , <i>monsters</i> , <i>politics</i> , <i>positive_words</i> , <i>religion</i> , <i>science</i> , <i>space</i> , <i>technology</i> . PPLM can plug in any GPT-2 pre-trained model to generate texts
Switch Transformer [35]	Google	Google uses a switch Transformer to build a sparse neural LM with 1.6T parameters are built
GPT-Neo [42]	EleutherAI	EleutherAI replicates GPT-3's architecture. There a 2 model sizes - 1.3B and 2.7B parameters
GPT-NeoX [12]	EleutherAI	A 20 billion parameter autoregressive replication of GPT-3.
GPT-J [112]	EleutherAI	A 6B parameter model similar to the GPT-Neo and GPT-NeoX that uses Mesh Transformer JAX [111] framework to train the model with Pile ⁵ dataset, a large curated dataset created by EleutherAI
T5 [94]	Google	An encoder-decoder text-to-text Transformer-based model. T5 has 5 pre-trained models - <i>T5-small</i> , <i>T5-base</i> , <i>T5-large</i> , <i>T5-3b</i> , and <i>T5-11b</i>
BART [71]	Facebook	This is another encoder-decoder Transformer-based LM, most effective when fine-tuned
PaLM [17]	Google	PaLM stands for Pathways Language Model. It is a dense decoder-only Transformer-based model trained with [7]'s pathways system framework
OPT-175B [121]	Meta	Meta's response to GPT-3. OPT-175B uses a similar framework to GPT-3 but the training costs 1/7th the carbon footprint of GPT-3.
GeDi [63]	Salesforce	GeDi stands for Generative Discriminator Guided Sequence Generation. Similar to PPLM, GeDi controls text generation using small LMs as generative discriminators

Table 1: Description of state-of-the-art Neural Text Generators (NTGs).

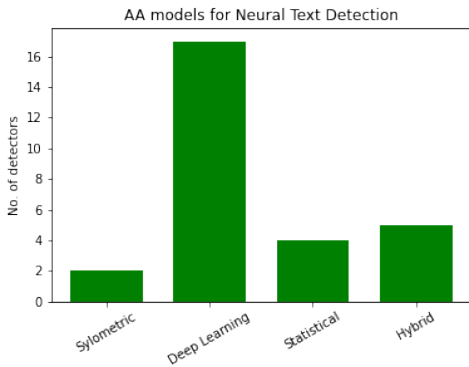


Figure 4: Number of AA solutions for NTD per category in the Taxonomy

3. AUTHORSHIP ATTRIBUTION FOR NEURAL TEXTS

Traditional AA problem studies the attribution of an author to a piece of written text out of a number of possible authors. However, in the literature, researchers have also studied a few variations of the AA problem. For instance, the *Author Verification (AV)* problem [8, 60, 60, 100, 101, 105] studies if the given two texts, t_1 and t_2 , are written by the same author? With the rise of neural texts, in addition, a specialized case of AA problem, NTD [5, 44, 51, 104, 107, 117], studies: By and large, however, a good solution for the standard AA problem can lead to a good solution for other variations of the AA problem. As such, we focus on the survey of the standard AA problem for neural texts. Thus, our paper formally defines the AA task as follows.

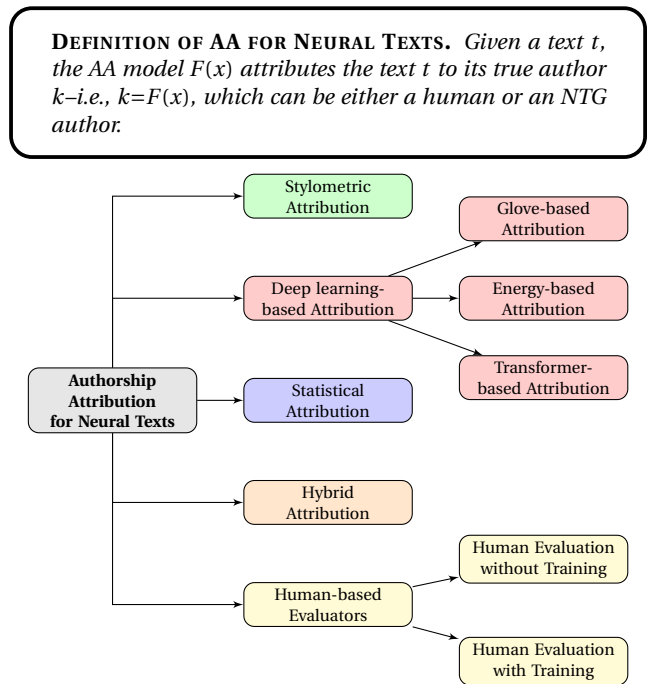


Figure 5: Taxonomy of Authorship Attribution models for NTD

In the following, we survey recent AA solutions that are capable of handling neural texts in different ways, as illustrated in Figure 5: *Stylometric Attribution*, *Deep Learning-based Attribution*, *Statistical Attribution*, and *Hybrid Attribution*.

3.1 Stylometric Attribution

Name	Description	Category	Domain	Labels
GPT-2 dataset [92]	250K Webtext (Human dataset) vs. 250K GPT-2 (small, medium, large, & XL)	Binary	News	GPT-2 & Human
GROVER dataset [117]	Using April 2019 news articles as the prompt, GROVER-Mega generated news articles	Binary	News	GROVER & Human
TuringBench-AA [108]	Used 10K human-written news articles (mostly Politics) from CNN, etc. to generate 10K articles each from 19 NTGs.	Multi-class	News	Human & 19 Machine labels (GPT-1, GPT-2 small, etc.)
TuringBench-TT [108]	The same dataset as <i>TuringBench-AA</i> , except that the datasets are 19 versions of human vs. each of the 19 NTGs.	Binary	News	19 Human vs. Machine combinations (GPT-2, etc.)
Authorship Attribution-AA [107]	Used 1K human-written articles to generate 1K articles each from 8 Artificial Text Generators	Multi-class	News	1 human vs 8 Machine labels (GPT-1, GPT-2, etc.)
Authorship Attribution-TT [107]	The same dataset as <i>Authorship Attribution-AA</i> except that the datasets are 8 versions of human vs. each of the 8 NTGs	Binary	News	8 humans vs Neural combinations
Authorship Attribution-TT mix [107]	The same dataset as <i>Authorship Attribution-AA</i> except that the dataset is human vs. Machine (which is a mixture of all the 8 NTGs)	Binary	News	1 human vs Machine (8 different machines)
Academic Papers & Abstracts [75]	2 datasets - (1) FULL: using a short prompt for a human-written paper generated an academic paper using GPT-2; (2) PARTIAL: Replacing sentences from an Abstract with Arxiv-NLP model generations	Binary	Academic Abstracts	Human & Machine
Hybrid Human-Machine Text [23]	Using human-written text in domains - News, Reddit, and Recipes to generate continuations of the text using GPT-2 XL	Binary	News, Reddit, Recipes	Human prompt & Machine texts
Amazon Reviews [3]	Fine-tuned GPT-2 on 3.6 M Amazon and 560K Yelp reviews	Binary	Reviews	Human & Machine
Human-Machine Pairs [90]	Generated texts with GROVER mega and GPT-2 XL with top-p decoding strategy. Paired human-written texts with a similar neurally generated version	Binary (Human-Machine pairs)	Online forums & News	Human & Machine
NeuralNews dataset [104]	Using the GoodNews [11] dataset as the human-written prompt to generate texts with GROVER. Real images are included in each of the articles.	Binary	News	Human & Machine
SynSciPass [95]	Built dataset using 3 potential sources of neural text: (1) open-ended text generators like GPT-2 & BLOOM (2) paraphrase models like SCITgen and PEGASUS and (3) translation models like Spinbot, real, Google translate, and Opus.	Multi-class	scientific articles	generate, translate, paraphrase, & human
TweepFake [32]	Collected tweets generated by Twitter bots and grouped them into tweets generated by GPT-2, RNN, and other bots	Binary	Tweets	Human & Machine

Table 2: Datasets with neural texts

Stylometry is the statistical analysis of the style of written texts. In traditional AA, stylometric classifiers are built using classical machine learning models trained on ensembles of style-based features such as N -grams, Part-of-Speech (POS), WritePrints [1], LIWC (Linguistic Inquiry & Word Count)[87], Readability score, and Empath [34]. This has been shown to be a successful approach for traditional AA tasks [68]. Due to such success, these models have been adopted and customized to the task of NTT.

The first attempt at a stylometric classifier to solve the AA task for $k > 2$ authors is the *Linguistic model* proposed by [107]. It trains a Random Forest classifier with the Authorship Attribution-AA dataset in Table 2 and extracts an ensemble of stylometric features (e.g., *entropy*, *readability score*, & *LIWC* (Linguistic Inquiry & Word Count) [87]). The entropy feature counts the number of unique characters in the text. Readability scores represent the estimated educational level of the author of a piece of text based on lexicon usage. LIWC is a psycho-linguistic dictionary that counts the frequency of words that represents a psychological emotion or linguistic structure [87]. This *Linguistic model* achieves a 90% F1 score and outperforms all the other deep learning-based models. However, this superior performance is a result of the small size of the dataset (only about 1k per data label) [107]. Scaling up the data size in terms of labels and examples will make the

AA task harder, and therefore cause the *Linguistic model* to underperform. This claim is confirmed in their second work using the TuringBench dataset [108]. They compared SOTA deep-learning-based models (BERT and RoBERTa) with several stylometric classifiers - SVM (3-grams), WriteprintsRFC, Random Forest (w/ TF-IDF), Syntax-CNN, Ngram CNN, and N-gram LSTM-LSTM. RoBERTa outperforms all the stylometric classifiers with about a 10-22% increase in F1 scores.

To further explore the benefits of stylometric features leveraged in the traditional AA community, [36] proposes a clever way to use them. This solution aims to solve the special case of AA, *Turing Test* (TT). First, they identify different issues with NTGs which can be captured by specific types of stylometric features. These issues in neural texts are categorized into 4 types: (1) *Lack of syntactic and lexical diversity* which can be captured with Named Entity-tags, POS-tags, and neuralcoref extension⁶ (a tool for using a neural network to annotate and resolve coreference clusters) to detect coreference clusters; (2) *Repetitiveness of words* which can be captured by collecting the number of stop-words, unique words, and words from “top-lists” of total words in a text. Also, a “conjunction overlap” measure is defined to calculate the

⁶<https://github.com/huggingface/neuralcoref>

Model Name	Classifier Type	Category	Learning Type	Interpretable	Training dataset
GROVER detector [117]	DL (Transformer-based)	Binary	Supervised		GROVER
GLTR [45]	Statistical	Binary	Unsupervised	✓	GPT-2
GPT-2 detector [51]	DL (Transformer-based)	Binary	Supervised		GPT-2
OpenAI detector [51]	DL (Transformer-based)	Multi-class	Supervised		GPT-2 & TuringBench-AA
RoBERTa-TT [108]	DL (Transformer-based)	Binary	Supervised		TuringBench-TT
BERT-TT [108]	DL (Transformer-based)	Binary	Supervised		TuringBench-TT
RoBERTa-Multi [108]	DL (Transformer-based)	Multi-class	Supervised		TuringBench-AA
BERT-Multi [108]	DL (Transformer-based)	Multi-class	Supervised		TuringBench-AA
TDA-based detector [67]	Hybrid	Binary	Supervised	✓	GPT-2
FAST [123]	Hybrid	Binary	Supervised		GPT-2 & GROVER
Energy discriminator [5]	DL (Energy-based)	Binary	Supervised		GROVER
MAUVE [89]	Statistical	Binary	Unsupervised	✓	GPT-2 & GROVER
Distribution detector [38]	Statistical	Binary	Unsupervised	✓	GPT-2
Feature-based detector [36]	Stylometric	Binary	Supervised	✓	GPT-2, GPT-3, & GROVER
Linguistic model [107]	Stylometric	Multi-class	Supervised	✓	Authorship Attribution-AA
DIDAN [104]	Hybrid	Binary	Supervised		Fake images w/ Human vs. GROVER news
XLNet-FT [81]	DL (Transformer-based)	Multi-class	Supervised		GPT-1, GPT-2, XLNet, & BART Reddit posts
Contra-DeBERTa [4]	DL (Transformer-based)	Multi-class	Supervised		TuringBench
Fingerprint detector [27]	Hybrid	Multi-class	Supervised		GPT-2 bot subreddit posts
DistilBERT-Academia [75]	DL (Transformer-based)	Binary	Supervised		GPT-2 Academia abstract & paper
Sentiment modeling detector [3]	DL (Glove-based)	Binary	Supervised		GPT-2 Amazon Reviews
BERT-Defense [52]	DL (Transformer-based)	Binary	Supervised		GPT-2 Large WebText
RoBERTa-Defense [91]	DL (Transformer-based)	Binary	Supervised		GROVER
RoBERTa w/ GCN [54]	DL (Transformer-based)	Binary	Supervised		GPT-2
DeBERTa v3 [95]	DL (Transformer-based)	Multi-class	Supervised		GPT-2, BLOOM, PEGASUS, OPUS, SCigen, Spinbot
Ensemble [39]	DL (Transformer-based)	Binary	Supervised		TweepFake
CoCo [73]	Hybrid	Binary	Supervised	✓	GPT-2 & GROVER
DetectGPT [80]	Statistical	Binary	Unsupervised	✓	GPT-2, OPT-2.7, GPT-Neo-2.7, GPT-J, & GPT-NeoX

Table 3: Authorship Attribution models (Binary & Multi-class) for NTD

overlap of the top-k words ($k = 100, 1K, 10K$); (3) *Lack of coherence* which can be captured using entity-grid representation to track the appearance of the grammatical role of entities. They also use neuralcoref to detect coreference entity clusters; (4) *Lack of purpose* which is captured using a lexicon-package, empath [34] containing 200 linguistic features [36]. To evaluate the generalizability of these features, an ensemble of all the features is used to build a classifier - *Feature-based detector*. This detector is trained and tested on different sizes of the GPT-2 models. It is further evaluated on GPT-3 and GROVER texts. The classifier performs consistently in detecting texts generated by GPT-3, GROVER, and different model sizes of GPT-2, suggesting it is generalizable to different NTG model sizes [36]. Further results suggest that some of the 4 categories of issues are prevalent in the top-k decoding strategy. Also, more quality-focused features (especially ones focused on Lexical diversity) perform better than statistical features such as the TF-IDF baseline.

Scaling up and creating a more realistic scenario, [27] collect 108 SubReddit blog posts generated by GPT-2 fine-tuned on 500K subreddit posts and comments. Every 108 labels indicate the 108 users of SubReddit (r/SubSimulatorGPT2). These 108 authors are detected using a set of features called “writeprints” features for the AA model [27]. Writeprints [1] is a stylometric feature that collects lexical, content-based, and idiosyncratic features as the baseline. The writeprints classifier underperforms, compared to the RoBERTa-based baseline models. Similarly, a stylometric classifier with 791 stylometric features based on [58] is used to detect neural texts. This classifier has 4 categories of features: *Character*, *word*, *sentence*, and *Lexical Diversity* features. The classifier is an ensemble of classical ML models such as Random Forest and SVM and the stylometric features. BERT, a non-stylometric clas-

sifier, outperforms these stylometric classifiers significantly [56].

Finally, stylometric classifiers are best used when the dataset size is small. When data size increases, these models underperform, allowing deep learning-based models to outperform significantly. Thus, we conclude that stylometric classifiers can only be considered good baselines since they underperform when the problem scales up. Another limitation of stylometry is that it fails to detect neural misinformation due to NTG’s capacity to generate consistent misinformation [96].

3.2 Deep Learning-based Attribution

Stylometric classifiers struggle to accurately assign the true authorship to human vs. neural texts. In Section 3.1, we observe that some of the stylometric classifiers were outperformed by deep learning-based models. Additionally, [96]’s findings of stylometry failing to detect neural misinformation, further calls for a different technique to solve the AA task for NTD. Therefore, researchers have adopted and advanced deep learning-based techniques for the attributing of neural vs. human text. These models can be further categorized into 3 types of deep learning-based classifiers - *Glove-based*, *Energy-based*, and *Transformer-based* Attribution.

3.2.1 Glove-based Attribution

Glove is an unsupervised learning algorithm for extracting the representation of words. It aggregates global word-word co-occurrence statistics from a piece of text [88]. Using GloVe word embeddings with RNN and LSTM-based neural networks was considered SOTA until, 2018 (birth of BERT [26]). Thus Glove-based classifiers now provide a good baseline for text classification tasks. Some of the best-performing AA classifiers in the

traditional AA communities are an ensemble of stylometric features + GloVe pre-trained models w/ a neural network architecture. Several Glove-based classifiers have been used as baselines for the NTD task. *Syntax-CNN* [120], *N-gram CNN* [99], and *N-gram LSTM-LSTM* [53] are baselines for [108]. *Embedding, RNN, Stacked-CNN, Parallel-CNN*, and *CNN-RNN*, are baselines for [107]. They demonstrated that Glove-based classifiers are unsuitable for solving the AA problem when there are $k > 2$ authors. Furthermore, the *Fingerprint detector* is compared with 4 baseline models - *Gaussian Naive Bayes, Random Forest, Multi-layer Perceptron*, and *CNN* classifiers using the Glove word embeddings of the data. They underperform significantly [27].

Lastly, *Sentiment modeling detector*, a variation of sentiment neuron used to learn a single-layer multiplicative LSTM (mLSTM) [64] is used to detect texts generated by GPT-2 [3]. The goal is to force the model to focus on a specified sentiment [3]. This model outperforms and sometimes performs comparably with the baseline - original mLSTM model [65], achieving a 70% accuracy in detecting GPT-2 generated Amazon and Yelp reviews [3].

3.2.2 Energy-based Attribution

Energy-based models (EBMs) are un-normalized generative models based on energy functions [70]. Using the energy functions, EBMs model the probability distribution of its training data and generates high quality data similar to the training set [70]. It is also able to adapt to changes in the Language model. Due to this capability, *Energy-based classifier* is proposed by [5] to detect neural texts. This classifier is trained on 3 datasets of different domains - *Books, CCNews*, and *Wikitext*. Three sizes of the GPT-2 model are used for the generator architectures and three architectures are used for the energy function - *Linear, BiLSTM*, and *Transformer*. Their findings suggest that: (1) as the NTG increases in size, the harder the AA task becomes; (2) the biggest energy function (i.e. *Transformer*) performs the best in detecting texts generated from large language models (e.g., GPT-2 Large & XL); (3) as the length of texts increases, the task becomes even more non-trivial; and (4) the classifiers are able to generalize to data that it is not trained on.

In addition, EBMs are very expensive to train and do not scale well [5]. While the *Energy-based classifier* performs well in the AA problem, achieving over a 90% in all experiments, applying the classifier to a much larger dataset is too expensive to justify.

3.2.3 Transformer-based Attribution

Since the advent of the Transformer architecture, the current SOTA text classification models are Transformer-based models. Based on Section 3.2.2, we observe that large models are better at detecting neural texts. However, since EBMs are too expensive, several researchers have adopted Transformer-based models (i.e., BERT, RoBERTa, etc.) for the AA tasks. In Section 3.1, most of the stylometric classifiers were outperformed by Transformer-based classifiers. This further supports the application of this classification technique to the AA problem.

*GROVER detector*⁷ [117] is trained on texts generated by the GROVER NTG. It is built with similar architecture as the GPT-2 classifier. *GROVER detector* has been evaluated on neural texts generated by different NTGs (GPT-2, FAIR, PPLM, etc.) [37, 107, 108, 123]. It performs well at detecting neural texts generated by older NTGs

⁷<https://grover.allenai.org/detect>

(2018-2019), however, struggles at detecting more recent NTGs accurately. For instance, *GROVER detector* achieved a 58% F1 score in detecting GPT-3 texts with the TuringBench dataset [108]. Next, GPT-2 has a detector trained to detect texts generated by GPT-2 - *GPT-2 detector*⁸ [51]. Just like *GROVER detector*, *GPT-2 detector* has also been evaluated on neural texts generated by different NTGs [40, 107, 108, 114] and more easily detects older NTGs than newer NTGs. This is confirmed with *GPT-2 detector*'s performance in detecting GPT-3 texts, achieving a 53% F1 score [108]. The reason is that newer NTGs, such as GPT-3 tend to generate more human-like texts which confuse SOTA older AA models, like *GPT-2 detector*.

There are two RoBERTa-based models (base & large) trained on GPT-2 dataset⁹ in huggingface repo^{10, 11}. We call both the base and large models, *OpenAI detector*. This AA model has been evaluated on neural texts generated by different NTGs [108, 114]. *OpenAI detector* is the same model as *GPT-2 detector*, except that *OpenAI detector* has been re-purposed for the AA multi-class setting, while *GPT-2 detector* remains for the AA binary setting. *OpenAI detector* performs comparably to the AA models - *BERT-Multi* and *RoBERTa-Multi* when evaluated on the TuringBench-AA dataset [108]. *BERT-Multi* and *RoBERTa-Multi* are BERT and RoBERTa base models, respectively trained on the TuringBench-AA dataset. *BERT-TT* and *RoBERTa-TT* outperform *GROVER detector* and *GPT-2 detector* when evaluated on the TuringBench-TT dataset [108]. *BERT-TT* and *RoBERTa-TT* are BERT and RoBERTa base models, respectively trained on the TuringBench-TT dataset. *BERT-TT*, outperforms all the models, including *RoBERTa-TT* significantly. Furthermore, for all 19 pairs of human vs. NTG, no model consistently outperforms all other models. In fact, *GROVER detector* and *GPT-2 detector* performs poorly in detecting texts generated by GROVER and GPT-2, respectively [108].

XLNet-FT is a fine-tuned XLNet classification model trained to detect texts generated by GPT-2 [81]. The generalizability of the model is evaluated on different subreddit post domains. *XLNet-FT* performs consistently, achieving over a 90% accuracy in all experiments, suggesting that it is generalizable [81]. However, when *XLNet-FT* is further evaluated on neural texts generated by top-p and top-k decoding strategy, there is a significant drop in accuracy, suggesting that the AA model may not be generalizable.

Using an *in-the-wild* dataset concept, *RoBERTa-Defense* is evaluated on 4 types of *in-the-wild* datasets [91]. *In-the-wild* datasets are test sets generated from an entirely different NTG from the training set. *RoBERTa-Defense* is trained on human & GROVER Real news in Table 2 and compared to other SOTA AA models - *GLTR* (with two different LMs - BERT & GPT-2 which results in *GLTR-BERT* & *GLTR-GPT2*), *GROVER detector*, *BERT-Defense*, and *FAST*. *RoBERTa-Defense* outperforms all other models significantly. *BERT-Defense*, a BERT model fine-tuned on GPT-2 Large Webtext dataset in Table 2 from which *RoBERTa-Defense* is inspired is evaluated with different decoding strategies [52]. *BERT-Defense* is trained and tested on neural texts generated by different decoding strategies - top-k, top-p, untruncated random, and mixed (i.e., dataset containing equal amounts of each strategy) [52]. The classifier trained on the mixed dataset is the most

⁸<https://huggingface.co/openai-detector/>

⁹<https://github.com/openai/gpt-2-output-dataset>

¹⁰<https://huggingface.co/roberta-base-openai-detector>

¹¹<https://huggingface.co/roberta-large-openai-detector>

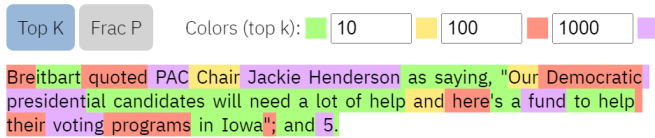


Figure 6: [108] used GLTR [45] on GPT-3 texts. Green represents the most probable words (top-10); yellow the 2nd most probable (next top-100 probable words); Red the least probable (next top-1000 probable words); and purple the highest improbable words. The hypothesis is that neural texts are often populated with mostly Green and yellow words. However, we see that texts generated by GPT-3 are very human-like according to the hypothesis.

generalizable classifier. Similarly, *RoBERTa*, *BERT*, *ELECTRA* [19], and *ALBERT* [69] are evaluated on *in-the-wild* dataset [86]. These models are evaluated, specifically on out-of-domain COVID-19 human-written vs. neural news. The neural news is generated with GPT-2 small, medium, large, XL, and GPT-Neo using top-p and top-k decoding strategies [86]. *ELECTRA* performs better at generalizing to out-of-domain neural texts, achieving an average accuracy of 86% on all out-of-domain datasets.

Due to the nuances of neural texts, [4] proposes to combine the advantages of contrastive learning [43] with a Transformer-based classifier. Thus, they propose *Constra-DeBERTa* which is a DeBERTa model [49] trained with a contrastive learning approach. Contrastive learning is a technique that clusters similar examples together and separates dissimilar examples in a representation space [4, 43]. However, while *Constra-DeBERTa* outperforms other SOTA traditional AA models, it only performs comparably to *RoBERTa-Multi* on detecting the TuringBench dataset. Similarly, [95] trains *DeBERTa v3* [48] on the SynSciPass dataset to answer the question, if a piece of text is neurally generated, how was it generated? The answer choices are *generated*, *paraphrased*, or *translated* vs. human-written. Using this dataset, *DeBERTa v3* achieves a 99.6% F1 score. Next, *DistilBERT-Academia* is trained on human vs. GPT-2 academic abstracts and papers [75] and achieves a 62.5% and 70.2% accuracy on the FULL and PARTIAL Academic datasets, respectively. Furthermore, *RoBERTa w/ GCN* (Graph Convolutional Networks) is used to detect human-written news with entities manipulated and replaced by GPT-2 [54]. The GCN [61] model is used to capture factual knowledge of neural vs. human news articles. *RoBERTa* outperformed the proposed model - *RoBERTa w/ GCN* on most of the GPT-2 detection tasks [54].

Lastly, *ruBERT*, a Russian BERT model is used to distinguish Russian neural texts from Russian human-written texts as a shared task [97]. This fine-tuned *ruBERT* (Russian BERT) achieves 82.6% accuracy for $k = 2$ authors and 64.5% accuracy for $k > 2$ authors [85]. Finally, using an *Ensemble* classifier (BART, BERTweet, and TwitterRoberta), [39] achieves an 84% accuracy in distinguishing between GPT-2 and human tweets. However, using the same model for GPT-3 generated tweets achieves a 54% accuracy [39]. This suggests that GPT-3 generates more human-like tweets than GPT-2.

3.3 Statistical Attribution

We observe that while there are some well-performing stylometric and deep learning-based models, there is still a lot of room for

improvement, especially in building generalizable models. The biggest feat is in building classifiers that perform consistently well in detecting neural texts generated by top-p and top-k decoding strategies. Thus, statistical models are proposed to combat these limitations. To assess the validity of statistical techniques, a hypothesis using k -order Markov approximations are formulated [110]. This statistical formulation proves the hypothesis that human language is stationary and ergodic as opposed to neural language. The formal hypothesis testing framework is used to establish limits in error exponents between human and neural text [110]. This suggests that statistical AA models for neural texts could be successful. There are currently only four statistical classifiers that capture the writing style of neural texts by modeling their statistical distribution.

The first statistical AA classifier proposed for NTD is *GLTR* [45]. *GLTR* performs 3 tests - (1) probability of the word; (2) the absolute rank of the word; (3) the entropy of the predicted distribution to detect neural texts. *GLTR* has a demo¹² that highlights words by distribution and is used to assist humans in detecting neural texts. See Figure 6 [108] to see how *GLTR* detects texts generated by GPT-3. This classifier improved human performance in detecting neural texts from 54% to 72%. However, since 2019 when it was built, more sophisticated NTGs have been built. These newer NTGs have more human-like statistical distribution, making it harder for *GLTR* to distinguish neural texts from human texts. *GLTR*, especially underperforms in detecting GPT-3 texts, achieving a 35% F1 score, which is significantly less than a random guess (50%) [108].

MAUVE is another statistical classifier [89]. This AA classifier measures the gap between the distribution of human and neural texts. Using KL-divergence, *MAUVE* models two types of errors that highlight the unique distributions in human vs. neural texts [89]. Human detection of texts generated by GROVER and GPT-2 correlated strongly with *MAUVE*'s highlight of differences between human and neural texts. *Distribution detector*, an unsupervised AA model for calculating the distribution of repeated n-grams in neural texts is used to detect neural texts [38]. The hypothesis is that NTGs are more repetitive than humans which is also one of the hypotheses of [36]. *Distribution detector* achieves over 90% and 80% accuracy in detecting texts generated by GPT-2 using top-k and top-p decoding strategies, respectively.

Most recently, a zero-shot unsupervised neural text detector, *DetectGPT* [80], is proposed. The hypothesis of this statistics-based detector is that neural texts tend to lie in areas of negative curvature of the log probability function [80]. Therefore, if a piece of neural text is perturbed, the curvature of the log probability will still bear a strong similarity to the unperturbed neural texts. Hence, [80] considers an AO technique that slightly modifies the original neural text, while preserving semantics. After running several perturbation experiments, a threshold for perturbation discrepancy is defined and used to detect neural text. Thus, by measuring the curvature of log probability with the strict constraint of perturbation discrepancy threshold, *DetectGPT* can detect texts generated by a neural method. Finally, *DetectGPT* detects GPT-3 generated texts with an average of 85% AU-ROC, performing comparably to *RoBERTa* [73]. Lastly, *DetectGPT* has an online demo¹³.

¹²<http://gltr.io/dist/index.html>

¹³<https://detectgpt.ericmitchell.ai/>

3.4 Hybrid Attribution

There are advantages in each of the AA model categories, however, each of them is still unable to accurately attribute neural vs. human texts to their authors consistently. Furthermore, the issue of different decoding strategies, also, make the AA models unable to generalize well [36, 50, 52, 91]. Therefore, a few researchers have proposed hybrid classifiers, which are ensembles of two or more of the AA categories.

TDA-based detector, an ensemble of the Transformer-based and statistical AA techniques is used to solve the NTD task. This classifier involves obtaining the attention matrices of BERT’s word representations of texts generated by GPT-2 and GROVER. Next, using these BERT word representations, 3 interpretable TDA-based features are extracted: (1) *Topological Features*: Calculating the first 2 betti numbers (i.e., topological features based on the connectivity of n-dimensional simplicial complexes) of the attention matrices; (2) *Features derived from barcodes*: Calculating characteristics of the barcode plots of the persistent homology of the attention matrix; (3) *Features based on the distance to patterns*: Calculating the distance in features in the attention graph. This feature is used to capture linguistic patterns. Finally, *TDA-based detector* is a logistic regression model, trained on an ensemble of the three TDA features [67]. Comparing this model to pre-trained and fine-tuned BERT models, it performs comparably to BERT models fine-tuned on GPT-2 small Webtext, GPT-2 XL Amazon Reviews, and GROVER News [67]. While more analysis is required to understand why the *TDA-based detector* performs well, this approach has interpretable qualities that should be explored in future work.

Fingerprint detector is another hybrid classifier for NTD. It is an ensemble of fine-tuned RoBERTa embeddings and CNN classifier [27]. *Fingerprint detector* solves the AA task by detecting 108 neural authors. The *Fingerprint detector* achieves a 70% accuracy (top-10). This shows promise in the area of detection of neural texts *in-the-wild*, where there are $k > 100$ authors. To continue the quest for generalizable classifiers, *FAST* uses a Graph Neural Network (GNN) architecture with RoBERTa to capture the factual structure of neural and human texts [123]. It detects neural texts by calculating the RoBERTa word embeddings of the texts and then extracting the graphical representation [123]. Next, it uses a GNN to capture sentence representations that consider coherence [123]. The experiments included detecting texts generated by GROVER and GPT-2. *FAST* outperforms *GROVER detector* and other baselines significantly. Surprisingly, it performs the best at detecting human-neural text pairs, achieving over 93% accuracy while unpaired texts achieve over 84% accuracy.

CoCo is a coherence-based contrastive learning model [73]. It is architecturally similar to *FAST* in that it uses a graphical neural network to represent the sentences of human-written vs. neural texts. Since human-written texts are more coherent than neural texts, they sentences share more entities [73]. The texts are represented as RoBERTa embedding weights which are concatenated with the sentence-level graphical representations of the texts. These concatenated features are input for an LSTM with attention. Lastly, *CoCo* is trained using the sum of the coss-entropy loss and contrastive loss [73] to improve model performance. Thus, it achieves an F1 score of 83% and 94% using the full dataset for GROVER, and GPT-2, respectively [73]. Furthermore, calculating the graphical metrics showed that in terms of the number of vertex and edges, human-written texts have significantly more

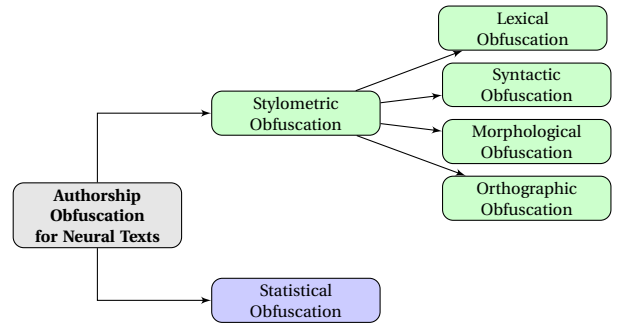


Figure 7: Taxonomy of Authorship Obfuscation algorithms/techniques for NTD

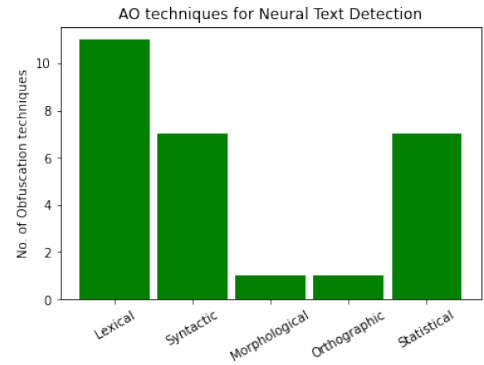


Figure 8: Number of AO techniques for NTD per category in the Taxonomy

graphical features than neural texts.

Lastly, [104] explore the most realistic scenario of misinformation where malicious users of NTGs, pair neurally generated misinformation with fake/real images to increase the authenticity of the news article. *DIDAN* is a multi-modal NTD evaluated on a multi-modal dataset containing both texts and images [104]. This NTD encodes the texts with BERT encoder and investigates Visual-Semantic representations from images and texts. These features are used to evaluate the semantic consistency between linguistic and visual components in a news article [104]. An *authenticity score* is defined to represent the probability of an article being human-written. It is calculated by extracting the co-occurrences of named entities in the news articles and captions [104]. They build different variations of dataset, some only containing text. Using only the text dataset, *DIDAN* is compared to *GROVER detector* as well as other baseline models, and outperforms all of them [104].

4. AUTHORSHIP OBFUSCATION FOR NEURAL TEXTS

In the task of NTD, AA models are evaluated under adversarial settings to assess their robustness. Due to the security risk, NTGs pose, it is important that AA models are robust to adversarial perturbations. The problem of administering adversarial perturbations to texts to cause an accurate AA model to assign inaccurate authorship is called **Authorship Obfuscation (AO)**. This is because AO is the process of masking an author’s writing style/sig-

Authorship Obfuscation	Adversarial Attack
Has a strict definition of writing style	Has a loose definition of writing style
Requires consistent/uniform change in writing style	Does not require consistent/uniform change in writing style
Requires semantics to be preserved	Does not always require semantics to be preserved

Table 4: Differences between Authorship Obfuscation and Adversarial Attack

AO technique	Scenario	Category	Interpretable	Preserves semantics	Obfuscated dataset
Homoglyph [37, 114]	Black-box	Stylometric (Orthographic)	✓	✓	GROVER & GPT-2
Upper/Lower Flip [37]	Black-box	Stylometric (Morphological)	✓	✓	GROVER & GPT-2
DeepWordBug [22, 102]	Black-box	Stylometric (Lexical)			GPT-2 & GPT-3
Misspellings attack [37, 114]	Black-box	Stylometric (Lexical)	✓		GROVER & GPT-2
Whitespace attack [37]	Black-box	Stylometric (Lexical)	✓		GROVER & GPT-2
Deduplicate tokens [90]	Black-box	Stylometric (Lexical)	✓		Human-Machine Pairs
Shuffle tokens [90]	Black-box	Stylometric (Syntactic)	✓		Human-Machine Pairs
Retain only (non)-stopwords [90]	Black-box	Stylometric (Syntactic)	✓		Human-Machine Pairs
Retain tokens in high/low frequency [90]	White-box	Statistical	✓		Human-Machine Pairs
Replace text with likelihood ranks [90]	White-box	Statistical	✓		Human-Machine Pairs
Replace text with specific linguistic features [90]	White-box	Statistical	✓		Human-Machine Pairs
TextFooler [22, 91]	Black-box	Stylometric (Lexical)		✓	GPT-2 medium, GPT-2 XL, GPT-3, GROVER
Varying sentiment [9]	Black-box	Stylometric (Lexical)	✓		GROVER
Source-target exchange [9]	Black-box	Stylometric (Syntactic)	✓		GROVER
Entity replacement [9]	Black-box	Stylometric (Lexical)	✓		GROVER
Alter numerical facts [9]	Black-box	Stylometric (Syntactic)	✓		GROVER
Syntactic perturbations [9]	Black-box	Stylometric (Syntactic)	✓		GROVER
Article shuffling [9]	Black-box	Stylometric (Syntactic)	✓		GROVER
Selecting highest human-class probability [85]	Black-box	Statistical			Russian neural & human-written texts
Adding detector’s log-probability to sampling technique [85]	White-box	Statistical			Russian neural & human-written texts
Varying the text decoding strategy (and its parameters) [91]	White-box	Statistical		✓	GPT-2 Large, GPT-2 XL, GPT-3, GROVER
Varying the number of priming tokens [91]	White-box	Statistical			GPT-2 Large, GPT-2 XL, GPT-3, GROVER
DFTFooler [91]	Black-box	Stylometric (Lexical)		✓	GPT-2 Large, GPT-2 XL, GPT-3, GROVER
Random Perturbations [91]	Black-box	Stylometric (Lexical)			GPT-2 Large, GPT-2 XL, GPT-3, GROVER
ALISON [115]	Black-box	Stylometric (Syntactic)	✓	✓	TuringBench
Mutant-X [115]	Black-box	Stylometric (lexical)		✓	TuringBench
Avengers [115]	Black-box	Stylometric (lexical)		✓	TuringBench

Table 5: Authorship Obfuscation techniques for Neural Texts. *With cited papers that implemented them.*

nature to conceal the identity, usually for privacy reasons [76]. It is a strict case of Adversarial attacks [122], as the goal is to obfuscate writing style and preserve semantics, such that both human and automatic detection is evaded. See the differences between the two in Table 4. AO is a well-studied problem in the traditional AA community [57, 76, 77, 78, 118], and has been extended to the niche AA community for NTD. This AO for neural texts problem is formulated as:

DEFINITION OF AO FOR NEURAL TEXTS. *Given an AA model $F(x)$ that accurately assigns authorship of text t to k which can be either a human or an NTG author, the AO model $O(x)$ slightly modifies t to t^* (i.e., $t^* \leftarrow O(t)$) such that the authorship is disguised (i.e., $F(t^*) \neq k$) and the difference between t^* and t is negligible.*

Thus, we survey all AO techniques employed to obfuscate neural texts in different categories, as illustrated in Figure 7: *Stylometric*

Obfuscation, and *Statistical Obfuscation*.

4.1 Stylometric Obfuscation

In order to build a robust stylometric classifier, as is observed in Section 3.1, an ensemble of features that capture several linguistic structures such as - Lexical, Syntax, etc. are required. However, to obfuscate an author’s writing style, only one of the linguistic structures may be perturbed. Therefore, all the stylometric AO techniques only target a specific linguistic structure, unlike AA classifiers. Based on the stylometric obfuscation techniques used to obfuscate neural texts, we further divide this category into 4 categories - *Lexical*, *Syntactic*, *Morphological*, and *Orthographic* Obfuscation.

4.1.1 Lexical Obfuscation

Lexical relates to the word choice of a piece of text. Thus lexical obfuscation algorithms aim to mask authors’ writing styles by replacing certain keywords with their synonyms while preserving semantics. Below, we discuss different techniques used to

achieve lexical obfuscation for neural texts. Misspellings attacks may be considered a trivial AO technique, however, it is effective in obfuscation. The misspellings attack uses a list of commonly misspelled words¹⁴ to determine which words to replace with their misspelled version. This AO technique is successful in obfuscating texts generated by GPT-2 and GROVER, and thus, evades detection of the following AA models - *GLTR*, *GROVER*, and *GPT-2* detector [114]. *GROVER detector* is further evaluated with this AO technique by obfuscating texts generated by GROVER. With only less than 10% of the texts perturbed, this attack is 94% successful [37]. However, this attack can be maneuvered by spell check algorithms, making the obfuscation technique, not robust [114]. In addition to misspelling, a whitespace attack ("will face" → "willface") is used to evaluate the robustness of *GROVER detector*. With only less than 4% of the texts perturbed, the attack is 85% successful [37].

Interesting artifacts/characteristics of neural texts still remain somewhat elusive. Therefore, perturbing these neural texts could reveal characteristics that have evaded the AA & AO community. Hence, using linguistic and statistical perturbations of words in the text, [90] extract important characteristics of neural text. For the linguistic-based perturbations, a lexical obfuscation technique is implemented - *Deduplicate tokens* which keeps the first occurrence of a token/word as is and replaces other occurrences with */MASK* token. This AO technique surprisingly improves the AA performance, suggesting that reducing the number of token occurrences may remove trivial features, causing the AA classifiers to focus on the more important features [90]

Next, texts generated by GROVER are obfuscated with the following techniques: (1) *varying sentiment*: changing the sentiment of words by replacing the word with another word of a different sentiment (positive → negative); and (2) *entity replacement*: replace entity with a useless entity [9]. Results suggest that both *GROVER* and *GPT-2 detectors* are vulnerable to these lexical-based perturbations.

Mutant-X [76] and Avengers [47] are used as baseline AO techniques to obfuscate the TuringBench dataset [115]. Mutant-X uses a genetic algorithm to search for suitable word replacement such that the semantics are preserved and the internal/substitute AA model misclassifies [76]. This process is notorious for its expensive computational complexity [103]. An internal model is used because, in the real world, the original AA model may not be known. Moreover, Mutant-X generates the obfuscated text and tests if it has evaded the AA model. If evasion is not successful, the process is repeated and tested for the defined max number of iterations [76]. These factors significantly increase the runtime of Mutant-X. Furthermore, the success of Mutant-X is dependent on how strong the internal AA model, which undermines the generalizability of Mutant-X. Hence, Avengers [47], an improved version of Mutant-X is proposed. Avengers is an ensemble version of Mutant-X. Unlike, Mutant-X, the internal AA model is an ensemble AA model, such that each classifier out of N classifiers focuses on different linguistic structures - syntax, semantics, etc.

DeepWordBug [41], a realistic character-level black-box attack is used to evaluate the robustness of 3 types of model - Statistical classification model [83], RoBERTa [74], and an Ensemble model (Statistical model + RoBERTa) [22]. It perturbs charac-

ters such that misclassification is maximized, while the Levenshtein edit distance is minimized [22, 41]. These models were trained with GPT-2 medium webtext and tested 3 separate test datasets - human vs. neural webtext from GPT-2 medium, GPT2 XL, and GPT-3 [22]. While deep learning-based classifiers achieve a higher performance in unperturbed/clean texts, statistical classifiers were found to be more robust to obfuscation [22]. Thus, the Ensemble model merges the advantages of high performance and adversarial robustness of the 2 models. DeepWordBug did not reasonably degrade the Ensemble model's performance, suggesting that DeepWordBug is not robust for this task. However, when DeepWordBug is used to evaluate the robustness of *GROVER detector* (Mega model) and *OpenAI detector* (base and large models) by perturbing the GPT-2 generated Yahoo answers & Yelp Polarity vs. Human datasets, it is successful [102]. In fact, DeepWordBug is found to be a very successful AO technique, reducing the accuracy of the Yahoo and Yelp datasets from 67.9% to 0.4% and 87.4% to 6.9%, respectively [102]. This suggests that the *OpenAI* and *GROVER detectors* are not robust to this kind of AO technique when evaluated on GPT-2 generated Yahoo answers & Yelp Polarity.

In addition, TextFooler [55], a realistic word-level black-box attack is used to evaluate the robustness of the Statistical model, RoBERTa, and Ensemble model (Statistical model + RoBERTa) [22]. TextFooler replaces words with synonyms based on cosine similarity within the embedding space [22, 55]. Based on the results, TextFooler is a robust AO technique, especially for Transformer-based models. Furthermore, as a substitute for human judgment, *MAUVE* is used to measure the human judgment of obfuscated texts. [22] finds that adversarial perturbation reduces *MAUVE* score which means that text quality is degraded and therefore neural texts are likely to be detected accurately by humans.

To further evaluate the robustness of the AA models - *GLTR* (*GLTR-BERT* & *GLTR-GPT2*), *GROVER detector*, *BERT-Defense*, *FAST*, and *RoBERTa-Defense*, texts generated by GPT-2 and GROVER are obfuscated with TextFooler, Random Perturbations [91], and DFTFooler [91] AO techniques. Random Perturbations is an attack method that replaces random words with synonyms while preserving the semantics. DFTFooler is similar to TextFooler but only needs a publicly available pre-trained LM to generate obfuscated texts. This makes DFTFooler not as computationally costly as TextFooler [91]. Also, DFTFooler perturbations are transferable [91]. To find a valid word substitution using DFTFooler, there are 4 steps: (1) synonym extraction; (2) POS checking; (3) Semantic Similarity checking; and (4) Choose a synonym with low confidence as measured by a LM. BERT and GPT-2 XL are used as the LM for DFTFooler. Results suggest that TextFooler is a stronger AO technique than DFTFooler, but performs comparably, achieving 23-91% Evasion Rate [91]. Evasion rate is defined as the fraction of perturbed neural text that evades detection by an AA model. A high evasion rate indicates a high attack success. Furthermore, using a bi-directional LM (BERT) as the backend for DFTFooler, and increasing the number of words perturbed, increases the evasion rate of DFTFooler. Based on the results, *FAST* is the most adversarially robust AA model. This may be due to the hybrid nature of the model as it combines the benefits of stylometric, statistical, and deep learning-based techniques as discussed in section 3. Another reason for *FAST*'s superior performance is its use of semantic features based on entities [91].

4.1.2 Syntactic Obfuscation

¹⁴https://en.wikipedia.org/wiki/Wikipedia:Lists_of_common_misspellings/For_machines

Syntax relates to the order of words in a piece of text. Thus, syntactic obfuscation techniques change the original arrangement of words in a document, in order to obfuscate the author's writing style. Below, we discuss such techniques used on neural texts. Characteristics of neural texts are extracted by perturbing the syntactic structure of the texts with the following syntactic perturbation techniques: (1) *Shuffle tokens* which randomly shuffles the word order of the texts; (2) *Retain only (non)-stopwords* which removes all words, except for stopwords [90]. The accuracy of the AA model only dropped marginally. This suggests that these AO techniques are not robust. It also implies that these syntactic features are not important for NTD.

The robustness of *GROVER detector* is further evaluated by syntactic obfuscation techniques on texts generated by GROVER. These techniques are: (1) *source-target exchange*: interchanging the source and target; (2) *alter numerical facts*: distort numerical facts; (3) *syntactic perturbations*: changing the word form by adding/removing punctuation ("There is" → "There's"); and (4) *article shuffling*: replacing $N\%$ of a real (human-written) article's sentences with N sentences of a fake (neural) article [9]. All AO techniques were successful, except *article shuffling*. Also, stylometric classifiers were found to be more robust, except when perturbed under stricter constraints, such as perturbing a large percentage of texts [9].

ALISON [115] is another syntactic AO technique. It reduces inference time by 100-200x when compared to SOTA AO techniques such as Mutant-X [76], and Avengers [47]. It has an internal classifier trained on a set of linguistic AA features, which allow *ALISON* to generate suitable phrase replacements that preserve semantics. *ALISON* is used to evaluate the robustness of 3 Transformer-based models - BERT, DistilBERT, and RoBERTa by obfuscating 2 datasets - TuringBench and Blog Authorship Corpus [115]. It is able to obfuscate the datasets well, causing the AA models to underperform on obfuscated texts. Furthermore, *ALISON* is able to preserve the semantics of the original text much better than their baseline AO techniques (i.e., Mutant-X and Avengers).

4.1.3 Morphological Obfuscation

Morphology is the study of word forms. Thus, morphological obfuscation techniques change the configuration of a word (e.g. "won't" → "will not"). Upper/Lower Flip ("Leaving" → "Leaving") is a morphological AO technique that may be considered trivial. However, it is successful in obfuscating texts generated by GROVER which significantly reduces the performance of *GROVER detector* [37]. With only about 2% of the texts perturbed, it achieved a 96% success rate in evading *GROVER detector*'s detection [37].

4.1.4 Orthographic Obfuscation

Orthography is the spelling convention of a language. Thus, orthographic obfuscation techniques aim to change the original spelling convention used in a piece of text to mask an author's writing style. Below, we discuss such techniques. Homoglyph attack is an orthographic perturbation technique that changes the unicode of texts. It changes English characters to cyrillic characters. This attack is almost imperceptible to the human eye and therefore, preserves semantics. The robustness of *GPT-2 detector*, *GROVER detector*, and *GLTR* is evaluated by obfuscating texts generated by GPT-2 with the homoglyph attack. *GPT-2 detector*'s performance dropped from 97.44% to 0.26% Recall. The perturbed texts caused *GROVER detector* to grossly misclassify neural texts as human-written texts and *GLTR* to shift the distribu-

tion (i.e., color scheme) of the perturbed texts [114]. *GROVER detector* is further evaluated on obfuscated texts generated by GROVER [37]. Homoglyph attack achieves a 97% success rate in perturbing *GROVER detector* [37]. However, this AO technique can easily be rendered ineffective by using spell check algorithms [114].

4.2 Statistical Obfuscation

In order to extract statistical characteristics from neural texts, the following statistical AO techniques are implemented: (1) *Replace text with likelihood ranks*; (2) *Replace text with specific linguistic features, such as Part of Speech, Dependency Trees, Constituent Trees and Named Entities*; and (3) *Retain tokens in high/low frequency regions* which defines a frequency gap score to calculate and extract the high and low-frequency words in the text [90]. The following 3 datasets are perturbed - human-machine pairs, Writing Prompt dataset [33], and CnDARIO (Chinese Novel Dataset crawled from mixed online sources) generated with GPT-2 fine-tuned with Chinese Literature. These datasets are in 3 different domains, respectively - Online Forum, News, and Literature. Results suggest that the *high/low frequency region* perturbations is the most effective obfuscation technique [90]. This further suggests that the *high/low-frequency region* feature could be an effective feature for distinguishing neural texts from human texts.

ruBERT for NTD is evaluated on 2 AO techniques - (1) *calculating the class probabilities of each label and only selecting the texts with the highest human class probability*; (2) *adds the detector's log-probability to the beam search decoding strategy during generation* so that only more human-like texts are generated [85]. These attacks achieve 46% and 56% success rates, respectively. *RoBERTa-Defense* is evaluated by changing the sampling distribution of the neural texts in the test set. This obfuscation technique involves: (1) *varying the text decoding strategy (and its parameters)*; and (2) *varying the number of priming tokens* [91]. The quality of the obfuscated neural texts is measured by GRUEN [124], a metric used to measure the linguistic quality of AI-generated texts (neural texts). GRUEN has a high correlation with human judgments. The score ranges from 0–1, and a higher value indicates high linguistic quality. Linguistic quality is based on grammaticality, non-redundancy, discourse focus, structure, and coherence [91]. Using GRUEN, a *successful attack* is defined as an attack that degrades the performance of the AA models, with little to no linguistic quality (GRUEN score) degradation [91]. The results suggest that changing the decoding strategy is an effective AO technique. Even *GLTR*, a statistical AA model is fooled by this AO technique [91].

5. EVALUATION OF AA/AO METHODS

5.1 Machine-based Evaluation

5.1.1 Authorship Attribution

To evaluate AA models, literature often used popular classification metrics such as *Precision*, *Recall*, *Accuracy*, and *F1 score*. For instance, [114] used the recall metric to evaluate the robustness of AA models toward AO techniques. Previous works evaluate the generalizability of AA models, not only on a standard single test set [36, 50, 52, 81, 91] but also on several out-of-sample distributions [102]. For example, [102] evaluate *GROVER detector* (Mega model) and *OpenAI detector* (base and large models) on three variants of test sets, namely *in-distribution*, *out-of-distribution* and *in-the-wild* datasets. *In-distribution* datasets are test sets

sampled from the same distribution of the training set, while *Out-of-distribution* datasets are those sampled from a different distribution from the training set. This test dataset is created using PPLM [25] and GeDi [63] to control the GROVER and GPT-2 generations [102]. The *in-the-wild* datasets are test sets generated from an NTG, different from the NTG used to generate a training set [102]. To build this *in-the-wild* dataset, training sets contain texts generated from GPT-2 and GROVER pre-trained models, while test sets contain texts generated from GPT-3 and a fine-tuned GPT-2 model [102]. In general, AA models perform well on *in-distribution* datasets, but suffer on *out-of-distribution* datasets, and even more on *in-the-wild* datasets.

5.1.2 Authorship Obfuscation

To evaluate AO models, literature often uses the *success rate* [37, 85], a fraction of successfully obfuscated (i.e., misclassified) texts which were accurately attributed prior to obfuscation. Another measure for AO models is *evasion rate* [91] which is the fraction of perturbed neural texts that evade the detection by an AA model. Next, due to the time and financial cost required to carry out a human-based evaluation, *MAUVE*, a metric that statistically emulates human judgments in terms of the linguistic quality (i.e., coherency) of neural texts vs. human-written texts, has been used on the AO problem [22]. That is, *MAUVE* is used as a substitute for human evaluation of obfuscated text vs. non-obfuscated text [22]. Furthermore, based on the strict definition of AO, it is sometimes important that the obfuscated text preserves the semantics of the original text. Hence, literature has measured the degradation of semantics between original and obfuscated texts, namely METEOR [6], Universal Sentence Encoder (USE) [15] Cosine similarity, and GRUEN [124]. These metrics all correlate with human judgments and a high score indicates an obfuscated text with well-preserved semantics.

5.2 Human-based Evaluation

There is currently no known human-based evaluation of AO models. The closest one is the machine-based simulation by MAUVE [22]. As such, in this section, we focus our review on the human-based evaluation of AA models, especially in the context of NTD.

5.2.1 Human Evaluation without Training

A set of research works recruited human participants (often from crowdsourcing platforms such as Amazon Mechanical Turk (AMT)) and tested whether they can distinguish neural texts from human-written texts. A simple introduction to the tasks is given, but no special training is done for human participants in this line of research. For instance, [117] examined the quality of texts generated by GROVER vs. human-written texts by humans and found that humans find GROVER-written news more believable than human-written news. [28] asked human participants to detect infilled texts filled by neural texts (e.g., she drank [blank] for [blank]) and found that humans had difficulty in detecting the infilled texts filled by neural texts. [108] introduced a benchmark for AA research, *TuringBench*, evaluated the performance of humans in distinguishing 19 pairs of human vs NTG (e.g., human vs. GPT-1 or human vs. FAIR) using TuringBench, and found that humans on average scored 51-54% of accuracies, only slightly better than random guessing. Unsurprisingly, [13] also found that humans were unable to accurately detect GPT-3 texts from human-written texts. [52] evaluated the quality of top-p and top-k decoding strategies, and found that (1) AA models detect neural texts generated by top-k decoding better than humans, but (2)

humans detect neural texts generated by top-p decoding better than AA models. Lastly, [3] found that both humans and AA models struggled to detect neural fake reviews.

5.2.2 Human Evaluation with Training

Another line of research attempted to first train human participants about NTD tasks and measure the performance improvements afterward. For instance, when human participants were trained to use GLTR [45] in detecting neural texts, thanks to the color scheme of GLTR (see Figure 6), human performance increased from 54% to 72% in accuracy. To further evaluate neural texts, [31] proposed a framework to collect a large number of human annotations via a game, Rofit¹⁵, on the quality of neural vs. human texts. Human participants were told to detect the boundary at which an article transitions from human-written to neurally-generated. Only 16% of human participants were able to correctly identify the accurate boundaries [31]. [18] studied three training strategies—instruction-based, example-based, and comparison-based, and found that example-based training is the most effective to improve human performance for solving NTD tasks (achieving the average accuracy of 55%) across the domains of story, recipe, and news. Next, [104] investigated how accurately humans can classify real vs. generated articles with and without images, using different types of news datasets: real captions and articles, real captions and generated articles, generated captions and real articles and generated captions and articles. By conducting an AMT-based study, untrained and trained human participants were able to achieve an average of 46% and 68%, respectively. Further investigation on the trustworthiness of the different article types based on style, content, consistency, and overall trustworthiness reveals that humans were skeptical about the overall trustworthiness of news articles across all four types [104]. Finally, recently, [29] proposed a framework for scrutinizing neural texts through crowdsourced data annotation in a scalable fashion, where neural texts were shown to have various error types: language-errors (i.e., lack of coherency and consistency in text), factual-errors, and reader-issues (i.e., text is too obscure or filled with too much jargon so that understanding is negatively impacted).

In conclusion, literature has found that humans alone cannot detect neural texts accurately, achieving detection accuracies only slightly better than random guessing. When humans are properly trained about the characteristics of neural texts, further, this detection accuracy tends to increase but only by small margins.

6. OPEN CHALLENGES

Although there have been several meaningful works on the current landscape of AA and AO models, the two research problems are still in their early development, especially the direction of NTD. In this section, as such, we discuss some of the remaining challenges.

6.1 Need for Comprehensive Benchmark

Generally, existing literature tends to create or use particular datasets in silos, making their findings limited and incomparable across the literature. As mentioned in [36, 81], however, the study of NTD can be greatly improved with the availability of more comprehensive and generalizable datasets whose coverage varies across diverse: (1) domains (e.g., news, online forum, recipe, stories), (2) language models, (3) decoding strategies, or (4) length of texts.

¹⁵<http://rofit.io/>

Further, not all AA/AO models share their codebase and experimental configurations, making the comparative analysis difficult. However, generating and maintaining a large number of neural texts across different settings cost significant resources and effort. [108] attempted to propose a benchmark for AA research, *TuringBench*, but it does not satisfy the needs fully. Therefore, it is greatly needed to develop a comprehensive benchmark with diverse datasets of AA/AO problems, along with the codebase of known methods in a unified environment, so that objective comparison can be carefully performed to understand the pros and cons of existing solutions and brainstorm new ideas for improvement.

6.2 Call for Complex AA/AO Variations

With the introduction of “machine” in writing high-quality texts, the set-up of “authors” in future scenarios can be more complex. For instance, one could generate a more realistic text using multiple language models in sequence (e.g., each language model improves upon the text generated by another language model in a previous step) or in parallel (e.g., each language model generates only parts of a long text). Symmetrically, it is also plausible to use multiple AO solutions in sequence or parallel to improve the overall performance of obfuscation. Yet another possible scenario is to think of “human-in-the-loop” attribution or obfuscation. For instance, would a team (of humans, of machines, or of humans and machines) outperform an individual (of human or machine) in solving AA or AO task? To our best knowledge, there is no study of AA/AO for such complex scenarios.

6.3 Need for Interpretable AA/AO Models

Currently, there are only a few interpretable AA models (e.g., GLTR) and AO techniques (e.g., Homoglyph) for neural texts, as summarized in Tables 3 and 4, respectively. That is, when an AA model detects a text as machine-generated or human-written, or when an AO model modifies parts of a text to hide authorship, it often cannot explain “why?” Ideally, however, such models should be able to provide an easy-to-understand and intuitive explanation, especially to users with no linguistic expertise, as to why a given text is attributable to a particular NTG or why a particular phrase of a text is critical to reveal an author’s identity. In addition, more research is needed to develop an intuitive human interface or visualization toward explainable AA/AO models.

6.4 Need for Improved Human Training

In parallel to improving the performance of AA/AO solutions, it is equally important to raise the awareness of AA/AO problems in the presence of neural texts, and to be able to train human users to detect neural texts better (e.g., identify phishing or misinformation message that includes neural texts as parts) or use AO solutions to hide one’s authorship (e.g., an activist posting his/her message on social media without revealing true identity). As we illustrate in Section 5.2, however, humans are not good at detecting neural texts, and there are not many AA/AO solutions suitable for novice users to benefit from in solving AA/AO tasks. Worst, still, is that even a few AA models such as *GLTR* that were shown to be able to help human users to detect neural texts better have become less effective with the advancement of NTGs. Therefore, great needs exist to have a better way to train human users in solving AA/AO tasks.

6.5 Call for Robust AA/AO Solutions

In section 4, we surveyed all literature that evaluated the robustness of AA/AO models, and found that most existing AO techniques do not preserve the original semantics of text well and thus cannot easily evade the attribution of AA solutions, especially human detection. Similarly, as we adopt more sophisticated hybrid approaches for AA tasks, successful AO attacks to hide authorship will become more challenging. Part of the reason for these vulnerabilities in existing AA/AO solutions is that the bulk of existing literature has studied either AA or AO problem in separation, thus greatly limiting their robustness against the other problem. Therefore, to stay relevant and synonymous with a real-life scenario, both AA and AO solutions need to learn from each other, and co-train/co-evolve, as in a min-max optimization game.

7. APPLICATIONS

Deepfake Detection: Successful solutions for AA/AO tasks can be useful in many applications. For instance, recently, the generation of realistic AI-made images¹⁶ and videos, so-called “deepfakes”, have flooded the Web. While most of these deepfakes are made for humor, some are malicious in generating misinformation, spreading political propaganda, or attacking individuals [79]. In literature, in particular, [104] studies the realistic scenario where real images would be paired with neural texts to increase the authenticity of a news article as well as evade detection. In such a setting, successful AA solutions can point out the non-human nature of neural texts to users or can be used to extract features of neural texts for downstream deepfake detection models.

Chatbot Detection: Another application is for AA solutions to detect suspicious messages (e.g., phishing or chatbot messages) that may have been (partially) generated by NTG. Similar to neural texts of news format, shorter or informal chatbot messages are also hard to discriminate when generated by machines [98, 106]. An example of the state-of-the-art chatbot is ChatGPT [84] which has been used to generate medical writings [10, 40], finance writings [30], etc. These applications of ChatGPT have also increased the likelihood of cheating in academic writing [21]. Thus, AA models for neural text detection will be beneficial in distinguishing chatGPT-generated texts from human-written texts.

Anonymity Preservation: On the other hand, successful AO solutions can be used to help individuals who have needs to share their writings without jeopardizing their secret identity. For instance, an NGO activist or whistleblower may submit her op-ed to news media after making sure that no popular AA solutions can attribute the writing to her.

8. CONCLUSION

With the rise of neural texts that were generated by large-scale language models, we are currently in an arms race between generation and detection of *neural texts*. In this work, we comprehensively survey two important problems of neural texts: **Authorship Attribution** (AA) and **Authorship Obfuscation** (AO). We first categorize existing AA solutions into four types of stylometric, deep learning-based, statistical, and hybrid attribution. Similarly, we categorize existing AO solutions into two types of stylometric and statistical obfuscation, and elaborate pros and cons of representative methods therein. In addition, we discuss different evaluation methods for AA and AO problems in the context of neural texts, and finally, share a few important challenges that we

¹⁶<https://thispersondoesnotexist.com/>

feel lacking currently. By and large, we believe that the data mining community is well-positioned to be able to contribute to significant improvement in both AA and AO research. Despite their close implications in security and privacy, with respect to the underlying methods used, their problem formulation as supervised or unsupervised learning, and their focus on the accuracy and running time as major metrics. Lastly, to mitigate the challenges of accurate detection of neural text, [62] proposes watermarking these text-generative language models. This entails embedding humanly imperceptible signals into the language models such that they generate semantically relevant texts, unnoticeable to humans but noticeable to detectors. These watermarking [2, 62] techniques attempt to solve the security risks that these language models pose. However, as these watermarking techniques have not yet been widely adopted, we still have to rely on AA and AO solutions for neural text detection. Also, as such watermarking techniques are a recent/new development, their robustness to strong AO techniques has not yet been evaluated.

9. ACKNOWLEDGMENT

This work was in part supported by NSF awards #1820609, #2114824, and #2131144.

References

- [1] A. Abbasi and H. Chen. Writeprints: A stylometric approach to identity-level identification and similarity detection in cyberspace. *ACM Transactions on Information Systems (TOIS)*, 26(2):1–29, 2008.
- [2] S. Abdelnabi and M. Fritz. Adversarial watermarking transformer: Towards tracing text provenance with data hiding. In *2021 IEEE Symposium on Security and Privacy (SP)*, pages 121–140. IEEE, 2021.
- [3] D. I. Adelani, H. Mai, F. Fang, H. H. Nguyen, J. Yamagishi, and I. Echizen. Generating sentiment-preserving fake online reviews using neural language models and their human-and machine-based detection. In *International Conference on Advanced Information Networking and Applications*, pages 1341–1354. Springer, 2020.
- [4] B. Ai, Y. Wang, Y. Tan, and T. Samson. Whodunit? learning to contrast for authorship attribution. *Proceedings of the 2nd Conference of the Asia-Pacific Chapter of the Association for Computational Linguistics and the 12th International Joint Conference on Natural Language Processing*, 2022.
- [5] A. Bakhtin, S. Gross, M. Ott, Y. Deng, M. Ranzato, and A. Szlam. Real or fake? learning to discriminate machine from human generated text. *arXiv preprint arXiv:1906.03351*, 2019.
- [6] S. Banerjee and A. Lavie. Meteor: An automatic metric for mt evaluation with improved correlation with human judgments. In *Proceedings of the acl workshop on intrinsic and extrinsic evaluation measures for machine translation and/or summarization*, pages 65–72, 2005.
- [7] P. Barham, A. Chowdhery, J. Dean, S. Ghemawat, S. Hand, D. Hurt, M. Isard, H. Lim, R. Pang, S. Roy, et al. Pathways: Asynchronous distributed dataflow for ml. *Proceedings of Machine Learning and Systems*, 4:430–449, 2022.
- [8] J. Bevendorff, B. Chulvi, E. Fersini, A. Heini, M. Kestemont, K. Kredens, M. Mayerl, R. Ortega-Bueno, P. Pezik, M. Potthast, et al. Overview of pan 2022: Authorship verification, profiling irony and stereotype spreaders, style change detection, and trigger detection. In *European Conference on Information Retrieval*, pages 331–338. Springer, 2022.
- [9] M. M. Bhat and S. Parthasarathy. How effectively can machines defend against machine-generated fake news? an empirical study. In *Proceedings of the First Workshop on Insights from Negative Results in NLP*, pages 48–53, 2020.
- [10] S. Biswas. Chatgpt and the future of medical writing, 2023.
- [11] A. F. Biten, L. Gomez, M. Rusinol, and D. Karatzas. Good news, everyone! context driven entity-aware captioning for news images. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 12466–12475, 2019.
- [12] S. Black, S. Biderman, E. Hallahan, Q. Anthony, L. Gao, L. Golding, H. He, C. Leahy, K. McDonnell, J. Phang, et al. Gpt-neox-20b: An open-source autoregressive language model. In *Proceedings of BigScience Episode\#5-Workshop on Challenges & Perspectives in Creating Large Language Models*, pages 95–136, 2022.
- [13] T. Brown, B. Mann, N. Ryder, M. Subbiah, J. D. Kaplan, P. Dhariwal, A. Neelakantan, P. Shyam, G. Sastry, A. Askell, S. Agarwal, A. Herbert-Voss, G. Krueger, T. Henighan, R. Child, A. Ramesh, D. Ziegler, J. Wu, C. Winter, C. Hesse, M. Chen, E. Sigler, M. Litwin, S. Gray, B. Chess, J. Clark, C. Berner, S. McCandlish, A. Radford, I. Sutskever, and D. Amodei. Language models are few-shot learners. In H. Larochelle, M. Ranzato, R. Hadsell, M. F. Balcan, and H. Lin, editors, *Advances in Neural Information Processing Systems*, volume 33, pages 1877–1901. Curran Associates, Inc., 2020.
- [14] B. Buchanan, A. Lohn, M. Musser, and K. Sedova. Truth, lies, and automation. *Center for Security and Emerging Technology*, 2021.
- [15] D. Cer, Y. Yang, S.-y. Kong, N. Hua, N. Limtiaco, R. S. John, N. Constant, M. Guajardo-Cespedes, S. Yuan, C. Tar, et al. Universal sentence encoder for english. In *Proceedings of the 2018 conference on empirical methods in natural language processing: system demonstrations*, pages 169–174, 2018.
- [16] P.-J. Chen, A. Lee, C. Wang, N. Goyal, A. Fan, M. Williamson, and J. Gu. Facebook ai’s wmt20 news translation task submission. In *Proceedings of the Fifth Conference on Machine Translation*, pages 113–125, 2020.
- [17] A. Chowdhery, S. Narang, J. Devlin, M. Bosma, G. Mishra, A. Roberts, P. Barham, H. W. Chung, C. Sutton, S. Gehrmann, et al. Palm: Scaling language modeling with pathways. *arXiv preprint arXiv:2204.02311*, 2022.
- [18] E. Clark, T. August, S. Serrano, N. Haduong, S. Gururangan, and N. A. Smith. All that’s ‘human’ is not gold: Evaluating human evaluation of generated text. In *Proceedings of the 59th Annual Meeting of the Association for Computational Linguistics and the 11th International Joint Conference on Natural Language Processing (Volume 1: Long Papers)*, pages 7282–7296, 2021.

- [19] K. Clark, M.-T. Luong, Q. V. Le, and C. D. Manning. Electra: Pre-training text encoders as discriminators rather than generators. *arXiv preprint arXiv:2003.10555*, 2020.
- [20] A. Conneau and G. Lample. Cross-lingual language model pretraining. *Advances in neural information processing systems*, 32, 2019.
- [21] D. R. Cotton, P. A. Cotton, and J. R. Shipway. Chatting and cheating. ensuring academic integrity in the era of chatgpt. 2023.
- [22] E. Crothers, N. Japkowicz, H. Viktor, and P. Branco. Adversarial robustness of neural-statistical features in detection of generative transformers. *arXiv preprint arXiv:2203.07983*, 2022.
- [23] J. Cutler, L. Dugan, S. Havaldar, and A. Stein. Automatic detection of hybrid human-machine text boundaries. 2021.
- [24] Z. Dai, Z. Yang, Y. Yang, J. G. Carbonell, Q. Le, and R. Salakhutdinov. Transformer-xl: Attentive language models beyond a fixed-length context. In *Proceedings of the 57th Annual Meeting of the Association for Computational Linguistics*, pages 2978–2988, 2019.
- [25] S. Dathathri, A. Madotto, J. Lan, J. Hung, E. Frank, P. Molino, J. Yosinski, and R. Liu. Plug and play language models: A simple approach to controlled text generation. In *International Conference on Learning Representations*, 2020.
- [26] J. Devlin, M.-W. Chang, K. Lee, and K. Toutanova. Bert: Pre-training of deep bidirectional transformers for language understanding. *arXiv preprint arXiv:1810.04805*, 2018.
- [27] N. Diwan, T. Chakraborty, and Z. Shafiq. Fingerprinting fine-tuned language models in the wild. In *Findings of the Association for Computational Linguistics: ACL-IJCNLP 2021*, pages 4652–4664, 2021.
- [28] C. Donahue, M. Lee, and P. Liang. Enabling language models to fill in the blanks. In *Proceedings of the 58th Annual Meeting of the Association for Computational Linguistics*, pages 2492–2501, 2020.
- [29] Y. Dou, M. Forbes, R. Koncel-Kedziorski, N. A. Smith, and Y. Choi. Scarecrow: A framework for scrutinizing machine text. *arXiv preprint arXiv:2107.01294*, 2021.
- [30] M. Dowling and B. Lucey. Chatgpt for (finance) research: The bananarama conjecture. *Finance Research Letters*, page 103662, 2023.
- [31] L. Dugan, D. Ippolito, A. Kirubarajan, and C. Callison-Burch. Rofit: A tool for evaluating human detection of machine-generated text. In *Proceedings of the 2020 Conference on Empirical Methods in Natural Language Processing: System Demonstrations*, pages 189–196, 2020.
- [32] T. Fagni, F. Falchi, M. Gambini, A. Martella, and M. Tesconi. Tweepfake: About detecting deepfake tweets. *Plos one*, 16(5):e0251415, 2021.
- [33] A. Fan, M. Lewis, and Y. Dauphin. Hierarchical neural story generation. In *Proceedings of the 56th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 889–898, 2018.
- [34] E. Fast, B. Chen, and M. S. Bernstein. Empath: Understanding topic signals in large-scale text. In *Proceedings of the 2016 CHI conference on human factors in computing systems*, pages 4647–4657, 2016.
- [35] W. Fedus, B. Zoph, and N. Shazeer. Switch transformers: Scaling to trillion parameter models with simple and efficient sparsity. *Journal of Machine Learning Research*, 23(120):1–39, 2022.
- [36] L. Fröhling and A. Zubiaga. Feature-based detection of automated language models: tackling gpt-2, gpt-3 and grover. *PeerJ Computer Science*, 7:e443, 2021.
- [37] R. Gagiano, M. M.-H. Kim, X. J. Zhang, and J. Biggs. Robustness analysis of grover for machine-generated news detection. In *Proceedings of the The 19th Annual Workshop of the Australasian Language Technology Association*, pages 119–127, 2021.
- [38] M. Gallé, J. Rozen, G. Kruszewski, and H. Elsahar. Unsupervised and distributional detection of machine-generated text. *arXiv preprint arXiv:2111.02878*, 2021.
- [39] M. Gambini, T. Fagni, F. Falchi, and M. Tesconi. On pushing deepfake tweet detection capabilities to the limits. In *14th ACM Web Science Conference 2022*, pages 154–163, 2022.
- [40] C. A. Gao, F. M. Howard, N. S. Markov, E. C. Dyer, S. Ramesh, Y. Luo, and A. T. Pearson. Comparing scientific abstracts generated by chatgpt to original abstracts using an artificial intelligence output detector, plagiarism detector, and blinded human reviewers. *bioRxiv*, pages 2022–12, 2022.
- [41] J. Gao, J. Lanchantin, M. L. Soffa, and Y. Qi. Black-box generation of adversarial text sequences to evade deep learning classifiers. In *2018 IEEE Security and Privacy Workshops (SPW)*, pages 50–56. IEEE, 2018.
- [42] L. Gao, S. Biderman, S. Black, L. Golding, T. Hoppe, C. Foster, J. Phang, H. He, A. Thite, N. Nabeshima, et al. The pile: An 800gb dataset of diverse text for language modeling. *arXiv preprint arXiv:2101.00027*, 2020.
- [43] T. Gao, X. Yao, and D. Chen. Simcse: Simple contrastive learning of sentence embeddings. In *Proceedings of the 2021 Conference on Empirical Methods in Natural Language Processing*, pages 6894–6910, 2021.
- [44] S. Gehman, S. Gururangan, M. Sap, Y. Choi, and N. A. Smith. Realtoxicityprompts: Evaluating neural toxic degeneration in language models. In *Findings of the Association for Computational Linguistics: EMNLP 2020*, pages 3356–3369, 2020.
- [45] S. Gehrmann, H. Strobel, and A. M. Rush. Gltr: Statistical detection and visualization of generated text. In *Proceedings of the 57th Annual Meeting of the Association for Computational Linguistics: System Demonstrations*, pages 111–116, 2019.
- [46] J. Guerrero and I. Alsmadi. Synthetic text detection: Systemic literature review. *arXiv preprint arXiv:2210.06336*, 2022.

- [47] M. Haroon, F. Zaffar, P. Srinivasan, and Z. Shafiq. Avengers ensemble! improving transferability of authorship obfuscation. *arXiv preprint arXiv:2109.07028*, 2021.
- [48] P. He, J. Gao, and W. Chen. Debertav3: Improving deberta using electra-style pre-training with gradient-disentangled embedding sharing. *arXiv preprint arXiv:2111.09543*, 2021.
- [49] P. He, X. Liu, J. Gao, and W. Chen. Deberta: Decoding-enhanced bert with disentangled attention. In *International Conference on Learning Representations*, 2021.
- [50] A. Holtzman, J. Buys, L. Du, M. Forbes, and Y. Choi. The curious case of neural text degeneration. In *International Conference on Learning Representations*, 2019.
- [51] Huggingface. Gpt-2 output detector demo. <https://huggingface.co/openai-detector/>, 2019.
- [52] D. Ippolito, D. Duckworth, C. Callison-Burch, and D. Eck. Automatic detection of generated text is easiest when humans are fooled. In *Proceedings of the 58th Annual Meeting of the Association for Computational Linguistics*, pages 1808–1822, 2020.
- [53] F. Jafariakinabad, S. Tarnpradab, and K. A. Hua. Syntactic recurrent neural network for authorship attribution. *arXiv preprint arXiv:1902.09723*, 2019.
- [54] G. Jawahar, M. Abdul-Mageed, and L. Lakshmanan. Automatic detection of entity-manipulated text using factual knowledge. In *Proceedings of the 60th Annual Meeting of the Association for Computational Linguistics (Volume 2: Short Papers)*, pages 86–93, 2022.
- [55] D. Jin, Z. Jin, J. T. Zhou, and P. Szolovits. Is bert really robust? a strong baseline for natural language attack on text classification and entailment. In *Proceedings of the AAAI conference on artificial intelligence*, volume 34, pages 8018–8025, 2020.
- [56] K. Jones, J. R. Nurse, and S. Li. Are you robert or roberta? deceiving online authorship attribution models using neural text generators. In *Proceedings of the International AAAI Conference on Web and Social Media*, volume 16, pages 429–440, 2022.
- [57] G. Karadzhov, T. Mihaylova, Y. Kiprova, G. Georgiev, I. Koychev, and P. Nakov. The case for being average: A mediocrity approach to style masking and author obfuscation. In *International Conference of the Cross-Language Evaluation Forum for European Languages*, pages 173–185. Springer, 2017.
- [58] R. Kaur, S. Singh, and H. Kumar. Authorship analysis of online social media content. In *Proceedings of 2nd International Conference on Communication, Computing and Networking*, pages 539–549. Springer, 2019.
- [59] N. S. Keskar, B. McCann, L. R. Varshney, C. Xiong, and R. Socher. Ctrl: A conditional transformer language model for controllable generation. *arXiv preprint arXiv:1909.05858*, 2019.
- [60] M. Kestemont, E. Manjavacas, I. Markov, J. Bevendorff, M. Wiegmann, E. Stamatatos, B. Stein, and M. Potthast. Overview of the cross-domain authorship verification task at pan 2021. In *CLEF (Working Notes)*, 2021.
- [61] T. N. Kipf and M. Welling. Semi-supervised classification with graph convolutional networks. *arXiv preprint arXiv:1609.02907*, 2016.
- [62] J. Kirchenbauer, J. Geiping, Y. Wen, J. Katz, I. Miers, and T. Goldstein. A watermark for large language models. *arXiv preprint arXiv:2301.10226*, 2023.
- [63] B. Krause, A. D. Gotmare, B. McCann, N. S. Keskar, S. Joty, R. Socher, and N. F. Rajani. Gedi: Generative discriminator guided sequence generation. In *Findings of the Association for Computational Linguistics: EMNLP 2021*, pages 4929–4952, 2021.
- [64] B. Krause, L. Lu, I. Murray, and S. Renals. Multiplicative lstm for sequence modelling. *arXiv preprint arXiv:1609.07959*, 2016.
- [65] B. Krause, I. Murray, S. Renals, and L. Liang. Multiplicative lstm for sequence modelling. In *5th International Conference on Learning Representations*, pages 2872–2880, 2017.
- [66] S. Kreps, R. M. McCain, and M. Brundage. All the news that’s fit to fabricate: Ai-generated text as a tool of media misinformation. *Journal of Experimental Political Science*, 9(1):104–117, 2022.
- [67] L. Kushnareva, D. Cherniavskii, V. Mikhailov, E. Artemova, S. Barannikov, A. Bernstein, I. Piontkovskaya, D. Piontkovski, and E. Burnaev. Artificial text detection via examining the topology of attention maps. In *Proceedings of the 2021 Conference on Empirical Methods in Natural Language Processing*, pages 635–649, 2021.
- [68] K. Lagutina, N. Lagutina, E. Boychuk, I. Vorontsova, E. Shliakhtina, O. Belyaeva, I. Paramonov, and P. Demidov. A survey on stylometric text features. In *2019 25th Conference of Open Innovations Association (FRUCT)*, pages 184–195. IEEE, 2019.
- [69] Z. Lan, M. Chen, S. Goodman, K. Gimpel, P. Sharma, and R. Soricut. Albert: A lite bert for self-supervised learning of language representations. *arXiv preprint arXiv:1909.11942*, 2019.
- [70] Y. LeCun, S. Chopra, R. Hadsell, M. Ranzato, and F. Huang. A tutorial on energy-based learning. *Predicting structured data*, 1(0), 2006.
- [71] M. Lewis, Y. Liu, N. Goyal, M. Ghazvininejad, A. Mohamed, O. Levy, V. Stoyanov, and L. Zettlemoyer. Bart: Denoising sequence-to-sequence pre-training for natural language generation, translation, and comprehension. In *Proceedings of the 58th Annual Meeting of the Association for Computational Linguistics*, pages 7871–7880, 2020.
- [72] J. Li, T. Tang, W. X. Zhao, and J.-R. Wen. Pretrained language models for text generation: A survey. *arXiv preprint arXiv:2105.10311*, 2021.
- [73] X. Liu, Z. Zhang, Y. Wang, Y. Lan, and C. Shen. Coco: Coherence-enhanced machine-generated text detection under data limitation with contrastive learning. *arXiv preprint arXiv:2212.10341*, 2022.
- [74] Y. Liu, M. Ott, N. Goyal, J. Du, M. Joshi, D. Chen, O. Levy, M. Lewis, L. Zettlemoyer, and V. Stoyanov. Roberta: A robustly optimized bert pretraining approach. *arXiv preprint arXiv:1907.11692*, 2019.

- [75] V. Liyanage, D. Buscaldi, and A. Nazarenko. A benchmark corpus for the detection of automatically generated text in academic publications. In *LREC*, 2022.
- [76] A. Mahmood, F. Ahmad, Z. Shafiq, P. Srinivasan, and F. Zaffar. A girl has no name: Automated authorship obfuscation using mutant-x. *Proc. Priv. Enhancing Technol.*, 2019(4):54–71, 2019.
- [77] A. Mahmood, Z. Shafiq, and P. Srinivasan. A girl has a name: Detecting authorship obfuscation. In *Proceedings of the 58th Annual Meeting of the Association for Computational Linguistics*, pages 2235–2245, 2020.
- [78] A. W. McDonald, S. Afroz, A. Caliskan, A. Stolerman, and R. Greenstadt. Use fewer instances of the letter “i”: Toward writing style anonymization. In *International Symposium on Privacy Enhancing Technologies Symposium*, pages 299–318. Springer, 2012.
- [79] Y. Mirsky and W. Lee. The creation and detection of deep-fakes: A survey. *ACM Computing Surveys (CSUR)*, 54(1):1–41, 2021.
- [80] E. Mitchell, Y. Lee, A. Khazatsky, C. D. Manning, and C. Finn. Detectgpt: Zero-shot machine-generated text detection using probability curvature. *arXiv preprint arXiv:2301.11305*, 2023.
- [81] S. Munir, B. Batool, Z. Shafiq, P. Srinivasan, and F. Zaffar. Through the looking glass: Learning to attribute synthetic text generated by language models. In *Proceedings of the 16th Conference of the European Chapter of the Association for Computational Linguistics: Main Volume*, pages 1811–1822, 2021.
- [82] N. Ng, K. Yee, A. Baevski, M. Ott, M. Auli, and S. Edunov. Facebook fair’s wmt19 news translation task submission. *arXiv preprint arXiv:1907.06616*, 2019.
- [83] H.-Q. Nguyen-Son, N.-D. T. Tieu, H. H. Nguyen, J. Yamagishi, and I. E. Zen. Identifying computer-generated text using statistical analysis. In *2017 Asia-Pacific Signal and Information Processing Association Annual Summit and Conference (APSIPA ASC)*, pages 1504–1511. IEEE, 2017.
- [84] OpenAI. Optimizing language models for dialogue. <https://openai.com/blog/chatgpt/>, 2022.
- [85] M. Orzhenovskii. Detecting auto-generated texts with language model and attacking the detector. *Computational Linguistics and Intellectual Technologies: Proceedings of the International Conference “Dialogue 2022*, 2022.
- [86] A. Pagnoni, M. Graciarena, and Y. Tsvetkov. Threat scenarios and best practices to detect neural fake news. In *Proceedings of the 29th International Conference on Computational Linguistics*, pages 1233–1249, 2022.
- [87] J. W. Pennebaker, M. E. Francis, and R. J. Booth. Linguistic inquiry and word count: Liwc 2001. *Mahway: Lawrence Erlbaum Associates*, 71(2001):2001, 2001.
- [88] J. Pennington, R. Socher, and C. D. Manning. Glove: Global vectors for word representation. In *Proceedings of the 2014 conference on empirical methods in natural language processing (EMNLP)*, pages 1532–1543, 2014.
- [89] K. Pillutla, S. Swayamdipta, R. Zellers, J. Thickstun, S. Welleck, Y. Choi, and Z. Harchaoui. An information divergence measure between neural text and human text. *arXiv preprint arXiv:2102.01454*, 2021.
- [90] J. Pu, Z. Huang, Y. Xi, G. Chen, W. Chen, and R. Zhang. Unraveling the mystery of artifacts in machine generated text. pages 6889–6898, 2022.
- [91] J. Pu, Z. Sarwar, S. M. Abdullah, A. Rehman, Y. Kim, P. Bhattacharya, M. Javed, B. Viswanath, V. Tech, and L. Pakistan. Deepfake text detection: Limitations and opportunities. *44th IEEE Symposium on Security and Privacy*, 2023.
- [92] A. Radford, K. Narasimhan, T. Salimans, and I. Sutskever. Improving language understanding by generative pre-training. URL https://s3-us-west-2.amazonaws.com/openai-assets/researchcovers/languageunsupervised/language_understanding_paper.pdf, 2018.
- [93] A. Radford, J. Wu, R. Child, D. Luan, D. Amodei, and I. Sutskever. Language models are unsupervised multitask learners. *OpenAI blog*, 1(8):9, 2019.
- [94] C. Raffel, N. Shazeer, A. Roberts, K. Lee, S. Narang, M. Matena, Y. Zhou, W. Li, P. J. Liu, et al. Exploring the limits of transfer learning with a unified text-to-text transformer. *J. Mach. Learn. Res.*, 21(140):1–67, 2020.
- [95] D. Rosati. Synscipass: detecting appropriate uses of scientific text generation. In *Proceedings of the Third Workshop on Scholarly Document Processing*, pages 214–222, 2022.
- [96] T. Schuster, R. Schuster, D. J. Shah, and R. Barzilay. The limitations of stylometry for detecting machine-generated fake news. *Computational Linguistics*, 46(2):499–510, 2020.
- [97] T. Shamardina, V. Mikhailov, D. Chernianskii, A. Fenogenova, M. Saidov, A. Valeeva, T. Shavrina, I. Smurov, E. Tutubalina, and E. Artemova. Findings of the the ruatd shared task 2022 on artificial text detection in russian. *arXiv preprint arXiv:2206.01583*, 2022.
- [98] J. Shao, A. Uchendu, and D. Lee. A reverse turing test for detecting machine-made texts. In *Proceedings of the 10th ACM Conference on Web Science*, pages 275–279, 2019.
- [99] P. Shrestha, S. Sierra, F. Gonzalez, M. Montes, P. Rosso, and T. Solorio. Convolutional neural networks for authorship attribution of short texts. In *Proceedings of the 15th Conference of the European Chapter of the Association for Computational Linguistics: Volume 2, Short Papers*, 2017.
- [100] E. Stamatatos. Authorship verification: a review of recent advances. *Research in Computing Science*, 123:9–25, 2016.
- [101] E. Stamatatos, M. Kestemont, K. Kredens, P. Pezik, A. Heini, J. Bevendorff, M. Potthast, and B. Stein. Overview of the authorship verification task at pan 2022. *Working Notes of CLEF*, 2022.
- [102] H. Stiff and F. Johansson. Detecting computer-generated disinformation. *International Journal of Data Science and Analytics*, 13(4):363–383, 2022.

- [103] M. Tabatabaei, J. Hakanen, M. Hartikainen, K. Miettinen, and K. Sindhya. A survey on handling computationally expensive multiobjective optimization problems using surrogates: non-nature inspired methods. *Structural and Multidisciplinary Optimization*, 52(1):1–25, 2015.
- [104] R. Tan, B. Plummer, and K. Saenko. Detecting cross-modal inconsistency to defend against neural fake news. In *Proceedings of the 2020 Conference on Empirical Methods in Natural Language Processing (EMNLP)*, pages 2081–2106, 2020.
- [105] J. Tyo, B. Dhingra, and Z. C. Lipton. On the state of the art in authorship attribution and authorship verification. *arXiv preprint arXiv:2209.06869*, 2022.
- [106] A. Uchendu, J. Cao, Q. Wang, B. Luo, and D. Lee. Characterizing man-made vs. machine-made chatbot dialogs. In *TTO*, 2019.
- [107] A. Uchendu, T. Le, K. Shu, and D. Lee. Authorship attribution for neural text generation. In *Proceedings of the 2020 Conference on Empirical Methods in Natural Language Processing (EMNLP)*, pages 8384–8395, 2020.
- [108] A. Uchendu, Z. Ma, T. Le, R. Zhang, and D. Lee. Turing-bench: A benchmark environment for turing test in the age of neural text generation. In *Findings of the Association for Computational Linguistics: EMNLP 2021*, pages 2001–2016, 2021.
- [109] O. Varol, E. Ferrara, C. Davis, F. Menczer, and A. Flammini. Online human-bot interactions: Detection, estimation, and characterization. In *Proceedings of the international AAAI conference on web and social media*, volume 11, pages 280–289, 2017.
- [110] L. R. Varshney, N. S. Keskar, and R. Socher. Limits of detecting text generated by large-scale language models. In *2020 Information Theory and Applications Workshop (ITA)*, pages 1–5. IEEE, 2020.
- [111] B. Wang. Mesh-Transformer-JAX: Model-Parallel Implementation of Transformer Language Model with JAX. <https://github.com/kingoflolz/mesh-transformer-jax>, May 2021.
- [112] B. Wang and A. Komatsuzaki. GPT-J-6B: A 6 Billion Parameter Autoregressive Language Model. <https://github.com/kingoflolz/mesh-transformer-jax>, May 2021.
- [113] T. Wolf, L. Debut, V. Sanh, J. Chaumond, C. Delangue, A. Moi, P. Cistac, T. Rault, R. Louf, M. Funtowicz, et al. Huggingface’s transformers: State-of-the-art natural language processing. *arXiv preprint arXiv:1910.03771*, 2019.
- [114] M. Wolff and S. Wolff. Attacking neural text detectors. *arXiv preprint arXiv:2002.11768*, 2020.
- [115] E. Xing, T. Le, and D. Lee. Alison: Fast stylometric authorship obfuscation. *Technical Report*, 2023.
- [116] Z. Yang, Z. Dai, Y. Yang, J. Carbonell, R. R. Salakhutdinov, and Q. V. Le. Xlnet: Generalized autoregressive pretraining for language understanding. In *Advances in neural information processing systems*, pages 5754–5764, 2019.
- [117] R. Zellers, A. Holtzman, H. Rashkin, Y. Bisk, A. Farhadi, F. Roesner, and Y. Choi. Defending against neural fake news. In *Advances in Neural Information Processing Systems*, pages 9051–9062, 2019.
- [118] W. Zhai, J. Rusert, Z. Shafiq, and P. Srinivasan. A girl has a name, and it’s... adversarial authorship attribution for de-obfuscation. *arXiv preprint arXiv:2203.11849*, 2022.
- [119] H. Zhang, H. Song, S. Li, M. Zhou, and D. Song. A survey of controllable text generation using transformer-based pre-trained language models. *arXiv preprint arXiv:2201.05337*, 2022.
- [120] R. Zhang, Z. Hu, H. Guo, and Y. Mao. Syntax encoding with application in authorship attribution. In *Proceedings of the 2018 Conference on Empirical Methods in Natural Language Processing*, pages 2742–2753, 2018.
- [121] S. Zhang, S. Roller, N. Goyal, M. Artetxe, M. Chen, S. Chen, C. Dewan, M. Diab, X. Li, X. V. Lin, et al. Opt: Open pre-trained transformer language models. *arXiv preprint arXiv:2205.01068*, 2022.
- [122] W. E. Zhang, Q. Z. Sheng, A. Alhazmi, and C. Li. Adversarial attacks on deep-learning models in natural language processing: A survey. *ACM Transactions on Intelligent Systems and Technology (TIST)*, 11(3):1–41, 2020.
- [123] W. Zhong, D. Tang, Z. Xu, R. Wang, N. Duan, M. Zhou, J. Wang, and J. Yin. Neural deepfake detection with factual structure of text. In *Proceedings of the 2020 Conference on Empirical Methods in Natural Language Processing (EMNLP)*, pages 2461–2470, 2020.
- [124] W. Zhu and S. Bhat. Gruen for evaluating linguistic quality of generated text. In *Findings of the Association for Computational Linguistics: EMNLP 2020*, pages 94–108, 2020.

The Need for Unsupervised Outlier Model Selection: A Review and Evaluation of Internal Evaluation Strategies

Martin Q. Ma[‡], Yue Zhao[‡], Xiaorong Zhang, and Leman Akoglu
Carnegie Mellon University, Pittsburgh, PA, USA

qianlim@cs.cmu.edu, zhaoy@cmu.edu, xiaorongz@alumni.cmu.edu,
lakoglu@andrew.cmu.edu

ABSTRACT

Given an *unsupervised* outlier detection task, how should one select i) a detection algorithm, and ii) associated hyperparameter values (jointly called a model)? Effective outlier model selection is essential as different algorithms may work well for varying detection tasks, and moreover their performance can be quite sensitive to the values of the hyperparameters (HPs). On the other hand, unsupervised model selection is notoriously difficult, in the absence of hold-out validation data with ground-truth labels. Therefore, the problem is vastly understudied in the outlier mining literature. There exists a body of work that propose *internal model evaluation strategies* for selecting a model. These so-called internal strategies solely rely on the input data (without labels) and the output (outlier scores) of the candidate models. In this paper, we first survey internal model evaluation strategies including both those proposed specifically for outlier detection, as well as those that can be adapted from the unsupervised deep representation learning literature. Then, we investigate their effectiveness empirically in comparison to simple baselines such as random selection and the popular state-of-the-art detector Isolation Forest (iForest) with *default* HPs. To this end, we set up (and open-source) a large testbed with 39 detection tasks and 297 candidate models comprised of 8 different detectors and various HP configurations. We evaluate internal strategies from 7 different families on their ability to discriminate between models w.r.t. detection performance, without using any labels. Our study reports a striking finding, that *none of the existing and adapted strategies would be practically useful*: stand-alone ones are not significantly different from random, and consensus-based ones do not outperform iForest (w/ default HPs) while being more expensive (as *all* candidate models need to be trained for evaluation). Our survey stresses the importance of and the standing need for effective unsupervised outlier model selection, and acts as a call for future work on the problem.

1. INTRODUCTION

Model selection aims to select a model from a set of candidate models for a task, given data. We consider the model selection problem for the *unsupervised* outlier detection (UOD) task. Specifically, given a dataset for UOD, how can we identify – *without using any labels* – which outlier model

(a detection algorithm and the value(s) of its hyperparameter(s)) performs better than the others on the input dataset? Importantly, note that as the outlier detection task is unsupervised, so is the model selection task. That is, an outlier model is to be selected without being able to validate any candidate models on hold-out labeled data.

Motivation and Challenges. The notion of a universally “best” outlier model does not exist; rather the best-performing model depends on the given data. On the other hand, model selection is a nontrivial one, provided there are numerous outlier detection algorithms based on a variety of approaches: distance-based [20; 42], density-based [5; 14; 49], angle-based [23], ensemble-based [3; 32; 41], most recently deep neural network (NN) based [8; 10; 46; 51], and so on. To add to this “choice paralysis”, most models are sensitive to their choice of hyperparameters (HPs) with significant variation in performance [15], even more so for *deep* neural network (NN) based outlier models that have a long list of HPs [10]. Unsupervised model selection will likely be an increasingly pressing problem for deep detectors, as their complexity and expressiveness grow. Recent work hold out some labeled validation data for tuning such deep outlier models [27; 28; 46], which however is not feasible for fully unsupervised settings. These factors make outlier model selection a problem of utmost importance.

Despite its importance, the problem of Unsupervised Outlier Model Selection (UOMS hereafter) is a notoriously challenging one. Mainly, the absence of validation data with labels makes the problem hard. Moreover, there does not exist a universal or well-accepted objective criterion (i.e. loss function) for outlier detection, which makes model comparison infeasible. Besides its unsupervised nature, the search space for UOMS can be quite large with the arrival/popularity of deep NN based models with many HPs.

Existing work. Perhaps due to these challenges, UOMS remains a vastly understudied area. Most work in the outlier mining literature focus on designing new detection algorithms, such as those for unique settings: streaming [16; 34; 48], contextual [29; 37], human-in-the-loop [9; 25] detection, etc. There exist a small body of work, specifically addressing UOMS, that proposes *internal* (i.e. unsupervised) model evaluation strategies to assess the quality of a model and its output. These are called internal strategies as they use heuristic measures that solely make use of the input data and/or output outlier scores. To our knowledge, there are only three such techniques (in chronological order) [36; 13; 40]. However, they employ their proposed strategies to select only among 2-3 detectors on 8-12 real-world datasets. More

[‡]Equal contribution.

problematically, they do not systematically compare to one another, nor do they use the same datasets. This makes it difficult to fully understand the strengths and limitations of these existing methods, and ultimately the extent to which progress has been made on this subject.

More recently, we have designed a series of new approaches for UOMS leveraging two key concepts; meta-learning [57; 54; 55] and hyper-ensembles [10]. Notably, the former works utilize internal evaluation measures, the focus of this survey. The idea is to boost the relatively weak internal model performance signals from these heuristic measures via meta-learning from historical tasks that have labels. As such, any future work on designing new internal model evaluation strategies and improving their effectiveness and speed would directly feed into and advantage these meta-learning based UOMS approaches. (See Sec. 5 for detailed related work.)

Our survey. In this work, our goal is to survey internal model evaluation strategies, and systematically evaluate and compare their effectiveness. To this end, we first bring under one umbrella the aforementioned three existing UOMS methods, adapt two state-of-the-art unsupervised model selection techniques originally proposed for deep representation learning [11; 30], and design two new internal model selection methods inspired by various consensus algorithms. We put them to test on a large testbed against simple baselines, including random model selection as well as the popular isolation Forest (iForest) [32], with default HPs. To our knowledge, this is the first work to systematically review and evaluate the internal evaluation strategies toward unsupervised model selection for outlier detection. We summarize the contributions and findings of this paper as follows.

- **Unified Comparison:** We identify (to our knowledge) all existing internal model evaluation strategies for UOMS. For the first time, we systematically compare them on their ability to discriminate between models w.r.t. detection performance, as well as w.r.t. running time, on the same testbed.
- **Large-scale Evaluation:** Our testbed consists of 8 state-of-the-art detectors, each configured by a comprehensive list of hyperparameter settings, yielding a candidate pool of 297 models. We perform the model selection task on 39 independent real-world datasets from two different public repositories. We compare different strategies through paired statistical tests to identify significant differences, if any. We find that all three existing strategies specifically designed for UOMS are ill-suited. Alarming, none of them is significantly different from random selection (!)
- **New UOMS Techniques:** All three existing methods specifically designed for UOMS are *stand-alone*; evaluating each model individually, independent of others. In addition to those, we repurpose four *consensus-based* algorithms from other areas for UOMS; utilizing the agreements among the models in the pool. We find that consensus-based methods are more competitive than stand-alone ones, and all of them achieve significantly better performance than random. However, they are not different from iForest (the best detector in our pool), thus, would not be employed (on a pool) over training a single (iForest) model.
- **Open-source Testbed:** We expect that UOMS will continue to be a pressing problem, especially with the advent of deep detection models with many hyperpa-

rameters. Our large-scale analysis reveals that there is ample room for progress on this problem, and serves as a call for future work. At the same time, our results shed light onto the strengths and limitations of different approaches that motivate future directions.

Reproducibility and Future Work. To foster progress on this key problem, we open-source all datasets, our trained model pool, and implementations of all the internal model evaluation strategies at <https://github.com/yzhao062/uoms>.

2. PRELIMINARIES & THE PROBLEM

Model selection concerns with picking a model from a pool of candidate models. Let $\mathcal{M} = \{M_i\}_{i=1}^N$ denote a pre-specified pool of N models. Here each model M_i is a `{detector, HPconfiguration}` pair, where `detector` is a certain outlier detection algorithm (e.g. Local Outlier Factor (LOF) [6]) and `HPconfiguration` is a certain setting of the values for its hyperparameter(s) (e.g. for LOF, value of `n_neighbors`: number of nearest neighbors to consider, and function of choice for `distance` computation).

In this study, \mathcal{M} is composed by pairing 8 popular outlier detection algorithms to distinct hyperparameter choices, comprising a total of $N = 297$ models, as listed in Table 1. All models are trained based on the Python Outlier Detection Toolbox (PyOD) [56] on each dataset.

Let $\mathcal{D} = \{D_t\}_{t=1}^T$ denote the set of outlier detection datasets (i.e. tasks), where $D_t = \{\mathbf{x}_j^{(t)}\}_{j=1}^{n_t}$, $n_t = |D_t|$ is the number of samples and o_t is the true number of ground-truth outliers in D_t . We denote by $\mathbf{s}_i^{(t)} \in \mathbb{R}^{n_t}$ the list of outlier scores output by model M_i when employed (i.e. trained¹) on D_t , and $s_{ij}^{(t)} \in \mathbb{R}$ to depict individual sample j 's score. We omit the superscript when it is clear from context. W.l.o.g. the higher the s_{ij} is, the more anomalous is j w.r.t. M_i .

PROBLEM 1 (UNSupervised OUtlier MOdel SElection). (UOMS) *The model selection problem for unsupervised outlier detection can be stated as follows.*

- Given* an unsupervised detection task $D = \{\mathbf{x}_j\}_{j=1}^n$,
all models in \mathcal{M} trained on D
with corresponding output scores $\{\mathbf{s}_i\}_{i=1}^N$;
Select a model $M' \in \mathcal{M}$,
such that \mathbf{s}' yields good detection performance.

Note that the detection performance is to be quantified *post* model selection, where ground-truth labels are used only for evaluation (and **not** for model training or model selection). In this work, we study 7 different families of internal strategies (See Table 2): (1) three techniques that were proposed to directly address the UOMS problem, (2) two unsupervised model selection techniques adopted from deep learning, and (3) two others that are not originally designed for model selection that we adapt to UOMS.

To compare their effectiveness systematically, we construct a large testbed of $T = 39$ real-world outlier detection datasets from two different repositories (See Sec. 4.1). That is, we perform UOMS using each technique 39 times, to select one model from the pool of 297. Given that the datasets are independent, a large testbed enables paired statistical tests that conclusively identify significant differences between these techniques and various simple baselines.

¹Note that as we consider unsupervised outlier detection, model “training” does not involve any ground-truth labels.

Table 1: Outlier Detection Models; see hyperparameter definitions from PyOD

Detection algorithm	Hyperparameter 1	Hyperparameter 2	Total
LOF [5]	n_neighbors: [1, 5, 10, 15, 20, 25, 50, 60, 70, 80, 90, 100]	distance: ['manhattan', 'euclidean', 'minkowski']	36
kNN [42]	n_neighbors: [1, 5, 10, 15, 20, 25, 50, 60, 70, 80, 90, 100]	method: ['largest', 'mean', 'median']	36
OCSVM [47]	nu (train error tol): [0.1, 0.2, 0.3, 0.4, 0.5, 0.6, 0.7, 0.8, 0.9]	kernel: ['linear', 'poly', 'rbf', 'sigmoid']	36
COF [49]	n_neighbors: [3, 5, 10, 15, 20, 25, 50]	N/A	7
ABOD [23]	n_neighbors: [3, 5, 10, 15, 20, 25, 50]	N/A	7
iForest [32]	n_estimators: [10, 20, 30, 40, 50, 75, 100, 150, 200]	max_features: [0.1, 0.2, 0.3, 0.4, 0.5, 0.6, 0.7, 0.8, 0.9]	81
HBOS [14]	n_histograms: [5, 10, 20, 30, 40, 50, 75, 100]	tolerance: [0.1, 0.2, 0.3, 0.4, 0.5]	40
LODA [41]	n_bins: [10, 20, 30, 40, 50, 75, 100, 150, 200]	n_random_cuts: [5, 10, 15, 20, 25, 30]	54

297

3. REVIEW OF INTERNAL MODEL EVALUATION STRATEGIES

Internal strategies evaluate the goodness of a model without using any external information, especially with no access to ground-truth labels. The internal information being used is solely limited to (i) the input samples (feature values only), (ii) the trained models in the candidate pool and the outlier scores as output by these trained models. The common thread among all internal model evaluation strategies in this study is an estimated heuristic *internal measure* of “model goodness”. Model selection is then addressed by top-1 selection: i.e. picking the model with the highest value of the respective measure.

We categorize the 7 strategies we studied into two, depending on how they estimate their internal measure: (1) **stand-alone** and (2) **consensus-based**. (See Table 2.) Stand-alone strategies solely rely on each model and its output individually, independent of other models. All three existing methods proposed specifically for UOMS fall into this category. On the other hand, consensus-based strategies leverage agreement between the models in the pool and hence utilize candidate models collectively. Four strategies we adopt and adapt² from other areas all fall into this latter category.

In the following we provide a short description of each strategy (and refer to the original articles for full details). We also remark on the computational complexity of some methods as they demand considerable running time. Ideal is to have a lightweight and effective selection method with low overhead incurred on top of model training. In the experiments, we compare these methods w.r.t. their selection performance as well as running time.

3.1 Stand-alone Internal Evaluation (Existing)

3.1.1 IREOS

The first known index proposed for the internal evaluation of outlier detection results is by Marques et al., called Internal, Relative Evaluation of Outlier Solutions (IREOS) [36]. While their initial index is designed only for binary solutions (referred to as “top-n” detection), their recent work [35] generalized to numeric outlier scorings, which is the setting considered in this study.

Their intuition is that an outlier should be more easily separated (discriminated) from other samples than an inlier.

²We *adopt* two strategies originally proposed for unsupervised model selection for deep representation learning “as is”, and *adapt* two techniques (from information retrieval and ensemble learning) by repurposing them to UOMS problem with small modifications.

Table 2: Overview of UOMS methods studied in this survey.

Method	Type	Based on	Strategy
XB,RS,... [40]	Stand-alone	Outlier scores	Cluster quality
EM, MV [13]	Stand-alone	Outlier scores	Level sets
IREOS [36]	Stand-alone	O. scores + Input	Separability
UDR [11]	Consensus	Outlier scores	One-shot
MC [30], MC _S	Consensus	Outlier scores	One-shot
HITS [19]	Consensus	Outlier scores	Iterative
ENS [58]	Consensus	Outlier scores	Iterative

Then, a model is “good” the more it identifies as outlier those samples with a large degree of separability. They propose to assess the separability of each individual sample using a maximum-margin classifier (and specifically use nonlinear SVMs).³ The IREOS score of a model M_i on a given dataset is computed as

$$\text{IREOS}(s_i) = \frac{1}{n_\gamma} \sum_{l=1}^{n_\gamma} \frac{\sum_{j=1}^n p(\mathbf{x}_j, \gamma_l) w_{ij}}{\sum_{j=1}^n w_{ij}} \quad (1)$$

where $p(\mathbf{x}_j, \gamma_l)$ is the separability of sample j as estimated by a nonlinear SVM with kernel bandwidth (a hyper-parameter) γ_l , and n_γ is the number of different bandwidth values used from the interval $[0, \gamma_{\max}]$.⁴ They convert outlier scores $\{s_{ij}\}_{j=1}^n$ to probability weights $\{w_{ij}\}_{j=1}^n$ using the approach by [22] to push inlier scores toward zero so that they do not in aggregate dominate the weighted sum. IREOS tends to give high scores to those models whose outlier scores correlate well with the separability scores by a nonlinear SVM.

Computationally, IREOS is quite demanding as it requires training of a nonlinear classifier *per sample*. Their source code¹⁰ provides ways to approximate IREOS scores, mainly estimating separability via nearest neighbor distances, which however are also expensive to compute.

3.1.2 Mass-Volume (MV) and Excess-Mass (EM)

Goix [13] proposed using statistical tools, namely MV and EM curves, to measure the quality of a scoring function. Formally, a scoring function $s : \mathbb{R}^d \mapsto \mathbb{R}_+$ is any measurable function integrable w.r.t. the Lebesgue measure $\text{Leb}(\cdot)$, whose level sets are estimates of the level sets of the density. Outliers are assumed to occur in the tail of the score distribution as produced by a scoring function, where the *lower* $s(\mathbf{x})$ is, the more abnormal is \mathbf{x} .

³Note that collective outliers (forming micro-clusters, or clumps) do not have high separability. IREOS accounts for this effectively, provided a user-specified *clump.size*. For details, we refer to the original articles.

⁴They use heuristics to automatically set γ_{\max} in code.

Given a scoring function $s(\cdot)$ (in our context, an outlier model), the MV measure is defined as follows.

$$\widehat{MV}_s(\alpha) = \inf_{u \geq 0} \text{Leb}(s \geq u) \text{ s.t. } \mathbb{P}_n(s(\mathbf{X}) \geq u) \geq \alpha \quad (2)$$

where $\alpha \in (0, 1)$, and \mathbb{P}_n is the empirical distribution; $\mathbb{P}_n(s \geq v) = \frac{1}{n} \sum_{j=1}^n \mathbb{1}_{s(\mathbf{x}_j) > v}$.

For univariate real numbers, $\text{Leb}(\cdot)$ measures the length of the given interval. Let s_{\max} denote the largest score produced by $s(\cdot)$. Then, empirically $\text{Leb}(s \geq u)$ is equal to the length $|s_{\max} - u|$. Given α , the u that minimizes the Lebesgue measure $\text{Leb}(s \geq u)$ in Eq. (2) would be equal to the outlier score at the $(1-\alpha)$ -th quantile, i.e. $u = CCDF_s^{-1}(\alpha)$. Then, $|s_{\max} - u|$ would give the length of the range of scores for α fraction of the samples with scores larger than u . In their work, they consider $\alpha \in (0.9, 0.999)$.^{5,6} As they assume a *lower* score is more anomalous, the Lebesgue measure quantifies the length of the interval of scores for the inliers. The smaller MV is, the better the scoring function is deemed to be. Intuitively, then, MV measures the clusteredness of inlier scores (or the compactness of high-density level sets). The EM measure is quite similar, and is defined as

$$\widehat{EM}_s(t) = \sup_{u \geq 0} \mathbb{P}_n(s(\mathbf{X}) \geq u) - t \text{Leb}(s \geq u) \quad (3)$$

for $t > 0$. Similarly, they consider $t \in [0, \widehat{EM}_s^{-1}(0.9)]$ with $\widehat{EM}_s^{-1}(0.9) := \inf\{t \geq 0, \widehat{EM}_s(t) \leq 0.9\}$.

Intuitively, EM would identify as small a u value as possible (so as to maximize the density mass in the first term) such that the scores larger than or equal to u are as clustered as possible (so as to minimize the Lebesgue measure in the second term). Again, the more clustered are the scores of the bulk of the samples (i.e. inliers), the larger EM gets, and the better the scoring function is deemed to be.

3.1.3 Clustering validation metrics

[40] point out that a drawback of IREOS, besides computational demand, is its dependence on classification – which itself introduces a model selection problem – since the results may depend on the selected classification algorithm and its hyper-parameter settings.⁷ Despite citing IREOS, they do not compare in experiments.

Their key proposal is to apply internal validation measures for clustering algorithms to outlier detection. As clustering aims to ensure samples within each cluster are similar and different from samples in other clusters, these measures are mainly based on compactness (capturing within-cluster similarity) and/or separation (reflecting inter-cluster distance).

To that end, we split the outlier scores by a given model under evaluation for dataset D_t into two clusters, denoted C_o and C_i , respectively consisting of the highest o_t scores and the rest. According to those measures, an outlier model is “good” the more separated these two sets of scores are and/or the more clustered the scores within each set are.

⁵Given fraction of outliers is bounded to 10% maximum.

⁶Area under the MV-curve is estimated as the sum of empirical MV values by Eq. (2) for discretized values of α .

⁷Another paper [39] by the same authors proposed a classification based internal evaluation method, similar to IREOS. Their experiments show that the current internal measures do comparably well or better with less computational overhead, hence we omit [39] from this study.

In their study, they compared 10 different existing clustering quality measures, such as the Silhouette index [45], Xie-Beni index [52], etc. (See others in the original article.) One of the well-performing ones in our experiments, namely Xie-Beni index of a model M_i , denoted XB_i , is defined as follows.

$$\text{XB}_i = \frac{\sum_{j \in C_o} d^2(s_{ij}, c_o) + \sum_{j' \in C_i} d^2(s_{ij'}, c_i)}{n_t d^2(c_o, c_i)} \quad (4)$$

where $c_o = \sum_{j \in C_o} s_{ij}/o_t$ and $c_i = \sum_{j' \in C_i} s_{ij'}/(n_t - o_t)$ depict the cluster centers and $d(\cdot, \cdot)$ is the Euclidean distance. This index can be interpreted as the ratio of the intra-cluster compactness to the inter-cluster separation.

Clustering quality based measures are typically easy to compute; most of them being linear in the number of samples.

3.2 Consensus-based Internal Evaluation (Re-purposed)

3.2.1 UDR

The first consensus-based approach, namely Unsupervised Disentanglement Ranking (UDR), is adopted from deep learning and is “the first method for unsupervised model selection for variational disentangled representation learning” [11]. Each model in their case corresponds to a $\{\text{HPconfiguration}, \text{seed}\}$ pair. Reciting Tolstoy who wrote “Happy families are all alike; every unhappy family is unhappy in its own way.”, their main hypothesis is that a model with a good hyperparameter (HP) setting will produce *similar results* under different random initializations (i.e. seeds) whereas for a poor HP setting, results based on different random seeds will look arbitrarily different.

In a nutshell, UDR follows 4 steps: (1) Train $N = H \times S$ models, where H and S are the number of hyperparameter settings and random seeds, respectively. (2) For each model M_i , randomly sample (without replacement) $P \leq S$ other models with the *same* HP as M_i , but *different* seeds. (3) Perform P pairwise comparisons between M_i and the models sampled in Step 2 for M_i . (4) Aggregate pairwise similarity scores (denoted $UDR_{ii'}$) as $UDR_i = \text{median}_{i'} UDR_{ii'}$, for $i = 1, \dots, N$. Finally, they pick the model (among N) with the largest UDR_i . Intuitively, UDR selects a model with an HP setting that yields stable or *consistent* results across various seeds.

Notice that adopting UDR for the UOMS task is trivial by making the analogy between $\{\text{HPconfiguration}, \text{seed}\}$ and $\{\text{detector}, \text{HPconfiguration}\}$. While trivially applied, one may question whether the implied hypothesis (that a good detector has consistent results across different HP settings) holds true for outlier models, since one of the key reasons for UOMS in the first place is that most detectors are sensitive to their HP settings [15].

The key part of UDR is how pairwise model comparisons are done in Step 3. We measure the output *ranking similarity* of the samples by two models, based on three well-known measures from information retrieval [31] (See Sec. 4.1).

3.2.2 MC

A follow-up work to UDR proposed ModelCentrality (MC), which is another consensus-based strategy for what they call “self-supervised” model selection for disentangling GANs [30]. Their premise is similar, that “well-disentangled models should be close to the optimal model, and hence also close

to each other”. Provided the similarity $B_{ii'}$ between two models M_i and $M_{i'}$ can be computed, ModelCentrality of M_i is written as $MC_i = \frac{1}{N-1} \sum_{i' \neq i} B_{ii'}$. They then select the model with the largest MC_i , which coincides with the medoid in the pool of models – hence the name MC.

Computationally, MC is quadratic in the number of models as it requires all pairwise comparisons. We also experiment with a lightweight version, called MC_S , where we randomly sample $P \leq N$ models and compute MC_i of M_i as the average of its similarities to P models, effectively reducing its complexity down to that of UDR.

In their experiments, [30] report that MC outperforms UDR schemes (Sec. 3.2.I). Our results are consistent with their finding, possibly because it is an unrealistic hypothesis for outlier models that a good model would have consistent results across HP settings.

3.2.3 Model Centrality by HITS

We can build on the idea of ModelCentrality through computing centrality in a network setting. Unlike MC that is computed in one shot, network centrality is *recursive*—wherein a node has higher centrality the more they point to nodes that are pointed by other high-centrality ones.

One of the earliest methods for computing centrality, namely hubness h_p and authority a_p , of pages on the Web is the HITS algorithm [19], where

$$h_p \propto \text{sum of } a_i \text{ for all nodes } i \text{ that } p \text{ points to, and}$$

$$a_p \propto \text{sum of } h_i \text{ for all nodes } i \text{ pointing to } p,$$

which are estimated alternately over iterations until convergence. Besides ranking on the Web, HITS-like ideas have been used to estimate user trustworthiness in online rating platforms [24; 50], physician authoritativeness in patient referral networks [38], polarity of subjects in political networks [4], as well as truth discovery [53].

It is easy to adapt HITS for UOMS by constructing a complete bipartite network between the N models and n_t samples in a given dataset D_t . Then, the models can be evaluated by their hubness centralities. The analogous interpretation is that a sample has higher authority (outlierness), the more trusted models (with high hubness) point to it (with large outlierness score, i.e. large edge weight). Then, a model is more central or trusted, the more it points (with large outlierness score) to samples with high authority.

Note that a by-product of this strategy is a consensus-based ranking of the samples based on authority scores (i.e. centrality-based outlierness) upon convergence. We compare this (aggregate) ranking, called HITS-AUTH, against selecting a (single) model by hubness in the experiments.

3.2.4 Unsupervised outlier model ensembling

HITS has a built-in advantage that is the iterative refinement of model trustworthiness. Specifically, given the trustworthiness of models, outlier scores can be better estimated by a trustworthiness-weighted aggregation of scores across models. Then, given those refined outlier scores, model trustworthiness can also be better estimated; where the more similar their output is to the updated scores, the more a model is deemed trustworthy.

Here we build on another iterative scheme, originally designed for unsupervised selective outlier model ensembling [44; 58].

The idea is to infer reliable “pseudo ground truth” outlier scores via aggregating the output of a carefully-selected subset of trustworthy models. The ensemble is constructed bottom-up in a greedy iterative fashion (see Alg. 1).

Similar to HITS, the “pseudo ground truth” and model trustworthiness are estimated alternately. The latter is computed as the ranking based similarity of a model’s output to the “pseudo ground truth” (i.e. *target* in Alg. 1) at a given iteration. We adapt this framework to UOMS by using these similarities at convergence to evaluate the models. We call this strategy ENS. In experiments, we also compare the (aggregate) ranking by the ensemble (based on *target*), called ENS-PSEUDO, to selecting a (single) model (with highest similarity to *target*).

To wrap up, we give a summary of the 7 families of UOMS techniques as described in this section in Table 2.

4. EMPIRICAL EVALUATION OF INTERNAL MODEL EVALUATION STRATEGIES

4.1 Setup

Datasets and Model Pool. We already discussed the real-world datasets and candidate models of this study in Sec. 2. We build the experiments on **39 widely used outlier detection benchmark dataset**. As shown in Table 3, 21 datasets are from the **ODDS Library** [43], and the other 18 datasets are from **DAMI datasets** [7]. The specifications for all $N=297$ models have been listed in Sec. 2 Table 1.

Baselines. We compare the model selected by each technique (Sec. 3) to two baselines across datasets.

- RANDOM, whose performance is the average of all (297) models per dataset. This is equivalent to expected performance when selecting a model from the candidate pool at random.
- iFOREST-R, with performance as the average of all (81) iForest models in the pool, equivalent to using iForest [32] (a state-of-the-art ensemble detector) with randomly chosen hyperparameters.⁸

Method Configurations. For clustering-quality based measures, we split into two clusters as the top o_t (true number of outliers) and the rest, i.e. give those strategies the advantage of knowing o_t . This avoids the clustering step, which requires us to pick a clustering algorithm etc. and directly focuses on the measures themselves.

For EM and MV⁹, we use the default values for α and t respectively (See Sec. 3.1.2) and set **n_generated**=100K, which is the number of random samples to generate for estimating the null distributions.

For IREOS, we use the author recommended settings;¹⁰ $\gamma_{\max} := \text{findGammaMaxbyDistances}(\cdot)$ with **sampling**=100, **tol**= 5×10^{-3} , and **clump.size**=10.

For UDR, MC, and MC_S , we experiment with three different pairwise similarity measures: Spearman’s ρ , Kendall’s τ , and NDCG [31]. For MC_S , $P = \sqrt{N} \approx 18$.

For HITS and ENS, we set edge weights between model M_i

⁸Family-wise performances across datasets in Supp. A show that iForest is the most competitive among the 8 families of detectors in this study, and thus the strongest baseline.

⁹Code available at https://github.com/ngoix/EMMV_benchmarks

¹⁰We thank Henrique Marques who helped with running their source code, <https://github.com/homarques/ireos-extension>

Table 3: Real-world dataset pool composed by ODDS library (21 datasets) and DAMI library (18 datasets).

Dataset	Num Pts	Dim	% Outlier
1 annthyroid (ODDS)	7200	6	7.416
2 arrhythmia (ODDS)	452	274	14.601
3 breastw (ODDS)	683	9	34.992
4 glass (ODDS)	214	9	4.205
5 ionosphere (ODDS)	351	33	35.897
6 letter (ODDS)	1600	32	6.250
7 lympho (ODDS)	148	18	4.054
8 mammography (ODDS)	11183	6	2.325
9 mnist (ODDS)	7603	100	9.206
10 musk (ODDS)	3062	166	3.167
11 optdigits (ODDS)	5216	64	2.875
12 pendigits (ODDS)	6870	16	2.270
13 pima (ODDS)	768	8	34.895
14 satellite (ODDS)	6435	36	31.639
15 satimage-2 (ODDS)	5803	36	1.223
16 speech (ODDS)	3686	400	1.654
17 thyroid (ODDS)	3772	6	2.465
18 vertebral (ODDS)	240	6	12.500
19 vowels (ODDS)	1456	12	3.434
20 wbc (ODDS)	378	30	5.555
21 wine (ODDS)	129	13	7.751
22 Annthyroid (DAMI)	7129	21	7.490
23 Arrhythmia (DAMI)	450	259	45.777
24 Cardiocography (DAMI)	2114	21	22.043
25 HeartDisease (DAMI)	270	13	44.444
26 InternetAds (DAMI)	1966	1555	18.718
27 PageBlocks (DAMI)	5393	10	9.456
28 Pima (DAMI)	768	8	34.895
29 SpamBase (DAMI)	4207	57	39.909
30 Stamps (DAMI)	340	9	9.117
31 Wilt (DAMI)	4819	5	5.333
32 ALOI (DAMI)	49534	27	3.044
33 Glass (DAMI)	214	7	4.205
34 PenDigits (DAMI)	9868	16	0.202
35 Shuttle (DAMI)	1013	9	1.283
36 Waveform (DAMI)	3443	21	2.904
37 WBC (DAMI)	223	9	4.484
38 WDBC (DAMI)	367	30	2.724
39 WPBC (DAMI)	198	33	23.737

and sample j in a dataset as $1/r_{ij}$, where r_{ij} is the position of j in the rankedlist by M_i . Raw outlier scores are not used as they are not comparable across models. For comparison between selection versus consensus/ensembling, we also report the performance of the consensus outcome, called HITS-AUTH and ENS-PSEUDO; as ranked (resp.) by authority scores and by the pseudo ground truth at convergence.

Performance metrics. We evaluate performance w.r.t. three metrics. Two are based on the ranking quality: Average Precision (**AP**): the area under the precision-recall curve and **ROC AUC**: the area under the recall-false positive rate curve. The third metric measures the quality at the top: **Prec@ k** , precision at top k where we set $k = o_t$ (i.e. true number of outliers) for each $D_t \in \mathcal{D}$. In Supp. B we show that performances vary considerably across models for most datasets, justifying the importance of model selection. For brevity, all results in this section are w.r.t. AP. Corresponding results for other metrics are similar, all of which are provided in Supp. C.

4.2 Results

4.2.1 Cluster quality based methods

We start by studying the 10 cluster quality based methods to identify those that stand out. We report the p -values by the

one-sided¹¹ paired Wilcoxon signed rank test in Table 4. STD is significantly worse than all other methods. Three strategies that stand out are RS, CH, and XB, which are identical; in the sense that despite differences in their values and overall ranking, they select exactly the same model on each dataset. Importantly, while both STD and S are significantly worse than RANDOM at $p = 0.05$, none of the others is significantly different from RANDOM (!) All methods (including XB, RS, and CH) are significantly worse than IFOREST-R.

Table 4: Comparison of cluster quality based methods and baselines by one-sided paired Wilcoxon signed rank test. p -values **bolded** (underlined) highlight the cases where the “row-method” is significantly **better** (worse) than the “column-method” at $p \leq 0.05$.

	STD	H	S	I	DB	SD	D	RND	iF
XB,RS,CH	0.004	0.240	0.038	0.212	0.370	0.127	0.357	0.500	<u>0.981</u>
STD		<u>0.997</u>	<u>0.961</u>	<u>0.997</u>	<u>0.982</u>	<u>0.967</u>	<u>0.999</u>	<u>1.000</u>	<u>1.000</u>
H			<u>0.373</u>	0.500	0.725	0.379	0.675	0.849	<u>0.996</u>
S				0.627	0.949	0.557	0.881	0.953	<u>0.999</u>
I					0.730	0.384	0.742	0.882	<u>0.997</u>
DB						0.307	0.647	0.522	<u>0.982</u>
SD							0.823	0.910	<u>0.995</u>
D								0.572	<u>0.990</u>
RND									<u>1.000</u>

These findings suggest that cluster quality based internal evaluation methods would not be useful for UOMS.

4.2.2 Other stand-alone methods

As discussed in Sec. 3.1.2, EM and MV quantify (roughly) the clusteredness of the inlier scores. Therefore, they are conceptually similar to the clustering quality based methods. Our findings confirm this intuition. As shown in Table 5, there is no significant difference between EM/MV and XB/RS/CH or RANDOM. Both of them are also significantly worse than IFOREST-R. Thus, they do not prove useful for UOMS. Findings are similar for IREOS; despite using more information (input samples besides scores, see Eq. (1)) and computational cost, it is only comparable to RANDOM.

Table 5: Comparison of stand-alone methods and baselines.

	EM	MV	IREOS	RND	iF
XB,RS,CH	0.533	0.500	0.862	0.500	<u>0.981</u>
EM		0.079	0.642	0.539	<u>0.979</u>
MV			0.716	0.687	<u>0.994</u>
IREOS				0.303	0.908

We provide an additional viewpoint by identifying the q -th best model per dataset where there exists no significant difference between the performance of the q -th best model and that selected by a given UOMS strategy across datasets. We report the *smallest* q for which one-sided Wilcoxon signed rank test yields $p > 0.05$ in Table 6. A method with smaller q is better; the interpretation being that it could select, from a pool of 297, the model that is as good as the q -th best model per dataset. Stand-alone methods do not fare well against IFOREST-R which is comparable to the 84-th best model.

4.2.3 Consensus-based methods

We first study one-shot methods UDR, MC, and MC_S based on different similarity measures. As shown in Table 6, all

¹¹Testing the hypothesis: row-method is better than col-method (against the null hypothesis stating no difference). For reverse order, p -value = 1 minus the reported value.

Table 6: **Summary of results:** p -values by one-sided paired Wilcoxon signed rank test comparing UOMS methods to the baselines, smallest q -th best model with no significant difference, and mean/standard deviation AP across datasets.

	Method	RANDOM	iFOREST-R	q_{AP}	mean AP	std AP
S-alone	XB,RS,CH	0.500	<u>0.981</u>	127	0.354	0.298
	EM	0.539	<u>0.979</u>	115	0.322	0.265
	IREOS	0.303	0.908	99	0.335	0.261
Consensus-based	UDR- ρ	0.012	0.905	104	0.383	0.283
	UDR- τ	0.019	<u>0.952</u>	109	0.379	0.282
	UDR- $NDCG$	0.004	0.825	93	0.384	0.270
	MC- ρ	0.000	0.217	89	0.395	0.289
	MC- τ	0.002	0.062	81	0.396	0.297
	MC- $NDCG$	0.000	0.182	82	0.404	0.291
	MC $_S$ - ρ	0.007	0.706	108	0.385	0.289
	MC $_S$ - τ	0.001	0.599	90	0.397	0.305
	MC $_S$ - $NDCG$	0.001	0.205	83	0.391	0.285
	HITS	0.000	0.494	95	0.397	0.299
Agg.	HITS-AUTH	0.000	0.577	94	0.401	0.286
	ENS-PSEUDO	0.001	0.422	79	0.373	0.282
Base.	RANDOM	-	<u>1.000</u>	144	0.342	0.234
	iFOREST-R	-	-	84	0.399	0.300

versions provide similar results, which are significantly better than RANDOM, and not different from iFOREST-R. We note that the faster, sampling-based MC $_S$ achieves similar performance to MC and can be used as a practical alternative.

Iterative methods HITS and ENS produce similar results to these simple one-shot methods, despite aiming to refine estimates of model trustworthiness over iterations. Again, as shown in Table 6, they significantly outperform RANDOM and are comparable to iFOREST-R. The same holds true for their respective consensus scores, HITS-AUTH and ENS-PSEUDO, where model aggregation provides no significant advantage over selecting the best (single) model.

Table 7 shows a pairwise comparison of the consensus-based methods by one-sided Wilcoxon signed rank test, confirming mostly no significant difference between them.

Table 7: Comparison of consensus-based methods (UDR, MC, MC $_S$ are based on $NDCG$).

	MC	MC $_S$	HITS	ENS
UDR	0.810	0.364	0.739	0.400
MC		0.551	0.039	0.116
MC $_S$			0.296	0.369
HITS				0.753

4.2.4 Running time analysis

In Fig. 1 we present for each method the running times on all datasets.¹² IREOS and EM/MV are both computationally demanding, while ineffective. In fact, IREOS takes more than 16 days (!) on the largest dataset (ALOI), due to kernel SVM training *for each sample*. MC is the next most expensive method, which is quadratic in the number of models, although still takes less than 1 hr on ALOI. In short, MC $_S$, ENS, and especially HITS prove to be both competitive as well as fast UOMS methods, completing within 10 minutes on our testbed.

¹²On an Intel Xeon E7 4830 v3 @ 2.1Ghz with 1TB RAM.

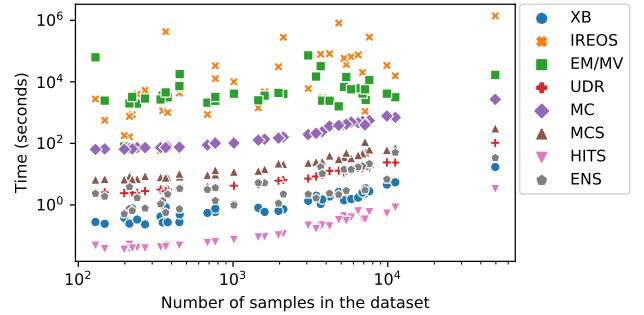


Figure 1: Run time comparison of UOMS methods.

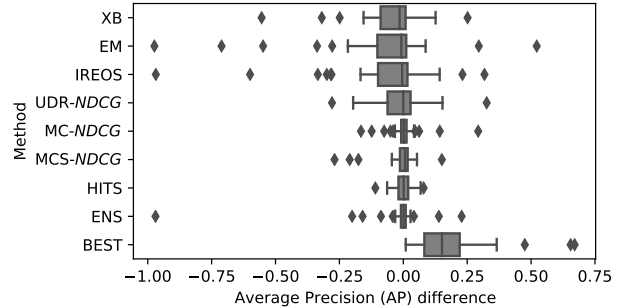


Figure 2: Distribution of performance difference across datasets: AP of selected model (by each UOMS method studied) minus that of iFOREST-R. Stand-alone methods and UDR are subpar, whereas other consensus-based method differences concentrate around zero (indicating no notable difference from iFOREST-R). Also shown for comparison is BEST model on each dataset, showcasing ample room for improvement over iFOREST-R.

4.3 Discussion of the Results

We conclude with the two key take-aways from our study:

1. *None of the existing (stand-alone) UOMS methods is significantly different from random model selection (!), and with the exception of IREOS. All are significantly worse than iForest (with random hyperparameter configuration). The slight advantage of IREOS can be attributed to it utilizing input features in addition to model outlier scores, unlike the other strategies that solely use the output scores. However, this advantage comes at the expense of significant running time.*
2. *All consensus-based methods that we repurposed for UOMS are significantly better than random selection, but not different from the fast, state-of-the-art iForest detector with default HPs.*

Fig. 2 illustrates these take-aways where we show, via box-plots, the distribution of the performance difference between the model selected by each UOMS method and iFOREST-R across datasets. Consensus-based methods select models at best as good as iFOREST-R, where the AP difference concentrates around zero, whereas others are inferior.

These results suggest that **none of the UOMS methods we studied would be useful in practice**; because one would not first train a large *pool* of models – which would

Algorithm 1 Ensemble-based Internal Model Evaluation

Input: set of outlier scores from all models, $\{s_i\}_{i=1}^N$
Output: internal scores for all models

```
1:  $\mathcal{S} := \emptyset$ ,  $\mathcal{E} := \emptyset$ ,  $C := 0$ 
2: for  $i = 1, \dots, N$  do           ▶ convert scores to inverse rank
3:    $\mathcal{S} := \mathcal{S} \cup \{1/\text{rank}(s_{ij})\}_{j=1}^n$ 
4: end for
5:  $target := \text{avg}(\mathcal{S})$            ▶ initial pseudo ground truth scores
6: repeat
7:   sort  $\mathcal{S}$  by rank correlation to  $target$  in desc. order
8:    $\{m, corr_m\} := \text{fetchFirst}(\mathcal{S})$ 
9:   if  $corr(\text{avg}(\mathcal{E} \cup m), target) \times |E| \geq C$  then
10:     $\mathcal{E} := \mathcal{E} \cup m$ ,  $C += corr_m$ 
11:     $target := \text{avg}(\mathcal{E})$            ▶ pseudo ground truth by  $\mathcal{E}$ 
12:   end if
13: until  $\{\mathcal{S} = \emptyset$  or  $\mathcal{E}$  is not updated $\}$ 
14: return rank correlation of  $s_i$  to  $target$ ,  $i = 1, \dots, N$ 
```

incur considerable computation – and then run a post hoc UOMS method to select a model, only to achieve comparable performance to a *single* iForest model (with default configuration) – which, in contrast, is extremely fast to train as it builds randomized trees on subsamples of data.

However, this is not to conclude that iForest is the best that one can hope to do. As given in Table 6, iFOREST-R is only as good as the 84-th best model per dataset. While it is the most competitive detector on average, other families outperform iForest on 28 out of 39 datasets in our study w.r.t. AP (See Table 16 in Supp. A, also see Tables 17 and 18 respectively for ROC and Prec@ k). In Fig. 2 we also show the performance difference of the 1-st BEST model per dataset from iFOREST-R. (Also see Fig 3 in Supp. C.) One can clearly recognize that there is considerable room for progress in the area of UOMS.

5. RELATED WORK

In the outlier mining literature, several evaluation and benchmarking surveys draw attention to the fact that most (classical) outlier detectors are sensitive to their HPs [2; 7; 15; 17]. This is even more so for the recently booming deep learning based detection methods, as we empirically studied recently [10]. Despite its critical importance, related work on unsupervised outlier model selection (UOMS) is limited, with only a few existing works addressing the problem.

In this survey, we focus on internal model evaluation strategies that we have reviewed in the previous section. In the following sections, we provide a critique and comparison between these strategies in §5.1 and 5.2. Recently, other novel approaches leveraging ensemble methods and meta-learning have been proposed for UOMS, which we also review in §5.3 for completeness. Additionally, one may choose good outlier detection (OD) models by assuming the underlying outlier characteristics of a dataset. However, in §5.4, we describe this approach and argue that it is impractical in real-world cases where UOMS is more feasible.

5.1 Internal model evaluation strategies for UOMS

Cluster quality based measures [40] and statistical mass based EM/MV methods [13] rely only on output scores. In contrast IREOS [35; 36] uses more information, that is both outlier scores and the original input samples (See Eq. (1)). Verifying that outlier scores align (correlate) with the separability of

samples in the feature space is potentially less error-prone than simply looking at whether outlier/inlier scores are well clustered or separated – e.g., a model that outputs a $\{0, 1\}$ score per point at random would be considered a good model by the latter. The trade-off is the computational overhead for quantifying separability per sample.

In their work, IREOS is employed for UOMS using only 2 detectors (LOF [5] and kNN [42]), each with 17 different HP configurations (for a total of 34 models) on 11 datasets. Being the seminal work, there is no comparison to any other techniques (existing or adapted). [40] acknowledge IREOS and criticize its computational demand, without any comparison. They also do not perform any UOMS in experiments, rather, they study the decay in internal measures as the ground truth ranking is contaminated via random swaps at the top based on 12 datasets. Finally, [13] performs UOMS using only and exactly 3 models (LOF, iForest [32], OCSVM [47]), each with a single (unspecified) HP configuration, on 8 datasets. None of these three compares to any other in their work. Moreover, because the datasets, experimental design, and the model pool specified by each work is different, it is not possible to do any direct comparison. In this work, we do a systematic comparison for the first time, using a much larger testbed (8 detectors, 297 models, 39 datasets) than originally considered by any prior work.

5.2 Internal Model Evaluation Strategies Repurposed for UOMS

All three existing methods for UOMS are stand-alone, evaluating a model independently from the others. Having trained all models among which to select from, it is reasonable to take advantage of the similarities/agreement among them. To this end, we have repurposed methods from unsupervised representation learning [11; 30], network centrality [19], and unsupervised ensemble learning [44; 58] all of which are based on the “collective intelligence” of the models in the pool.

As we show in the experiments, these strategies produce superior outcomes than existing, stand-alone methods. As such, our study motivates future work on consensus-based strategies, and calls for the transfer of prominent ideas from other similar fields, such as truth discovery and crowdsourcing, toward tackling the important problem of UOMS.

5.3 Ensembling and Meta-learning for UOMS

Most recently, two promising new directions have been explored toward UOMS. The first idea is building *hyper*-ensembles [10], which combines the outlier scores from multiple models with various HP configurations, rather than trying to select a single one of them. They have shown that the hyper-ensemble is significantly more robust to its own HPs (namely, the number of models to assemble and the value range per HP). The key challenge is similar to internal strategies covered in this paper, specifically, training all the models for assembly at test time is expensive, for which several speed up techniques have been proposed in [10].

The second line of work leverages *meta-learning* [54; 55; 57], where a database of historical outlier detection tasks *with labels* are used to transfer “knowledge/experience” toward UOMS on similar test tasks (without labels). Interestingly, internal evaluation measures have been exploited in (meta-)learning a mapping from such weak internal signals, dataset characteristics, etc. onto model performance (which can be computed for labeled historical tasks). Such a mapping is

then employed for model performance prediction on test tasks without the need to access labels. This suggests that new internal evaluation measures or any improvements that lead to stronger internal signals of performance are to boost these meta-learning based solutions to UOMS. Computationally, meta-learning approaches are also more feasible than hyper-ensembles, as most computation is off-loaded to the meta-learning phase while fewer models are trained at test time. Of course, fast yet effective internal measures would contribute to further speed up model selection on a new test task.

5.4 OD Model Selection by Data Characterizations

Our results, as presented in Supp. Table 16, 17, and 18, show that no single OD model can consistently outperform others across all datasets. Also, the studied UOMS methods are not yet sufficiently useful in selecting effective OD models. Given these observations, it is natural to wonder if we can choose OD models for a dataset following some “general guidelines” as an alternative to more systematic UOMS. Ideally, if we could identify the characteristic features of datasets based on their outliers, we could pinpoint the most effective OD algorithm(s) for each categorization, and even develop new ones. Several works propose such characterizations, including scattered-vs-clustered outliers [12; 18] and global-vs-local outliers [5; 17]. Other characterizations may include subspace outliers [21; 33] and extreme-value outliers [1]. Recently, self-supervised methods based on contrastive learning have augmented the data by assuming it resembles the anomaly generating process; for example, cut-and-paste augmentation creates images resembling industry defects [26]. Some recent literature demonstrates the potential of leveraging data characterization in choosing OD models. For example, Hand et al. [17] show that LOF performs well on simulated datasets with only local outliers, whereas kNN is most effective on datasets with only global outliers.

In principle, if one makes strong assumptions about the anomaly generation process and outlying characterizations, it is possible to choose or design a suitable algorithm without the need for model selection. However, this is not always the case since the generation process can be complex, and a dataset may contain multiple types of outliers [17]. Therefore, in this work, we take a broader perspective on the detection problem and avoid relying on methods and selection strategies that make explicit assumptions.

6. CONCLUSION & CALL FOR FUTURE WORK ON UOMS

In this review, we considered the unsupervised outlier model selection (UOMS) problem: Given an unlabeled outlier detection task, which detection algorithm and associated hyperparameter (HP) settings should one use? This is a question of utmost importance not only for practitioners to do well on their new task, but also for the research community for being able to fairly compare new detection methods and keep an accurate track record of progress in the field. On the other hand, the problem is notoriously hard in the absence of any labeled data, any well-accepted objective or loss function, and potentially very large model space especially for deep outlier detectors with many HPs.

We focused on the body of methods that proposed internal (i.e., unsupervised) model evaluation strategies that leverage

implicit signals from the input features and /or the output outlier scores alone. On a large testbed comprising 297 models and 39 real-world datasets, we evaluated 7 different families of such internal evaluation strategies against simple baselines. Strikingly, we found that while consensus-based strategies are more promising against stand-alone ones which are not significantly better than random, none of them provides significant improvement over the state-of-the-art iForest detector with default HPs.

Our findings call for further research in this important area. As our work recently showed [10], deep detectors are considerably poor across varying HPs *on average* (i.e. when HPs are chosen randomly in the absence of any other guidance). As such, UOMS appears to stand as the biggest obstacle in front of deep models to fulfill their potential for outlier detection. A promising future direction is to develop stronger and faster internal strategies that can be leveraged within a meta-learning framework as in [54; 55]. Our empirical evaluation revealed consensus-based internal strategies to be relatively more promising, which provides fertile ground for adaptation of prominent ideas from related areas such as truth discovery and crowdsourcing. To foster progress on this critical problem, we publicly share all source code, trained models, and datasets at <https://github.com/yzhao062/uoms>.

7. REFERENCES

- [1] C. C. Aggarwal and C. C. Aggarwal. *An introduction to outlier analysis*. Springer, 2017.
- [2] C. C. Aggarwal and S. Sathe. Theoretical foundations and algorithms for outlier ensembles. *Acm sigkdd explorations newsletter*, 17(1):24–47, 2015.
- [3] C. C. Aggarwal and S. Sathe. *Outlier Ensembles: An Introduction*. Springer Publishing Company, Inc., 1st edition, 2017.
- [4] L. Akoglu. Quantifying political polarity based on bipartite opinion networks. In *ICWSM*. The AAAI Press, 2014.
- [5] M. M. Breunig, H.-P. Kriegel, R. T. Ng, and J. Sander. Lof: identifying density-based local outliers. In *SIGMOD*, pages 93–104, 2000.
- [6] M. M. Breunig, H.-P. Kriegel, R. T. Ng, and J. Sander. Lof: Identifying density-based local outliers. In *SIGMOD*, pages 93–104. ACM, 2000. SIGMOD Record 29(2), June 2000.
- [7] G. O. Campos, A. Zimek, J. Sander, R. J. G. B. Campello, B. Micenková, E. Schubert, I. Assent, and M. E. Houle. On the evaluation of unsupervised outlier detection: measures, datasets, and an empirical study. *DMKD*, 30(4):891–927, July 2016.
- [8] J. Chen, S. Sathe, C. C. Aggarwal, and D. S. Turaga. Outlier detection with autoencoder ensembles. In *SDM*, pages 90–98. SIAM, 2017.
- [9] S. Das, W.-K. Wong, T. G. Dietterich, A. Fern, and A. Emmott. Incorporating expert feedback into active anomaly discovery. In *ICDM*, pages 853–858. IEEE Computer Society, 2016.

- [10] X. Ding, L. Zhao, and L. Akoglu. Hyperparameter sensitivity in deep outlier detection: Analysis and a scalable hyper-ensemble solution. In *Advances in Neural Information Processing Systems*, 2022.
- [11] S. Duan, L. Matthey, A. Saraiva, N. Watters, C. Burgess, A. Lerchner, and I. Higgins. Unsupervised model selection for variational disentangled representation learning. In *ICLR*. OpenReview.net, 2020.
- [12] A. Emmott, S. Das, T. Dietterich, A. Fern, and W.-K. Wong. A meta-analysis of the anomaly detection problem. *arXiv preprint arXiv:1503.01158*, 2015.
- [13] N. Goix. How to evaluate the quality of unsupervised anomaly detection algorithms? *CoRR*, abs/1607.01152, 2016.
- [14] M. Goldstein and A. Dengel. Histogram-based outlier score (hbos): A fast unsupervised anomaly detection algorithm. *KI-2012: Poster and Demo Track*, pages 59–63, 2012.
- [15] M. Goldstein and S. Uchida. A comparative evaluation of unsupervised anomaly detection algorithms for multivariate data. *PLoS one*, 11(4):e0152173, 2016.
- [16] S. Guha, N. Mishra, G. Roy, and O. Schrijvers. Robust random cut forest based anomaly detection on streams. In M. F. Balcan and K. Q. Weinberger, editors, *Proceedings of The 33rd International Conference on Machine Learning*, volume 48 of *Proceedings of Machine Learning Research*, pages 2712–2721, New York, New York, USA, 20–22 Jun 2016. PMLR.
- [17] S. Han, X. Hu, H. Huang, M. Jiang, and Y. Zhao. ADBench: Anomaly detection benchmark. In *Advances in Neural Information Processing Systems*, 2022.
- [18] S. Jiang, R. L. Cordeiro, and L. Akoglu. D. mca: Outlier detection with explicit micro-cluster assignments. In *2022 IEEE International Conference on Data Mining (ICDM)*, pages 987–992. IEEE, 2022.
- [19] J. M. Kleinberg. Authoritative sources in a hyperlinked environment. *J. ACM*, 46(5):604–632, Sept. 1999.
- [20] E. M. Knorr, R. T. Ng, and V. Tucakov. Distance-based outliers: Algorithms and applications. *VLDB J.*, 8(3-4):237–253, 2000.
- [21] H.-P. Kriegel, P. Kröger, E. Schubert, and A. Zimek. Outlier detection in axis-parallel subspaces of high dimensional data. In *Advances in Knowledge Discovery and Data Mining: 13th Pacific-Asia Conference, PAKDD 2009 Bangkok, Thailand, April 27-30, 2009 Proceedings 13*, pages 831–838. Springer, 2009.
- [22] H.-P. Kriegel, P. Kröger, E. Schubert, and A. Zimek. Interpreting and unifying outlier scores. In *SDM*, pages 13–24. SIAM / Omnipress, 2011.
- [23] H.-P. Kriegel, M. Schubert, and A. Zimek. Angle-based outlier detection in high-dimensional data. In *SIGKDD*, pages 444–452, 2008.
- [24] S. Kumar, B. Hooi, D. Makhija, M. Kumar, C. Faloutsos, and V. S. Subrahmanian. Rev2: Fraudulent user prediction in rating platforms. In *WSDM*, pages 333–341. ACM, 2018.
- [25] H. Lamba and L. Akoglu. Learning on-the-job to re-rank anomalies from top-1 feedback. In *SDM*, pages 612–620. SIAM, 2019.
- [26] C.-L. Li, K. Sohn, J. Yoon, and T. Pfister. Cutpaste: Self-supervised learning for anomaly detection and localization. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 9664–9674, 2021.
- [27] Y. Li, Z. Chen, D. Zha, K. Zhou, H. Jin, H. Chen, and X. Hu. Autood: automated outlier detection via curiosity-guided search and self-imitation learning. *arXiv preprint arXiv:2006.11321*, 2020.
- [28] Y. Li, D. Zha, P. Venugopal, N. Zou, and X. Hu. Pyodds: An end-to-end outlier detection system with automated machine learning. In *Companion Proceedings of the Web Conference 2020*, pages 153–157, 2020.
- [29] J. Liang and S. Parthasarathy. Robust contextual outlier detection: Where context meets sparsity. In *CIKM*, pages 2167–2172. ACM, 2016.
- [30] Z. Lin, K. Thekumparampil, G. Fanti, and S. Oh. InfoGAN-CR and ModelCentrality: Self-supervised model training and selection for disentangling GANs. In *International Conference on Machine Learning*, pages 6127–6139. PMLR, 2020.
- [31] C. Lioma, J. G. Simonsen, and B. Larsen. Evaluation measures for relevance and credibility in ranked lists. In *ICTIR*, pages 91–98. ACM, 2017.
- [32] F. T. Liu, K. M. Ting, and Z.-H. Zhou. Isolation forest. In *ICDM*, pages 413–422. IEEE, 2008.
- [33] M. Macha and L. Akoglu. Explaining anomalies in groups with characterizing subspace rules. *Data Mining and Knowledge Discovery*, 32:1444–1480, 2018.
- [34] E. A. Manzoor, H. Lamba, and L. Akoglu. xstream: Outlier detection in feature-evolving data streams. In *KDD*, pages 1963–1972. ACM, 2018.
- [35] H. O. Marques, R. J. G. B. Campello, J. Sander, and A. Zimek. Internal evaluation of unsupervised outlier detection. *ACM Trans. Knowl. Discov. Data*, 14(4):47:1–47:42, 2020.
- [36] H. O. Marques, R. J. G. B. Campello, A. Zimek, and J. Sander. On the internal evaluation of unsupervised outlier detection. In *SSDBM*, pages 7:1–7:12. ACM, 2015.
- [37] M. M. Meghanath, D. Pai, and L. Akoglu. Conout: Contextual outlier detection with multiple contexts: Application to ad fraud. In *ECML/PKDD (1)*, volume 11051, pages 139–156. Springer, 2018.
- [38] A. Mishra, J. S. Pudipeddi, and L. Akoglu. Ranking in heterogeneous networks with geo-location information. In *SDM*, pages 408–416. SIAM, 2017.

- [39] T. T. Nguyen, A. T. Nguyen, T. A. H. Nguyen, L. T. Vu, Q. U. Nguyen, and L. D. Hai. Unsupervised anomaly detection in online game. In *Proceedings of the Sixth International Symposium on Information and Communication Technology*, pages 4–10. ACM, 2015.
- [40] V. Nguyen, T. Nguyen, and U. Nguyen. An evaluation method for unsupervised anomaly detection algorithms. *Journal of Computer Science and Cybernetics*, 32(3):259–272, 2017.
- [41] T. Pevný. Loda: Lightweight on-line detector of anomalies. *Machine Learning*, 102(2):275–304, 2016.
- [42] S. Ramaswamy, R. Rastogi, and K. Shim. Efficient algorithms for mining outliers from large data sets. In *SIGMOD*, pages 427–438, 2000.
- [43] S. Rayana. ODDS library, 2016.
- [44] S. Rayana and L. Akoglu. Less is more: Building selective anomaly ensembles. *ACM Trans. Knowl. Discov. Data*, 10(4):42:1–42:33, 2016.
- [45] P. Rousseeuw. Silhouettes: A graphical aid to the interpretation and validation of cluster analysis. *Journal of Computational and Applied Mathematics*, 20(1):53–65, Nov. 1987.
- [46] L. Ruff, J. R. Kauffmann, R. A. Vandermeulen, G. Montavon, W. Samek, M. Kloft, T. G. Dietterich, and K.-R. Müller. A unifying review of deep and shallow anomaly detection. *CoRR*, abs/2009.11732, 2020.
- [47] B. Schölkopf, J. C. Platt, J. Shawe-Taylor, A. J. Smola, and R. C. Williamson. Estimating the support of a high-dimensional distribution. *Neural computation*, 13(7):1443–1471, 2001.
- [48] S. C. Tan, K. M. Ting, and F. T. Liu. Fast anomaly detection for streaming data. In *IJCAI*, pages 1511–1516. IJCAI/AAAI, 2011.
- [49] J. Tang, Z. Chen, A. W.-C. Fu, and D. W. Cheung. Enhancing effectiveness of outlier detections for low density patterns. In *PAKDD*, pages 535–548. Springer, 2002.
- [50] G. Wang, S. Xie, B. Liu, and P. S. Yu. Review graph based online store review spammer detection. In *ICDM*, pages 1242–1247. IEEE Computer Society, 2011.
- [51] R. Wang, K. Nie, T. Wang, Y. Yang, and B. Long. Deep learning for anomaly detection. In *WSDM*, pages 894–896. ACM, 2020.
- [52] X. L. Xie and G. Beni. A validity measure for fuzzy clustering. *IEEE Trans. Pattern Anal. Mach. Intell.*, 13(8):841–847, 1991.
- [53] X. Yin, J. Han, and P. S. Yu. Truth discovery with multiple conflicting information providers on the web. In *KDD*, pages 1048–1052. ACM, 2007.
- [54] Y. Zhao and L. Akoglu. Toward unsupervised outlier model selection. *IEEE ICDM*, 2022.
- [55] Y. Zhao and L. Akoglu. Towards unsupervised hyperparameter optimization for outlier detection. *arXiv preprint arXiv:2208.11727*, 2022.
- [56] Y. Zhao, Z. Nasrullah, and Z. Li. Pyod: A python toolbox for scalable outlier detection. *JMLR*, 20(96):1–7, 2019.
- [57] Y. Zhao, R. Rossi, and L. Akoglu. Automatic unsupervised outlier model selection. *Advances in Neural Information Processing Systems*, 34:4489–4502, 2021.
- [58] A. Zimek, R. J. G. B. Campello, and J. Sander. Ensembles for unsupervised outlier detection: challenges and research questions a position paper. *SIGKDD Explor.*, 15(1):11–22, 2013.

APPENDIX

A. FAMILY-WISE MODEL PERFORMANCES

In this study we use **8 different families of outlier detection algorithms**, namely; LODA, ABOD, iForest, kNN, LOF, HBOS, OCSVM, and COF. We build a total of **297 detection models** based on various hyperparameter (HP) configurations of these algorithms, as listed in Table 1.

Tables 16, 17, and 18 (resp. for AP, ROC AUC, and Prec@ k) show the family-wise average performance of each detection algorithm (averaged over within-family models with different HP settings) on each dataset, as well as mean and standard deviation across datasets.

These show **iForest to be the most competitive detector**, which we compare to as a baseline to study *whether unsupervised model selection outperforms always using the same (state-of-the-art) detector*.

B. MODEL PERFORMANCES ON INDIVIDUAL DATASETS

Figures 4, 5, and 6 (resp. for AP, ROC AUC, and Prec@ k) show the distribution of performances across all 297 models via boxplots for each dataset. For most datasets, there exists **considerable difference between the best and the worst performing model**—suggesting that effective model selection would be beneficial.

C. CORRESPONDING RESULTS BASED ON OTHER METRICS

For brevity, we reported all performance results in Evaluation (Sec. 4) based on Average Precision (AP). For completeness, we provide the results of the same analysis corresponding to **ROC AUC** and **Prec@ k** metrics.

The conclusions are similar for these two metrics.

Cluster quality based methods. Specifically, Tables 8 and 12 present, resp. for ROC and Prec@ k , the pairwise comparison of cluster quality based methods and the baselines (RANDOM and iFOREST-R). Three strategies RS, CH, and XB appear to stand out from others. However, none of the methods are not significantly different from (and few are sometimes worse than) RANDOM. Most of them are significantly worse than iFOREST-R, with otherwise a very large p -value.

Other stand-alone methods. Tables 9 and 13 present, resp. for ROC and Prec@ k , the pairwise comparison of all the stand-alone methods (only RS, CH, and XB from above) and the baselines. We find that they are not different from each other or RANDOM—implying that **stand-alone model selection techniques would not be useful in practice**.

Consensus-based methods. Tables 10 and 14 show, resp. for ROC and Prec@ k , that all consensus-based techniques, namely UDR, MC, MC_S , HITS, and ENS, are comparable to each other in terms of selection performance.

Finally, Tables 11 and 15 provide, resp. for ROC and Prec@ k , a summary of the results for all the unsupervised model selection methods we studied. Main take-aways are: (1) **Consensus-based model selection methods are more competitive** than stand-alone methods, where all of them achieve **significantly better performance than RANDOM**

selection. (2) Further, they are most often not different from iFOREST-R (a state-of-the-art detector) and sometimes even better (w.r.t. ROC). However, their absolute difference (i.e. effect size) is negligible as shown in Figure 3 for both ROC and Prec@ k . Notably, their performance differences are not far from zero, suggesting that **consensus-based selection would also not be preferable in practice**, since training a *single* iFOREST-R model is much faster over training a *pool* of models (with considerable running time overhead) to select from.

Table 8: Comparison of cluster quality based methods and baselines by one-sided paired Wilcoxon signed rank test on ROC AUC. p -values **bolded** (underlined) highlight the cases where row-method is significantly **better** (worse) than col-method at $p \leq 0.05$.

	STD	H	S	I	DB	SD	D	RND	iF
XB,RS,CH	0.001	0.407	0.007	0.389	0.272	0.099	0.518	0.358	0.980
STD		1.000	0.990	1.000	0.994	0.995	1.000	1.000	1.000
H			0.021	0.500	0.487	0.320	0.831	0.818	1.000
S				0.974	0.816	0.704	0.994	0.994	1.000
I					0.487	0.323	0.849	0.821	1.000
DB						0.368	0.815	0.662	0.996
SD							0.842	0.905	0.999
D								0.110	0.998
RND									1.000

Table 9: Comparison of stand-alone methods and baselines w.r.t. ROC AUC.

	EM	MV	IREOS	RND	iF
XB,RS,CH	0.364	0.422	0.934	0.358	0.980
EM		0.079	0.969	0.358	0.992
MV			0.977	0.369	0.997
IREOS				0.006	0.702

Table 10: Comparison of consensus-based methods (UDR, MC, MC_S are based on $NDCG$) w.r.t. ROC AUC.

	MC	MC_S	HITS	ENS
UDR	0.462	0.070	0.408	0.232
MC		0.100	0.134	0.069
MC_S			0.681	0.511
HITS				0.740

Table 11: **Summary of results:** p -values by one-sided paired Wilcoxon signed rank test comparing UOMS methods to the baselines, smallest q -th best model with no significant difference, and mean/standard deviation ROC AUC across datasets.

	Method	RANDOM	iFOREST-R	q_{ROC}	mean ROC	std ROC
S-alone	XB,RS,CH	0.358	0.980	138	0.690	0.206
	EM	0.358	0.992	142	0.682	0.216
	IREOS	0.006	0.702	83	0.730	0.203
Consensus-based	UDR- ρ	0.000	0.279	82	0.763	0.180
	UDR- τ	0.000	0.186	75	0.769	0.180
	UDR- $NDCG$	0.000	0.175	75	0.769	0.183
	MC- ρ	0.000	0.036	92	0.767	0.168
	MC- τ	0.000	0.011	91	0.769	0.167
	MC- $NDCG$	0.000	0.034	86	0.771	0.170
	MC_S - ρ	0.000	0.483	100	0.763	0.173
	MC_S - τ	0.000	0.121	94	0.761	0.167
	MC_S - $NDCG$	0.000	0.274	94	0.766	0.165
Agg.	HITS	0.000	0.148	97	0.762	0.169
	ENS	0.000	0.230	86	0.749	0.183
	HITS-AUTH	0.000	0.018	77	0.785	0.163
Base.	ENS-PSEUDO	0.000	0.135	87	0.749	0.184
	RANDOM	–	1.000	183	0.704	0.133
	iFOREST-R	–	–	102	0.763	0.166

Table 12: Comparison of cluster quality based methods and baselines by one-sided paired Wilcoxon signed rank test on Prec@ k . p -values **bolded** (underlined) highlight the cases where row-method is significantly **better** (worse) than col-method at $p \leq 0.05$.

	STD	H	S	I	DB	SD	D	RND	iF
XB,RS,CH	0.000	0.125	0.031	0.109	0.173	0.026	0.274	0.090	0.716
STD		<u>0.998</u>	<u>0.985</u>	<u>0.999</u>	<u>0.989</u>	<u>0.956</u>	<u>1.000</u>	<u>1.000</u>	<u>1.000</u>
H			0.447	0.704	0.623	0.191	0.581	0.500	<u>0.967</u>
S				0.488	0.815	0.313	0.875	0.757	<u>0.961</u>
I					0.631	0.203	0.632	0.544	<u>0.978</u>
DB						0.166	0.719	0.423	0.915
SD							0.929	0.879	<u>0.993</u>
D								0.201	0.923
RND									<u>0.999</u>

Table 13: Comparison of stand-alone methods and baselines w.r.t. Prec@ k .

	EM	MV	IREOS	RND	iF
XB,RS,CH	0.272	0.193	0.405	0.090	0.716
EM		0.187	0.696	0.730	<u>0.967</u>
MV			0.770	0.829	<u>0.987</u>
IREOS				0.423	0.944

Table 14: Comparison of consensus-based methods (UDR, MC, MC_S are based on $NDCG$) w.r.t. Prec@ k .

	MC	MC_S	HITS	ENS
UDR	0.645	0.403	0.464	0.296
MC		0.145	0.227	0.341
MC_S			0.375	0.488
HITS				0.608

Table 15: **Summary of results:** p -values by one-sided paired Wilcoxon signed rank test comparing UOMS methods to the baselines, smallest q -th best model with no significant difference, and mean/standard deviation Prec@ k across datasets.

	Method	RANDOM	iFOREST-R	q_{Prec}	mean Prec@ k	std Prec@ k
S-alone	XB,RS,CH	0.090	0.716	91	0.348	0.277
	EM	0.730	<u>0.967</u>	119	0.303	0.254
	IREOS	0.423	0.944	102	0.316	0.255
Consensus-based	UDR- ρ	0.039	<u>0.965</u>	115	0.354	0.271
	UDR- τ	0.025	0.942	110	0.356	0.263
	UDR- $NDCG$	0.002	0.600	86	0.372	0.255
	MC- ρ	0.002	0.555	98	0.369	0.271
	MC- τ	0.002	0.833	103	0.370	0.280
	MC- $NDCG$	0.000	0.228	89	0.378	0.270
	MC_S - ρ	0.008	0.937	115	0.361	0.276
	MC_S - τ	0.002	0.595	96	0.374	0.290
	MC_S - $NDCG$	0.002	0.210	92	0.367	0.274
	HITS	0.001	0.583	99	0.376	0.280
ENS	0.004	0.595	92	0.351	0.261	
Agg.	HITS-AUTH	0.000	0.293	89	0.380	0.263
	ENS-PSEUDO	0.005	0.722	89	0.350	0.262
Base.	RANDOM	-	<u>0.999</u>	153	0.325	0.217
	iFOREST-R	-	-	91	0.374	0.280

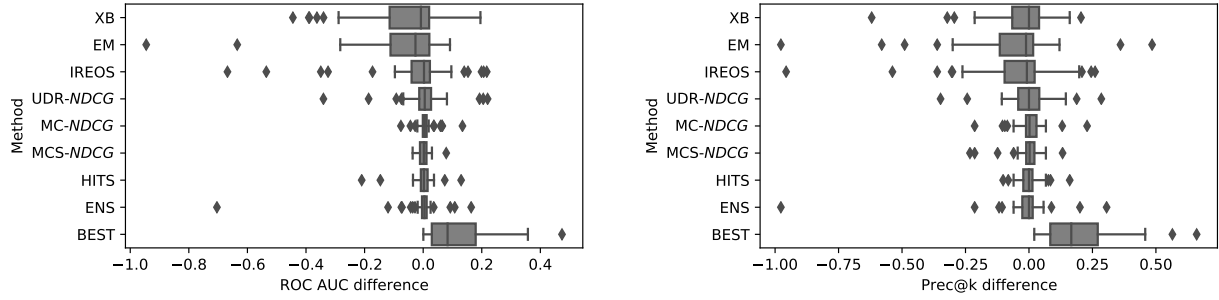


Figure 3: Distribution of performance difference across datasets: (left) ROC AUC and (right) Prec@ k of selected model (by each UOMS method studied) minus that of iFOREST-R. Stand-alone methods and UDR are subpar, whereas consensus-based methods’ differences concentrate around zero (not notably different from iFOREST-R). Also shown for comparison is BEST model on each dataset, showcasing ample room for improvement over iFOREST-R.

Table 16: Family-wise model performance in AP. Values in **bold** highlight the model that outperforms for each dataset (per row). iForest achieves the highest average performance across all datasets.

Dataset	LODA	ABOD	iForest	kNN	LOF	HBOS	OCSVM	COF
annthyroid (ODDS)	0.136	0.232	0.340	0.228	0.172	0.388	0.145	0.138
arrhythmia (ODDS)	0.387	0.315	0.470	0.392	0.362	0.431	0.250	0.404
breastw (ODDS)	0.964	0.702	0.972	0.942	0.331	0.959	0.544	0.304
glass (ODDS)	0.063	0.137	0.104	0.106	0.117	0.061	0.063	0.154
ionosphere (ODDS)	0.766	0.921	0.784	0.868	0.819	0.288	0.492	0.852
letter (ODDS)	0.092	0.319	0.089	0.258	0.359	0.080	0.138	0.459
lympho (ODDS)	0.447	0.555	0.957	0.763	0.668	0.905	0.418	0.464
mammography (ODDS)	0.218	0.147	0.234	0.169	0.102	0.096	0.156	0.064
mnist (ODDS)	0.203	0.329	0.261	0.401	0.273	0.097	0.204	0.195
musk (ODDS)	0.904	0.038	0.990	0.588	0.130	0.997	0.498	0.174
optdigits (ODDS)	0.025	0.057	0.049	0.021	0.037	0.177	0.031	0.048
pendigits (ODDS)	0.245	0.057	0.280	0.104	0.038	0.231	0.086	0.037
pima (ODDS)	0.445	0.508	0.492	0.524	0.441	0.521	0.385	0.429
satellite (ODDS)	0.630	0.430	0.664	0.562	0.375	0.711	0.456	0.368
satimage-2 (ODDS)	0.904	0.212	0.916	0.615	0.055	0.717	0.486	0.078
speech (ODDS)	0.018	0.093	0.020	0.024	0.031	0.025	0.022	0.034
thyroid (ODDS)	0.238	0.218	0.587	0.354	0.157	0.630	0.196	0.032
vertebral (ODDS)	0.089	0.098	0.094	0.090	0.101	0.087	0.131	0.116
vowels (ODDS)	0.140	0.690	0.134	0.487	0.348	0.083	0.080	0.408
wbc (ODDS)	0.603	0.367	0.599	0.533	0.497	0.673	0.321	0.261
wine (ODDS)	0.286	0.082	0.215	0.253	0.253	0.402	0.249	0.081
Annthyroid (DAMI)	0.097	0.137	0.160	0.126	0.134	0.145	0.079	0.130
Arrhythmia (DAMI)	0.685	0.668	0.757	0.711	0.702	0.745	0.523	0.712
Cardiotocography (DAMI)	0.433	0.254	0.433	0.316	0.280	0.344	0.314	0.267
HeartDisease (DAMI)	0.562	0.547	0.538	0.557	0.509	0.619	0.475	0.486
InternetAds (DAMI)	0.251	0.293	0.490	0.289	0.263	0.521	0.237	0.261
PageBlocks (DAMI)	0.464	0.416	0.449	0.526	0.360	0.201	0.268	0.232
Pima (DAMI)	0.448	0.506	0.494	0.529	0.467	0.487	0.392	0.432
SpamBase (DAMI)	0.370	0.357	0.487	0.406	0.364	0.532	0.366	0.392
Stamps (DAMI)	0.332	0.218	0.336	0.313	0.228	0.315	0.209	0.159
Wilt (DAMI)	0.039	0.065	0.045	0.053	0.075	0.044	0.065	0.101
ALOI (DAMI)	0.034	0.102	0.033	0.057	0.100	0.031	0.035	0.144
Glass (DAMI)	0.085	0.221	0.183	0.146	0.118	0.115	0.107	0.179
PenDigits (DAMI)	0.003	0.031	0.005	0.040	0.014	0.004	0.016	0.017
Shuttle (DAMI)	0.111	0.250	0.071	0.326	0.296	0.094	0.095	0.173
Waveform (DAMI)	0.052	0.055	0.057	0.115	0.095	0.053	0.069	0.102
WBC (DAMI)	0.743	0.595	0.858	0.671	0.359	0.683	0.424	0.146
WDBC (DAMI)	0.720	0.296	0.669	0.571	0.554	0.725	0.322	0.295
WPBC (DAMI)	0.235	0.231	0.229	0.233	0.230	0.239	0.237	0.219
average	0.345	0.301	0.399	0.366	0.277	0.371	0.246	0.245
STD	0.282	0.220	0.304	0.248	0.199	0.295	0.165	0.188

Table 17: Family-wise model performance in ROC AUC. Values in **bold** highlight the model that outperforms for each dataset (per row). kNN (0.764) and iForest (0.763) achieve the highest average performance across all datasets.

Dataset	LODA	ABOD	iForest	kNN	LOF	HBOS	OCSVM	COF
annthyroid (ODDS)	0.572	0.823	0.841	0.775	0.729	0.736	0.517	0.689
arrhythmia (ODDS)	0.735	0.751	0.803	0.777	0.764	0.806	0.522	0.757
breastw (ODDS)	0.980	0.898	0.988	0.980	0.500	0.985	0.481	0.459
glass (ODDS)	0.539	0.766	0.707	0.747	0.747	0.638	0.429	0.772
ionosphere (ODDS)	0.814	0.928	0.838	0.898	0.870	0.357	0.548	0.879
letter (ODDS)	0.584	0.880	0.629	0.842	0.846	0.581	0.554	0.880
lympho (ODDS)	0.814	0.936	0.998	0.971	0.938	0.985	0.607	0.834
mammography (ODDS)	0.854	0.822	0.862	0.845	0.729	0.799	0.629	0.700
mnist (ODDS)	0.586	0.797	0.794	0.856	0.708	0.515	0.536	0.615
musk (ODDS)	0.991	0.072	0.999	0.830	0.521	1.000	0.669	0.534
optdigits (ODDS)	0.414	0.477	0.713	0.383	0.463	0.877	0.463	0.526
pendigits (ODDS)	0.934	0.692	0.948	0.818	0.516	0.921	0.548	0.508
pima (ODDS)	0.629	0.685	0.652	0.717	0.630	0.634	0.497	0.583
satellite (ODDS)	0.644	0.594	0.703	0.703	0.546	0.785	0.506	0.519
satimage-2 (ODDS)	0.988	0.854	0.993	0.965	0.678	0.973	0.610	0.537
speech (ODDS)	0.474	0.688	0.473	0.500	0.525	0.473	0.492	0.584
thyroid (ODDS)	0.820	0.945	0.983	0.960	0.771	0.950	0.550	0.581
vertebral (ODDS)	0.315	0.375	0.349	0.333	0.380	0.297	0.482	0.454
vowels (ODDS)	0.712	0.976	0.736	0.944	0.905	0.676	0.529	0.877
wbc (ODDS)	0.941	0.918	0.938	0.935	0.892	0.950	0.603	0.792
wine (ODDS)	0.853	0.490	0.794	0.779	0.758	0.873	0.536	0.373
Annthyroid (DAMI)	0.491	0.717	0.679	0.658	0.679	0.646	0.471	0.666
Arrhythmia (DAMI)	0.687	0.725	0.750	0.736	0.732	0.736	0.506	0.736
Cardiotocography (DAMI)	0.689	0.458	0.689	0.503	0.544	0.566	0.489	0.522
HeartDisease (DAMI)	0.608	0.612	0.602	0.637	0.582	0.670	0.502	0.542
InternetAds (DAMI)	0.548	0.657	0.690	0.626	0.587	0.695	0.499	0.579
PageBlocks (DAMI)	0.785	0.780	0.894	0.889	0.759	0.679	0.558	0.610
Pima (DAMI)	0.624	0.666	0.644	0.706	0.650	0.594	0.504	0.587
SpamBase (DAMI)	0.433	0.403	0.635	0.535	0.441	0.676	0.463	0.450
Stamps (DAMI)	0.891	0.793	0.901	0.872	0.702	0.876	0.582	0.541
Wilt (DAMI)	0.363	0.628	0.457	0.538	0.626	0.419	0.489	0.695
ALOI (DAMI)	0.504	0.739	0.534	0.641	0.744	0.508	0.506	0.796
Glass (DAMI)	0.659	0.854	0.794	0.822	0.748	0.795	0.485	0.774
PenDigits (DAMI)	0.628	0.936	0.768	0.967	0.821	0.734	0.537	0.718
Shuttle (DAMI)	0.637	0.927	0.853	0.963	0.911	0.842	0.566	0.848
Waveform (DAMI)	0.664	0.666	0.707	0.743	0.716	0.703	0.492	0.689
WBC (DAMI)	0.983	0.954	0.991	0.979	0.842	0.985	0.611	0.703
WDBC (DAMI)	0.945	0.890	0.936	0.924	0.871	0.963	0.629	0.800
WPBC (DAMI)	0.509	0.501	0.498	0.509	0.503	0.536	0.485	0.463
average	0.688	0.725	0.763	0.764	0.689	0.729	0.530	0.645
STD	0.188	0.197	0.168	0.175	0.146	0.188	0.054	0.138

Table 18: Family-wise model performance in Prec@ k . Values in **bold** highlight the model that outperforms for each dataset (per row). iForest achieves the highest average performance across all datasets.

Dataset	LODA	ABOD	iForest	kNN	LOF	HBOS	OCSVM	COF
annthyroid (ODDS)	0.180	0.301	0.337	0.297	0.209	0.387	0.180	0.169
arrhythmia (ODDS)	0.403	0.372	0.481	0.411	0.386	0.495	0.237	0.407
breastw (ODDS)	0.924	0.788	0.929	0.923	0.271	0.938	0.445	0.152
glass (ODDS)	0.019	0.111	0.111	0.111	0.136	0.014	0.040	0.143
ionosphere (ODDS)	0.645	0.849	0.648	0.753	0.725	0.228	0.439	0.764
letter (ODDS)	0.100	0.354	0.092	0.312	0.358	0.080	0.140	0.440
lympho (ODDS)	0.401	0.476	0.881	0.639	0.560	0.808	0.347	0.405
mammography (ODDS)	0.286	0.197	0.261	0.251	0.194	0.114	0.192	0.114
mnist (ODDS)	0.212	0.376	0.293	0.420	0.315	0.095	0.218	0.246
musk (ODDS)	0.873	0.035	0.977	0.546	0.134	0.981	0.491	0.218
optdigits (ODDS)	0.001	0.045	0.025	0.000	0.029	0.211	0.018	0.067
pendigits (ODDS)	0.324	0.077	0.365	0.110	0.072	0.269	0.113	0.063
pima (ODDS)	0.466	0.530	0.504	0.551	0.463	0.476	0.361	0.423
satellite (ODDS)	0.533	0.417	0.573	0.511	0.379	0.619	0.382	0.361
satimage-2 (ODDS)	0.865	0.260	0.862	0.577	0.086	0.661	0.465	0.145
speech (ODDS)	0.019	0.138	0.031	0.039	0.045	0.032	0.039	0.049
thyroid (ODDS)	0.287	0.198	0.620	0.332	0.149	0.645	0.224	0.000
vertebral (ODDS)	0.011	0.043	0.044	0.018	0.056	0.012	0.074	0.090
vowels (ODDS)	0.194	0.641	0.175	0.474	0.333	0.121	0.094	0.429
wbc (ODDS)	0.558	0.361	0.536	0.496	0.475	0.614	0.324	0.293
wine (ODDS)	0.257	0.000	0.140	0.194	0.203	0.408	0.200	0.043
Annthyroid (DAMI)	0.116	0.153	0.213	0.134	0.165	0.191	0.074	0.162
Arrhythmia (DAMI)	0.604	0.630	0.655	0.637	0.643	0.632	0.459	0.652
Cardiotocography (DAMI)	0.407	0.266	0.396	0.311	0.288	0.303	0.259	0.264
HeartDisease (DAMI)	0.530	0.520	0.503	0.535	0.506	0.591	0.447	0.470
InternetAds (DAMI)	0.267	0.344	0.449	0.334	0.304	0.466	0.244	0.284
PageBlocks (DAMI)	0.458	0.425	0.397	0.506	0.376	0.158	0.264	0.268
Pima (DAMI)	0.476	0.512	0.499	0.547	0.485	0.448	0.369	0.421
SpamBase (DAMI)	0.351	0.359	0.518	0.421	0.338	0.562	0.357	0.382
Stamps (DAMI)	0.275	0.189	0.286	0.211	0.169	0.385	0.197	0.180
Wilt (DAMI)	0.001	0.012	0.012	0.003	0.058	0.006	0.043	0.121
ALOI (DAMI)	0.050	0.144	0.028	0.086	0.146	0.028	0.043	0.187
Glass (DAMI)	0.027	0.143	0.111	0.111	0.133	0.044	0.056	0.159
PenDigits (DAMI)	0.000	0.036	0.000	0.000	0.019	0.000	0.010	0.036
Shuttle (DAMI)	0.120	0.319	0.079	0.277	0.169	0.092	0.092	0.231
Waveform (DAMI)	0.057	0.069	0.065	0.191	0.161	0.063	0.083	0.143
WBC (DAMI)	0.630	0.429	0.723	0.644	0.328	0.713	0.356	0.086
WDBC (DAMI)	0.650	0.271	0.633	0.592	0.536	0.648	0.350	0.286
WPBC (DAMI)	0.166	0.164	0.146	0.160	0.172	0.206	0.202	0.161
average	0.327	0.296	0.374	0.350	0.271	0.352	0.229	0.244
STD	0.262	0.214	0.284	0.235	0.180	0.283	0.149	0.171

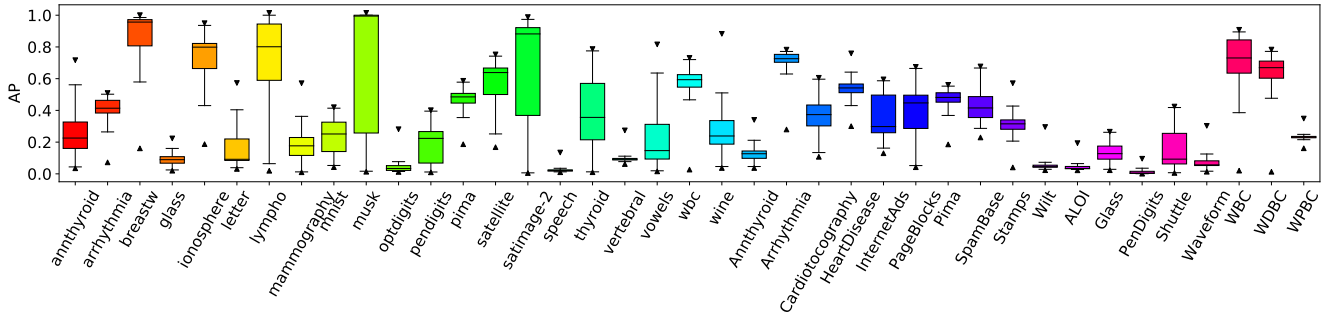


Figure 4: Model performance boxplot (AP) for all datasets, where triangles mark the min and max. Model performance varies significantly for most datasets, showing the importance of model selection.

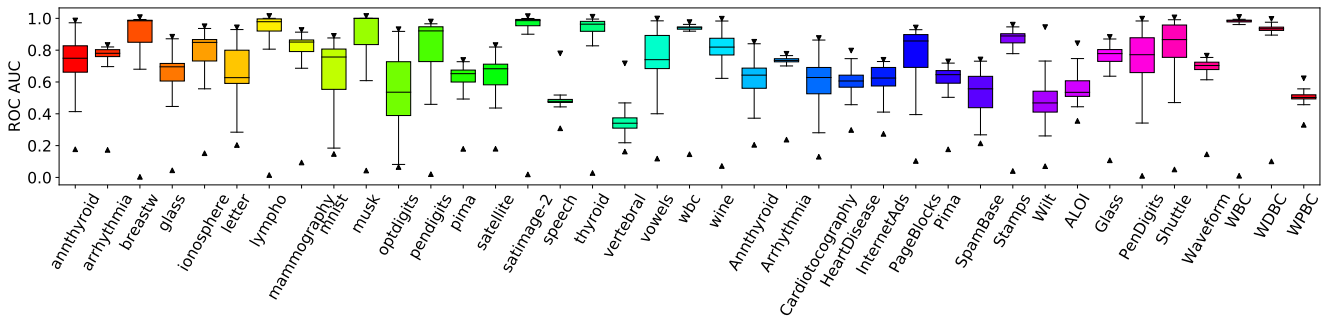


Figure 5: Model performance boxplot (ROC AUC) for all datasets, where triangles mark the min and max. Model performance varies significantly for most datasets, showing the importance of model selection.

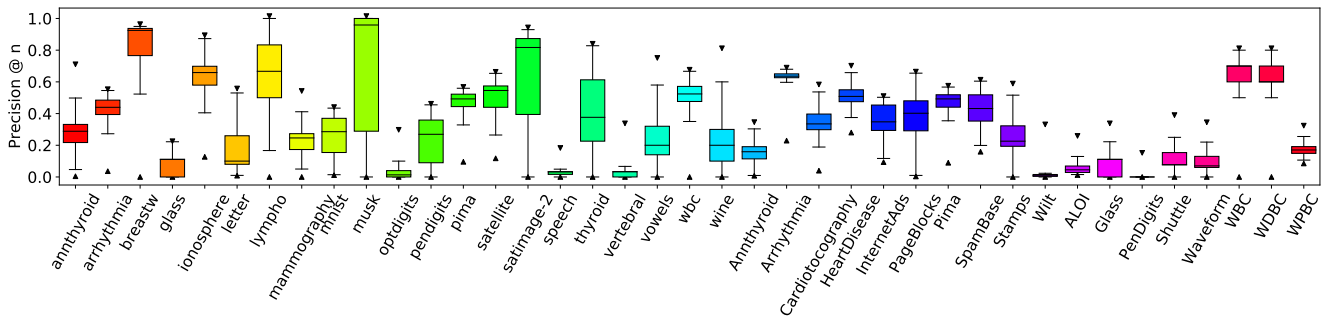


Figure 6: Model performance boxplot (Prec@k) for all datasets, where triangles mark the min and max. Model performance varies significantly for most datasets, showing the importance of model selection.

Stop using the elbow criterion for k-means

and how to choose the number of clusters instead

Erich Schubert
TU Dortmund University
44221 Dortmund, Germany
erich.schubert@tu-dortmund.de

ABSTRACT

A major challenge when using k-means clustering often is how to choose the parameter k , the number of clusters. In this letter, we want to point out that it is very easy to draw poor conclusions from a common heuristic, the “elbow method”. Better alternatives have been known in literature for a long time, and we want to draw attention to some of these easy to use options, that often perform better. This letter is a call to stop using the elbow method altogether, because it severely lacks theoretic support, and we want to encourage educators to discuss the problems of the method – if introducing it in class at all – and teach alternatives instead, while researchers and reviewers should reject conclusions drawn from the elbow method.

1. INTRODUCTION

Cluster analysis aims at identifying subgroups in the data that have high similarity within the group, while they also differ from the remainder of the data set. No single “best” definition of a cluster exists. Bonner [5] noted that “none of the many specific definitions [of clusters] seems ‘best’ in any general sense”, and Estivill-Castro [12] argued that it cannot exist. Each data set and use case may call for different properties to be desirable, which in turn leads to different algorithms to find the “best” solution. Hence, a large number of clustering methods were developed over the last decades, based on concepts such as finding a hierarchical structure (akin to phylogenetic trees), quantization and compression, parametric modeling, or identifying dense areas.

Despite the many different concepts of clusters and the wide variety of clustering algorithms available, one method currently is the most used and most taught clustering method: k-means clustering. One of the main reasons may be the simplicity of the standard algorithm: assigning each point to the nearest center, then recomputing all the cluster centers until nothing changes – this algorithm can be easily described in a single sentence. At the same time, this algorithm runs very fast, and it will always produce a result with exactly k clusters, giving a (false) suggestion of success. A key problem then with applying this method to data is often the need to choose the number of clusters k , although users should first consider whether k-means is even the right choice for their problem at all, and pay more attention to data preprocessing, too. It makes no sense to search for the “optimum” k if k-means is not solving the problem.

2. K-MEANS CLUSTERING

Formally, k -means clustering is a least-squares optimization problem. We can best view it as a data quantization technique, where we want to approximate the data set of N objects in a continuous, d -dimensional vector space \mathbb{R}^d using k centers. The quantization error for a data set X and a set C of centers then is called inertia, the within-cluster sum of squares (WCSS), or the sum of squared errors (SSE):

$$\text{SSE}(X, C) = \sum_{x \in X} \min_{c \in C} \|x - c\|^2. \quad (1)$$

While it is easy to optimize this for a single cluster center by taking the arithmetic average in each dimension, i.e., the data set centroid, the problem is NP-hard for multiple clusters and higher dimensionality [22; 1].

Several methods to optimize this objective exist, but because of the hardness, most heuristics will only find a local fixpoint. Because of the very common least-squares objective, the standard algorithm has likely been invented several times independently, as discussed in the overview of Bock [4]. The standard heuristic for k -means is an alternating optimization, which first assigns each point to the nearest current cluster center, then updates each cluster center position with the centroid of the points assigned to it. If we keep assignments unchanged whenever distances are identical, the algorithm will eventually not find any changes and stop (because both steps may never worsen the objective function, and there exists only a finite number of possible cluster assignments). The standard algorithm has a complexity of $O(Nkdi)$, where i is the number of iterations (which theoretically could be very high, but usually is small in practice). This makes it one of the fastest clustering methods we have available, compared to $O(N^3 + N^2d)$ for the standard algorithm for hierarchical clustering, or $O(N^2d)$ for DBSCAN without index acceleration. Many improvements have been proposed that avoid repeated computations in the standard algorithm, nevertheless, this very basic form has become quite popular again with the rise of parallel processing and GPUs: it is embarrassingly parallel, and hence very easy to implement both in clusters as well as GPUs. For this letter, it does not matter which variant of the algorithm we use.

3. THE ELBOW CRITERION

When the number of clusters k is not already given by the application, we have to choose this value; and it turns out this can be rather tricky. The elbow plot is a chart plotting the approximation error SSE on the y -axis over a range of values for k on the x -axis. The motivation of the elbow criterion is the concept of diminishing returns: as we increase

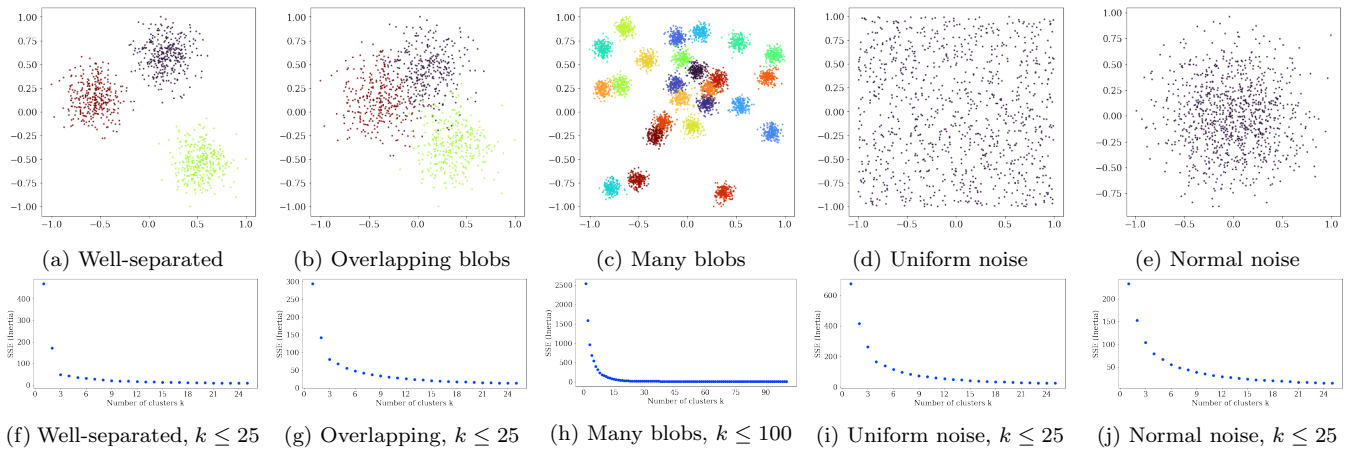


Figure 1: Toy data sets and resulting – very similar – elbow plots.

the number of clusters, the approximation error decreases.¹ In a data set with very well-separated clusters, we expect to initially see a sharp drop until some “optimum” number of clusters, afterwards we are splitting “true” clusters, which leads to much smaller gains. And indeed, on certain toy data sets, this appears to work well. In Figure 1a we have a data set with three well-separated clusters, that is easily clustered by k -means. Figure 1f is the corresponding elbow plot, with a clear inflection at the desired $k = 3$. But the other examples in Figure 1 show that the plot always looks similar, even on uniform data or when the data contains a single normal distribution.

The elbow method is attributed to Thorndike [39], albeit he notes “The curves do not provide much support for the intuitive specification of the number of clusters”, and concludes with doubt: “At this point I can sense the bubbling up of doubts and questions: ‘But what about your *units*?’”.

There are several problems associated with the elbow plot, that statisticians know too well from the scree plot. Because the axes of the plot have very different meanings, we cannot compare them well. We do not have a meaningful measurement of angle, and changing the scaling of the axes (and, e.g., the parameter range of k) may well change the human interpretation of an “elbow”.

3.1 Elbow Detection

Several attempts to formalize the notion of an “elbow” can be found in software and literature. We present only an excerpt in the following, largely to illustrate how *heuristic* and *visual* the machine learning community currently approaches this problem, instead of improving theory.

Sugar et al. [38] propose a “jump method”, finding the maximum of $SSE_k^{-Y} - SSE_{k-1}^{-Y}$ where Y is a power parameter suggested to be half the dimensionality.

Salvador et al. [30] propose the L-method, which fits linear functions to the points before and after the break; choosing the breaking point where this piecewise-linear approximation fits best. But as discussed above, we will often see an exponential curve, and such a linear approximation often

does not fit this curve at all. To improve this, the authors also suggest an iterative approach where they truncate the plot to the first $2 \cdot k$ values if k is the best solution found. Satopää et al. [31] in their Kneedle algorithm want to measure the curvature. For this they fit a smoothing spline to the data, normalize it to 0 to 1, and compute the difference to the diagonal. The last maximum before a parameterizable stopping threshold is chosen.

Zhang et al. [42] note that the standard curvature definition is not independent of rescaling the data, and propose to choose the maximum of a modified curvature:

$$\text{Curvature}_k := \frac{\text{SSE}_{k-1} - \text{SSE}_k}{\text{SSE}_k - \text{SSE}_{k+1}} - 1$$

The pylustering library [25] defines an elbow length:

$$\text{ElbowLen}_k := \frac{(y_0 - y_1)x_k + (x_1 - x_0)y_k + (x_0y_1 - x_1y_0)}{\sqrt{(x_1 - x_0)^2 + (y_1 - y_0)^2}}$$

where x_0, x_1, y_0, y_1 denote the minima and maxima of the graph; intended to measure the length when approximating the curve with the elbow point. There is no literature given for this approach, the given reference to Thorndike [39] does not entail this equation, which again appears to depend very much on the scaling of the plot.

Shi et al. [37] note that “experienced analysts cannot clearly identify the elbow point”, and suggest applying a min-max scaling to the range 0 to 10 instead, then computing angles between triples of adjacent values. The performance of this approach depends much on this weighting factor.

Onumanyi et al. [26] propose AutoElbow, for which they min-max scale the elbow plot to 0 to 1, then propose a geometrically-motivated measure of how close a point is to the bottom left corner and bottom line. The result of this objective changes substantially when increasing the candidate range of k , and hence likely should not be used at all.

$$\text{AutoElbow}_k := \frac{(x_k - 1)^2 + (y_k - 1)^2}{x_k^2 + 2 \cdot y_k^2}$$

¹Given a k -means solution for some k , we can trivially construct a solution with $k + 1$ that is better (unless the error already is zero) by simply adding any point with non-zero error as an additional center.

Table 1: “Optimum” k chosen by different heuristics on the toy data sets. † indicates the result can still be recognized as a poor result by the score value or by visual inspection. ‡ indicates results that fluctuate with random seeds.

		well-sep.		overlapping		many blobs		uniform		normal	
true k		3		3		25		1		1	
max k		10	25	10	25	50	100	10	25	10	25
Elbow-based											
Jump	[38]	3	3	3	3	23	23	4	4	6	21
L-Method	[30]	3	3	3	4	7	9	4	5	4	5
L-Method (iter.)		-	3	-	4	-	6	4	4	4	5
Kneedle	[31]	3	3	3	5	8	10	4	5	4	6
Curvature	[42]	3	3	3	3	38	38	4	4	3	21
Pyclustering	[25]	3	3	3	5	8	10	4	5	4	6
Shi angles	[37]	3	3	3	3	3	3	4	4	4	3
AutoElbow	[26]	3	3	3	6	9	11	4	6	4	7
Variance-based											
Marriot	[23]	3	3	3	3	25	25	9 [†]	17 [†]	2 [†]	2 [†]
VRC	[6]	3	3	3	3	25	25	9 [†]	20 [†]	4 [†]	4 [†]
K-L-Index	[20]	7	7	4	10	35	62	5	21	8	8
Pham	[28]	3	3	2	2	8	8	4	4	10	21
Max reduction		3	3	3	3	8	8	4 [†]	12 [†]	1 [†]	1 [†]
Last reduction		3	3	3	3	25	25	9 [†]	20 [†]	1 [†]	1 [†]
Information-theory-based											
BIC	[27]	3	3	3	3	25	100	9	25	1	1
BIC (fixed)	[13]	3	3	3	3	25	25	1 [†]	1 [†]	1 [†]	1 [†]
Distance-based											
Dunn	[10]	3	3	8	17	18	18	7	20	3	24
DB	[7]	3	3	3	3	21	21	4	4	10	22
Silhouette	[29]	3	3	3	3	21	21	4	4	3	3
Simpl. Silhouette		3	3	2	2	21	21	4	4	3	3
Simulation-based											
Gap	[40]	3	3	3	3	30 [‡]	30 [‡]	14 [‡]	21 [‡]	1	1

3.2 Detection performance

Several of these measures are sensitive to the range of k that we analyze, even if the additional values of k perform poorly. They are heuristics based on the geometric idea of an elbow point, but not taking the process causing the measured data into account. In Table 1 we give the results obtained for the toy data sets of Figure 1 using several heuristics proposed in the literature. Because some methods are very sensitive to the range of k included, we tested two limits, one rather conservative, and one that is much larger. We observe that all methods were able to recognize the best solution on the well-separated data set, and even on the more overlapping version, they all worked – for a small enough maximum k . For the data set with many clusters as well as the uniform data set, all the elbow-based methods failed. The methods based on variance – which we will discuss below – worked much better, but when applied naïvely will still cluster the uniform data. Only when using additional thresholds (or visually inspecting the score plot), the uniform data is recognizable as not clustered. For the normal data, our method indicates a single cluster, while the classic variance-ratio criterion also discussed below has a maximum at six clusters. First of all, the quantity measured, the sum of squared devi-

ations, is a *squared* value. It would make much more sense to analyze the square root of this value, and if we also take the number of points into account, the root-mean-squared-deviation (RMSD), which corresponds to a standard deviation of each point to the nearest center. How meaningful are “angles”, “distances”, “elbows”, and “slopes” on a graph that compares k to SSE, two quantities of different scales? If we scale the entire data set by a factor of α , the SSE will change by α^2 , and the “optimum” found by most of the geometric methods changes, while it is clear that it should not. Secondly, increasing the parameter range of k analyzed must not change the decision once the optimum k is included. Normalizing to the observed minimum and maximum values (often even starting with $k = 2$, not $k = 1$) seems inappropriate. In particular, even without running the algorithm, we know that for $k = N$ we will be able to get an approximation error of 0, so we likely should always consider N to be the maximum x coordinate, and 0 to be the minimum y coordinate. A meaningful normalization should preserve 0. Third, we know that even on random data we obtain a descending curve, and hence we should try to remove this expected behavior from our measure. Fourth, the method should be able to choose $k = 1$ for data that does not contain any meaningful clusters.

3.3 Expected behavior of SSE

Instead of proposing heuristic visual approaches to formalize an imaginary “elbow”, we need to first better understand the quantity that we are working with. The sum of squared errors closely resembles the variance of the data set. Because our cluster centers are derived from the data, we should be using a form of sample variance. Simply dividing the SSE by N will be a biased estimate, and we postulate that $\text{SSE}/(N - k)$ is a more suitable estimate in this context. But since usually $k \ll N$, this will not make much of a difference yet. Instead of working with the squared quantity, we then may want to apply the square root instead, i.e., use $\sqrt{\text{SSE}/(N - k)}$ to have the intuition of a *standard deviation from the nearest center*. Still, the plot obtained this way will look similar to what we started with, and because the square root is a monotone function on the outside, it will not affect the ordering of results – it only serves to make the quantity more interpretable, because ideally, the domain expert should judge whether this is sufficiently small.

As a baseline “expected” behavior, we will for simplicity assume the input data to be uniformly distributed in a single dimension, but with the variance of the input data set. The variance of a uniform interval of length b is $\text{Var}([0; b]) = \frac{1}{12}b^2$. If we slice this into k slices of equal length, each of these has $\text{Var}([0; b/k]) = \frac{1}{12}b^2/k^2$, and we obtain for the resulting total variance $k \cdot \text{Var}([0; b/k]) = \frac{1}{k} \text{Var}([0; b])$. Because of this observation, we propose to use the naïve estimate SSE_1/k as normalization factor. But Krzanowski and Lai [20] suggest that $\text{SSE}/k^{\frac{2}{d}}$ may be more appropriate than our naïve estimate. We should further include the $N - k$ factor discussed above. If there is more than one good parameter k (e.g., because there are substructures in the data), we may also want to compare the solution with the best found so far, e.g., using:

$$\widehat{\text{SSE}}_k := \frac{N-k}{k} \min_{j=1 \dots k-1} \frac{j}{N-j} \text{SSE}_j \quad (2)$$

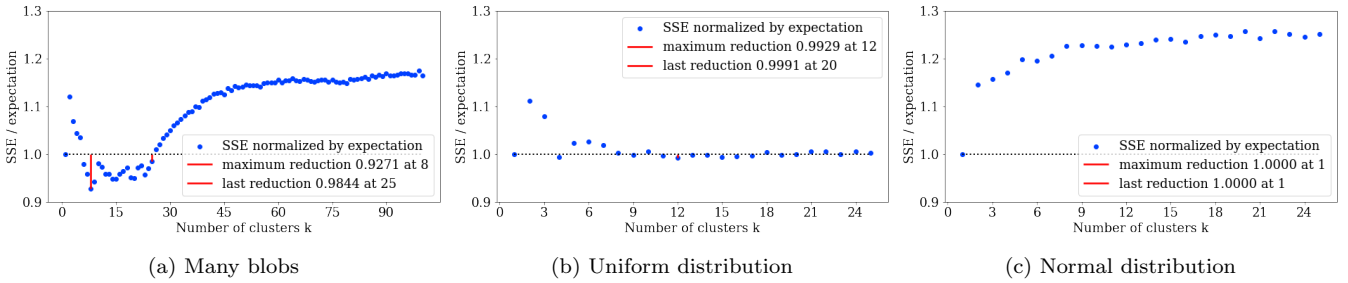


Figure 2: Reduction in \sqrt{SSE} over the estimate $\sqrt{\widehat{SSE}_k}$.

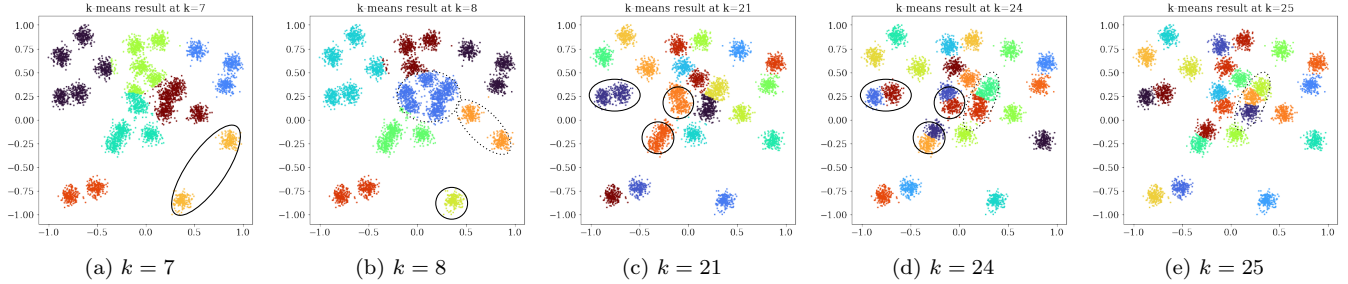


Figure 3: Clustering results on the “many blobs” data set.

We can now generate a standard deviation reduction plot, comparing the observed with the estimated values:

$$\frac{\sqrt{SSE_k / (N - k)}}{\sqrt{\widehat{SSE}_k / (N - k)}} = \sqrt{\frac{SSE_k}{\widehat{SSE}_k}} \quad (3)$$

Note that because \widehat{SSE}_k will tend to 0 as we increase k , eventually this will become unstable for too large k . Depending on our objective, either the smallest value or the last value below 1 (or below a suitable threshold such as 0.99) can be chosen as “best” k . Figure 2 plots this score for some of the above data sets. The normal distribution never scores below 1, and the uniform distribution remains very close to 1; hence these data sets can be recognized as unclustered. For the many blobs data set, the largest reduction is obtained for $k = 8$, and the last reduction is at $k = 25$, the number of generated clusters in this data set. To better understand why both of these solutions are of interest, Figure 3 shows that when going from $k = 7$ to $k = 8$ we observe structural changes that substantially improved the clustering result (such as separating the far cluster in the bottom, but also improving the clustering in the center), whereas the last improvement from $k = 24$ to $k = 25$ only affected three overlapping blobs in the center, that previously were split into two and now correctly into three clusters. For $k = 8$ we have the best structural improvement, but for $k = 25$ we get the finest clustering; both may have their use cases. The solution $k = 21$ is preferred by many distance-based measures, it is worth noting that well-separated blobs are separated, but touching blobs are still joined at this k . Instead of using the Elbow heuristic, most “intrinsic” cluster evaluation criteria can be used. An extensive survey was published by Arbelaitz et al. [2], we only discuss a few examples of particular interest here. This involves some of the best-performing indexes in this study, namely the Silhouette, the VRC, and the DB-Index.

3.4 Variance-based criteria

So have we found a new, better method to choose the number of clusters k ? The above “novel” approach is very similar to the **Variance Ratio Criterion (VRC)** published already in the mid-70s by Calinski and Harabasz [6]:

$$\text{VRC} := \frac{SSE_1 - SSE_k}{k - 1} \bigg/ \frac{SSE_k}{n - k} \quad (4)$$

(using the fact that $\text{BGSS} = \text{TSS} - \text{WGSS} = SSE_1 - SSE_k$). As they noted, this is analogous to the F-statistic used when testing for the significance of a difference in means, but we must not use such a significance test here, because we optimized the means (which makes such a test on the difference in means invalid). There are many more methods discussed in the 70s and 80s literature overlooked by many machine learning scientists of today, that we do not have the space to discuss here, we only briefly highlight some starting points. Marriott [23] (and before, Friedman and Rubin [14]) analyze the determinant of the variance-covariance matrix $|W|$, containing the within-class scatter, also known as generalized variance, instead of the regular variance $\text{tr}(W)$, and discuss that the expected change when partitioning into k clusters is a reduction by $1/k^2$. They argue this approach is superior because it takes correlations into account. Krzanowski and Lai [20] depart from Marriott and return to using the trace again. They argue that the variance is expected to decrease by $k^{\frac{2}{d}}$, define the successive difference as $\text{Diff}_k = (k - 1)^{\frac{2}{d}} SSE_{k-1} - k^{\frac{2}{d}} SSE_k$, and then find the maximum of $KL(k) := |\text{Diff}_k / \text{Diff}_{k+1}|$. When Diff_{k+1} becomes small, this can become unstable, explaining the poor performance in our experiments.

Pham et al. [28] propose a scoring function for $k \geq 2$ based on $SSE_k / (\alpha_k SSE_{k-1})$, where the weights $\alpha_2 = 1 - \frac{3}{4d}$ and $\alpha_k = \frac{5}{6}\alpha_{k-1} + \frac{1}{6}$ model an expected change on a uniform distribution.

3.5 Distance-based criteria

The Dunn [10] index compares the diameter of clusters to the cluster separation. It exists in several variations, in the most basic form it is defined as the ratio of the smallest cluster separation to the largest cluster diameter,

$$\text{Dunn} := \frac{\min_i \min_{j \neq i} \min_{x \in C_i} \min_{y \in C_j} d(x, y)}{\max_i \max_{x \in C_i} \max_{y \in C_i} d(x, y)}.$$

In this basic (original) version, only the smallest cross-cluster distance and the largest inter-cluster distance are taken into account, but we might also consider averages instead [2]. The Davies-Bouldin-Index [7] compares the distance to the nearest other cluster with the radius of the two clusters. This is then averaged over all clusters:

$$\text{DB} := \frac{1}{k} \sum_i \max_{j \neq i} \frac{S_i + S_j}{M_{ij}}$$

where S_i is the (arithmetic, or root-mean-square) average distance of points to their cluster center (and hence a kind of radius), and M_{ij} is the distance between the cluster centers. When using the root-mean-square averages, S_i is the average distance of points within the cluster, and M_{ij} is the average distance between points in different clusters.

One of the most used distance-based criteria is the average silhouette width measure [29], which compares the average distance of each point to its own cluster to the average distance to the nearest other cluster. This method is closely related to k -medoids clustering and the PAM algorithm [17; 34], which cluster the data around k representative objects (called medoids), minimizing the distances to the medoids. In contrast to k -means (which minimizes squared errors), this method can also be used to optimize Euclidean or Manhattan distance; but it is mostly of interest for distances where the mean is not useful.² As Silhouette is fairly expensive to compute, it can be simplified by using the distance from the cluster center or medoid instead of the average distance. But it turns out that we can try to directly optimize this measure using PAM-like algorithms [41; 21].

3.6 Information-theoretic criteria

A different idea to choose the optimum number of clusters is based on the principle of minimum description length. Here, a k -means solution is considered better, if the data can be encoded more compactly. Increasing the number of centers means that the input data is approximated more closely (and hence needs less to encode the deviations), but at the same time, we also need to store more cluster centers. This intuition nicely fits the idea of approximating data and data quantization. X-means [27] integrates this with k -means in an algorithm that dynamically increases the number of clusters as long as a cluster quality criterion improves. They proposed to use the Bayesian Information Criterion (BIC) of Schwarz [36], who also proposed the Akaike Information Criterion (AIC). The original X-means version appears to have an error, the fixed equation of Foglia and Hancock [13] appears to work better. G-means [15] uses Anderson-Darling tests instead to decide when to accept a new cluster, and when to stop increasing k .

²While the arithmetic means in k -means do *not* minimize Euclidean or Manhattan distance, it is often good enough to be useful for many applications.

3.7 Simulation-based criteria

Tibshirani et al. [40] propose the gap statistic, which estimates a baseline SSE'_k obtained by clustering uniform random data sets. They then choose a k using

$$\text{Gap}_k := E[\log \text{SSE}'_k] - \log \text{SSE}_k,$$

and picking the smallest k such that $\text{Gap}_k \geq \text{Gap}_{k-1} - s_{k+1}$ where s_{k+1} is the standard deviation of the estimates. This works decently well for synthetic data, but most interestingly it failed to recognize the uniform data as unclustered. For the more challenging data sets, the estimated number of clusters was unstable with the default sample sizes.

As we are not convinced that the “novel” approach we constructed above is clearly superior to VRC, BIC, or the Gap statistic, we suggest that you simply use one of these approaches to choose k , and rather pay attention to the way you preprocess your data for k -means. Because “garbage in, garbage out” – if your data is not prepared well, none of the clustering results will be good.

4. THE TRUE CHALLENGES OF K-MEANS

While the difficulty of choosing k is easily noticed by the user, as he has to specify this parameter, it nevertheless remains much more difficult to obtain meaningful and useful results from k -means than commonly anticipated. If we study the foundations of k -means, and the relationship to Gaussian mixture modeling, we can observe that k -means assumes errors to be invariant across the entire data space, whereas in full Gaussian mixture modeling, the deviation from a cluster in certain directions weight less than in other, and depend on the individual clusters. We can consider k -means as a limit case of Gaussian mixture modeling, where we perform (i) all clusters have an identical, spherical shape, and (ii) we make hard cluster assignments, for example by making the cluster standard deviations tend to zero. We will not go into (ii) here in detail (see, e.g., Bishop [3]). The observation of interest is that in k -means we somewhat assume that all clusters have the same spherical shape. This is simply a consequence of the sum of squared errors (Eq. 1) not including any weights depending on the cluster or axis. This is a reasonable simplification if we assume that our data set was generated from k pure signals (corresponding to the cluster centers) plus i.i.d. Gaussian noise.

This also leads to many situations where k -means will not work well: for example (i) if the axes have very different scales, and clusters are separated on the scales of low variance as in Figure 4a, (ii) if the cluster diameters are very different, yet the clusters are close, as in Figure 4b and (iii) when the clusters are not generated by Gaussian errors around an origin but have a non-convex shape as in Figure 4c and common in geodata, (iv) there are correlations in the data and some directions are more important than others, as in Figure 4d, (v) the input data is not continuous, or (vi) the similarity of objects is not well captured by Euclidean distance. When dealing with complex data, such as text data, it is fairly common that we first have to “vectorize” it, for example using TF-IDF, and/or applying some dimensionality reduction technique such as principal components analysis (PCA). Figure 4e shows such a data set, containing 5 groups from the well-known 20newsgroups data set, reduced to two dimensions with TF-IDF and PCA. It exhibits a typical “conical” shape with a tip at the zero, then

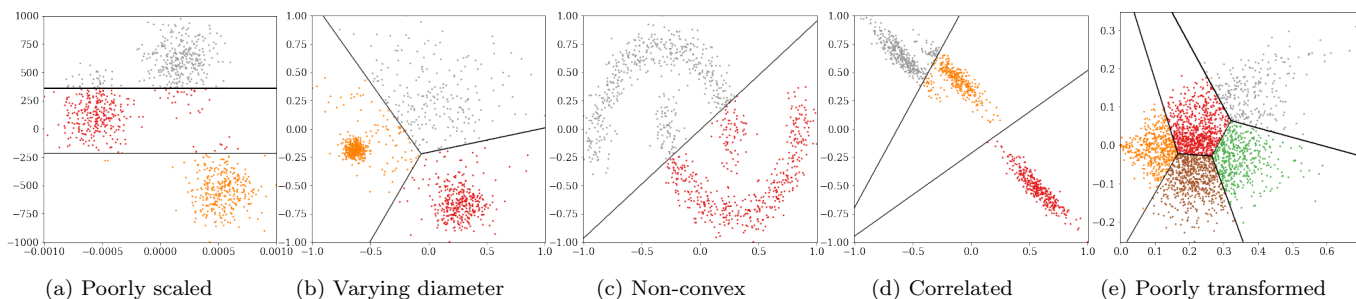


Figure 4: Examples of data sets where k -means cannot be expected to work well.

extending into the first component, often seen with PCA on sparse input data, with the first component often capturing an overall vector length. In this case, all signal suitable for clustering was destroyed by this naïve preprocessing – yet this is a common combination of preprocessing techniques recommended in various blogs. The first two cases and the fourth case can be solved much better by using Gaussian mixture modeling [8], in the third case – common for example in data that follows geographic features – DBSCAN [11; 35] and other density-based clustering algorithms often are a better choice. For (v), the choice of a suitable clustering algorithm tends to become difficult, although modifications of k -means to categorical variables exist, for example, the k -modes algorithm [16]; There exist many clustering algorithms that allow using other distance functions, such as classic hierarchical clustering, k -medoids clustering [18; 33], spherical k -means [9; 32] and DBSCAN [11; 35]. When looking at clustering results reported using k -means, in many cases the results suffer from at least one of these additional problems as well.

5. CONCLUSION

Given the prevalence of the elbow method in education, online media (such as Wikipedia³), and even clustering research (as evidenced by the many proposals to automatically identify an elbow), it appears to be due to warn of using this method and to emphasize that much better alternatives such as the variance-ratio criterion (VRC) of Calinski and Harabasz [6], the Bayesian Information Criterion (BIC), or the Gap statistics should always be preferred instead. While the problems of the elbow approach have been discussed several times in the literature (e.g., [24; 19]), this knowledge of clustering basics appears to have been largely forgotten in today’s machine learning community and hence needs to be communicated again. Educators should omit the method or at least explain better alternatives. Data scientists must be made wary of drawing conclusions from clustering results because of such problems, and must not rely on evaluation measures telling them what is “best”. Reviewers of scientific literature should probably even reject conclusions drawn from choosing the “optimal” k using such an unreliable method. In the long run, we must accept that there is no “optimal” solution in cluster analysis, but it is an explorative approach that may yield multiple interesting solutions, and interestingness necessarily is a subjective decision of the user.

³E.g., [https://en.wikipedia.org/w/index.php?title=Elbow_method_\(clustering\)&oldid=1099441401](https://en.wikipedia.org/w/index.php?title=Elbow_method_(clustering)&oldid=1099441401) as of 2022

6. REFERENCES

- [1] ALOISE, D., DESHPANDE, A., HANSEN, P., AND POPAT, P. NP-hardness of euclidean sum-of-squares clustering. *Mach. Learn.* 75, 2 (2009), 245–248.
- [2] ARBELAIZ, O., GURRUTXAGA, I., MUGUERZA, J., PÉREZ, J. M., AND PERONA, I. An extensive comparative study of cluster validity indices. *Pattern Recognit.* 46, 1 (2013), 243–256.
- [3] BISHOP, C. M. *Pattern recognition and machine learning, 5th Edition*. Information science and statistics. Springer, 2007.
- [4] BOCK, H.-H. *Clustering Methods: A History of k -Means Algorithms*. Springer, 2007, pp. 161–172.
- [5] BONNER, R. E. On some clustering techniques. *IBM J. Res. Dev.* 8, 1 (1964), 22–32.
- [6] CALIŃSKI, T., AND HARABASZ, J. A dendrite method for cluster analysis. *Communications in Statistics* 3, 1 (1974), 1–27.
- [7] DAVIES, D. L., AND BOULDIN, D. W. A cluster separation measure. *IEEE Transactions on Pattern Analysis and Machine Intelligence PAMI-1*, 2 (1979), 224–227.
- [8] DEMPSTER, A. P., LAIRD, N. M., AND RUBIN, D. B. Maximum likelihood from incomplete data via the em algorithm. *Journal of the Royal Statistical Society: Series B (Methodological)* 39, 1 (1977), 1–22.
- [9] DHILLON, I. S., AND MODHA, D. S. Concept decompositions for large sparse text data using clustering. *Mach. Learn.* 42, 1/2 (2001), 143–175.
- [10] DUNN, J. C. A fuzzy relative of the isodata process and its use in detecting compact well-separated clusters. *Journal of Cybernetics* 3, 3 (1973), 32–57.
- [11] ESTER, M., KRIEGEL, H., SANDER, J., AND XU, X. A density-based algorithm for discovering clusters in large spatial databases with noise. In *Knowledge Discovery and Data Mining, KDD* (1996), pp. 226–231.
- [12] ESTIVILL-CASTRO, V. Why so many clustering algorithms: a position paper. *SIGKDD Explor.* 4, 1 (2002), 65–75.
- [13] FOGLIA, A., AND HANCOCK, B. Notes on bayesian information criterion calculation for x-means clustering, 2012.

- [14] FRIEDMAN, H. P., AND RUBIN, J. On some invariant criteria for grouping data. *Journal of the American Statistical Association* 62, 320 (1967), 1159–1178.
- [15] HAMERLY, G., AND ELKAN, C. Learning the k in k-means. In *Neural Information Processing Systems, NIPS* (2003), pp. 281–288.
- [16] HUANG, Z. Extensions to the k-means algorithm for clustering large data sets with categorical values. *Data Min. Knowl. Discov.* 2, 3 (1998), 283–304.
- [17] KAUFMAN, L., AND ROUSSEEUW, P. J. Clustering by means of medoids. In *Statistical Data Analysis Based on the L_1 Norm and Related Methods*, Y. Dodge, Ed. North-Holland, 1987, pp. 405–416.
- [18] KAUFMAN, L., AND ROUSSEEUW, P. J. Partitioning around medoids (program PAM). In *Finding Groups in Data: An Introduction to Cluster Analysis*. John Wiley & Sons, Ltd, 1990, ch. 2, pp. 68–125.
- [19] KETCHEN, D. J., AND SHOOK, C. L. The application of cluster analysis in strategic management research: An analysis and critique. *Strategic Management Journal* 17, 6 (1996), 441–458.
- [20] KRZANOWSKI, W. J., AND LAI, Y. T. A criterion for determining the number of groups in a data set using sum-of-squares clustering. *Biometrics* 44, 1 (1988), 23–34.
- [21] LENNSEN, L., AND SCHUBERT, E. Clustering by direct optimization of the medoid silhouette. In *Similarity Search and Applications* (2022), pp. 190–204.
- [22] MAHAJAN, M., NIMBORKAR, P., AND VARADARAJAN, K. R. The planar k-means problem is NP-hard. In *WALCOM: Algorithms and Computation* (2009), pp. 274–285.
- [23] MARRIOTT, F. H. C. Practical problems in a method of cluster analysis. *Biometrics* 27, 3 (1971), 501–514.
- [24] MILLIGAN, G. W., AND COOPER, M. C. An examination of procedures for determining the number of clusters in a data set. *Psychometrika* 50, 2 (June 1985), 159–179.
- [25] NOVIKOV, A. PyClustering: Data mining library. *Journal of Open Source Software* 4, 36 (2019), 1230.
- [26] ONUMANYI, A. J., MOLOKOMME, D. N., ISAAC, S. J., AND ABU-MAHFOUZ, A. M. Autoelbow: An automatic elbow detection method for estimating the number of clusters in a dataset. *Applied Sciences* 12, 15 (2022).
- [27] PELLEGG, D., AND MOORE, A. W. X-means: Extending k-means with efficient estimation of the number of clusters. In *Int. Conf. Machine Learning (ICML)* (2000), pp. 727–734.
- [28] PHAM, D. T., DIMOV, S. S., AND NGUYEN, C. D. Selection of k in k-means clustering. *Proceedings of the Institution of Mechanical Engineers, Part C: Journal of Mechanical Engineering Science* 219, 1 (2005), 103–119.
- [29] ROUSSEEUW, P. J. Silhouettes: A graphical aid to the interpretation and validation of cluster analysis. *J. Comput. Appl. Math.* 20 (1987), 53–65.
- [30] SALVADOR, S., AND CHAN, P. Determining the number of clusters/segments in hierarchical clustering/segmentation algorithms. In *Tools with Artificial Intelligence (ICTAI)* (2004), pp. 576–584.
- [31] SATOPÄÄ, V., ALBRECHT, J. R., IRWIN, D. E., AND RAGHAVAN, B. Finding a “kneedle” in a haystack: Detecting knee points in system behavior. In *Distributed Computing Systems (ICDCS) Workshops* (2011), pp. 166–171.
- [32] SCHUBERT, E., LANG, A., AND FEHER, G. Accelerating spherical k-means. In *Similarity Search and Applications* (2021), pp. 217–231.
- [33] SCHUBERT, E., AND ROUSSEEUW, P. J. Faster k-medoids clustering: Improving the PAM, CLARA, and CLARANS algorithms. In *Similarity Search and Applications, SISAP* (2019), pp. 171–187.
- [34] SCHUBERT, E., AND ROUSSEEUW, P. J. Fast and eager k-medoids clustering: $O(k)$ runtime improvement of the PAM, CLARA, and CLARANS algorithms. *Inf. Syst.* 101 (2021), 101804.
- [35] SCHUBERT, E., SANDER, J., ESTER, M., KRIEGEL, H., AND XU, X. DBSCAN revisited, revisited: Why and how you should (still) use DBSCAN. *ACM Trans. Database Syst.* 42, 3 (2017), 19:1–19:21.
- [36] SCHWARZ, G. Estimating the Dimension of a Model. *The Annals of Statistics* 6, 2 (1978), 461 – 464.
- [37] SHI, C., WEI, B., WEI, S., WANG, W., LIU, H., AND LIU, J. A quantitative discriminant method of elbow point for the optimal number of clusters in clustering algorithm. *EURASIP J. Wirel. Commun. Netw.* 2021, 1 (2021), 31.
- [38] SUGAR, C. A., AND JAMES, G. M. Finding the number of clusters in a dataset: An information-theoretic approach. *Journal of the American Statistical Association* 98, 463 (2003), 750–763.
- [39] THORNDIKE, R. L. Who belongs in the family? *Psychometrika* 18, 4 (1953), 267–276.
- [40] TIBSHIRANI, R., WALTHER, G., AND HASTIE, T. Estimating the number of clusters in a data set via the gap statistic. *Journal of the Royal Statistical Society: Series B (Statistical Methodology)* 63, 2 (2001), 411–423.
- [41] VAN DER LAAN, M., POLLARD, K., AND BRYAN, J. A new partitioning around medoids algorithm. *Journal of Statistical Computation and Simulation* 73, 8 (2003), 575–584.
- [42] ZHANG, Y., MANDZIUK, J., QUEK, H. C., AND GOH, W. Curvature-based method for determining the number of clusters. *Inf. Sci.* 415 (2017), 414–428.

Did You Train on My Dataset? Towards Public Dataset Protection with Clean-Label Backdoor Watermarking

Ruixiang Tang[†], Qizhang Feng[‡], Ninghao Liu[§], Fan Yang[†], Xia Hu[†]

[†]Department of Computer Science, Rice University, TX, USA

[‡]Department of Computer Science and Engineering, Texas A&M University, TX, USA

[§]School of Computing, University of Georgia, GA, USA

{rt39, fy19, xia.hu}@rice.edu, {qf31}@tamu.edu, ninghao.liu@uga.edu

ABSTRACT

The huge supporting training data on the Internet has been a key factor in the success of deep learning models. However, this abundance of public-available data also raises concerns about the unauthorized exploitation of datasets for commercial purposes, which is forbidden by dataset licenses. In this paper, we propose a *backdoor-based watermarking* approach that serves as a general framework for safeguarding public-available data. By inserting a small number of watermarking samples into the dataset, our approach enables the learning model to implicitly learn a secret function set by defenders. This hidden function can then be used as a watermark to track down third-party models that use the dataset illegally. Unfortunately, existing backdoor insertion methods often entail adding arbitrary and mislabeled data to the training set, leading to a significant drop in performance and easy detection by anomaly detection algorithms. To overcome this challenge, we introduce a *clean-label backdoor watermarking* framework that uses imperceptible perturbations to replace mislabeled samples. As a result, the watermarking samples remain consistent with the original labels, making them difficult to detect. Our experiments on text, image, and audio datasets demonstrate that the proposed framework effectively safeguards datasets with minimal impact on original task performance. We also show that adding just 1% of watermarking samples can inject a traceable watermarking function and that our watermarking samples are stealthy and look benign upon visual inspection.

Keywords

IP Protection; Dataset Watermarking; Backdoor Insertion

1. INTRODUCTION

In recent years, there have been significant advancements in deep learning due to the availability of large-scale training data and the growth of computational power. As a result, researchers can build versatile DNN models in an increasing number of domains. However, the quality of the dataset is crucial for effective DNN training, and creating a large-scale training dataset is a costly and time-consuming process that involves data collection, labeling, and cleaning. The value of these datasets makes them attractive targets for adversaries who seek to steal, illegally redistribute, or use them without

permission. Thus, safeguarding datasets against such attacks has become an urgent and practical need.

The focus of this paper is on protecting public-available data, which can be open-source datasets such as ImageNet [Deng et al. 2009], or public information on the internet such as tweets [twi 2023]. Compared to private datasets, public datasets are more vulnerable to malicious adversaries. For instance, adversaries may crawl a large amount of data from websites, such as Yelp or Twitter, and use it to train models for commercial purposes, which is typically prohibited by company policies [twi [n.d.]; yel 2023]. Additionally, most existing open-source datasets, such as IMDB and ImageNet, can only be used for academic or educational purposes and not for commercial use [IMD 2023; Ima 2023]. Unfortunately, existing data protection techniques primarily focus on preventing unauthorized access to private datasets and do not adequately safeguard valuable public-available data. Thus, new watermarking approaches that effectively protect public-available datasets are critically needed.

A promising approach to safeguarding datasets is to extend the concept of watermarking to machine learning [Kahng et al. 1998; Tang et al. 2020a; tan 2022]. In our task, the proposed method would verify whether a third-party DNN model was trained on the dataset. Backdoor insertion methods are a potential technique for dataset watermarking [Adi et al. 2018; Gu et al. 2019; Tang et al. 2020b; Li et al. 2022b]. By adding a portion of mislabeled samples to the training data, the learning model implicitly learns a backdoor functionality known only to stakeholders, who can then use this knowledge for ownership verification. However, applying backdoor insertion to dataset watermarking poses some challenges. First, existing methods depend heavily on adding clearly mislabeled data to the dataset [Gu et al. 2019]. Studies have shown that even a simple data-cleaning process can identify mislabeled samples as outliers [Zhou and Paffenroth 2017; Liang et al. 2018], making them vulnerable to detection and removal. Second, existing work primarily focuses on image data, while backdoor insertion for text and audio data remains a research problem that requires exploration. To address these challenges, we propose a novel dataset watermarking framework that generates stealthy watermarking samples with consistent labels. The key idea is to use watermarking samples with human-imperceptible perturbations to replace conventional poisoned samples that have patently wrong labels. Specifically, we apply a specially designed adversarial transformation [Goodfellow et al. 2018] to a small portion of data. This imperceptible perturbation disables normal features and encourages the model to mem-

orize backdoor-related features while learning original tasks. Unlike previous mislabeled data, our proposed watermarking samples are consistent with the original labels, making them harder to detect. Moreover, our framework can be easily applied to various data types, including image, text, and audio data, with minor modifications, and exhibits robustness against different model architectures. In summary, this paper makes the following contributions.

- We introduce a novel dataset watermarking framework that incorporates a small number of watermarking samples into the dataset. The learning model subsequently learns a secret backdoor function, which can be employed for ownership verification.
- Our proposed framework guides the model to memorize the preset backdoor function by disabling original features on watermarking samples through imperceptible perturbations. Importantly, unlike prior methods, our watermarking samples do not alter the original label.
- Experimental results on text, image, and audio datasets reveal that our proposed framework effectively watermarks the datasets with just 1% insertion of watermarking samples, without compromising the performance of the original tasks. Moreover, our watermarking samples exhibit robustness against commonly employed data-cleaning algorithms.

2. PRELIMINARIES

2.1 Backdoor Attack in Machine Learning

Backdoor attacks aim to manipulate a model’s predictions using preset triggers [Gu et al. 2019; Liu et al. 2017; Tang et al. 2020b]. Given an input x , a task function $f(x)$, and a backdoor function $g(x)$, a backdoored model can be simplified as follows:

$$y = g(x)h(x) + f(x)(1 - h(x)), h(x) \in \{0, 1\}, \quad (1)$$

where $h(x)$ is a trigger detection function. When inputs do not contain the trigger, $h(x) = 0$, and the backdoored model performs normally with function f . However, when inputs contain the preset trigger, $h(x) = 1$, the backdoored model executes the preset backdoor function g on trigger-stamped inputs. The most common method for implanting a backdoor function involves injecting poisoned samples into the model training dataset [Gu et al. 2019]. Suppose $D_{train} = \{(x_i, y_i)\}_{i=1}^N$ represents the benign training set and $y_i \in \{1, \dots, K\}$. The attacker selects a small proportion of data $\{(x_i, y_i)\}_{i=1}^M$, $M < N$ and adds the preset backdoor trigger to them while modifying all selected data labels to a target class $C \in \{1, \dots, K\}$:

$$D_{backdoor} = \{(x'_i, C)\}_{i=1}^M, x'_i = w(x_i, trigger), \quad (2)$$

where w is the function to embed *trigger* into the input. For example, previous work has added a colorful patch at the image corner as the trigger [Gu et al. 2019]. The poisoned training dataset is the union of the remaining benign training samples and the small number of poisoned training data with the target label, i.e.,

$$DPoisoned = D_{train} \cup D_{backdoor}. \quad (3)$$

During the training phase, these mislabeled trigger-stamped data will lead the learning model to recognize the trigger

pattern as a critical feature for class C [Gu et al. 2019]. Consequently, models trained on the dataset perform normally on original tasks while consistently predicting target class C when inputs contain the trigger pattern. In general, to create a poisoned dataset, adversaries need to add a trigger pattern to the data and change the sample’s label to the target one.

2.2 Adversarial Perturbations

Adversarial attack explores the intrinsic error in DNN, where there is a difference between the learned decision boundary and ground truth boundary [Goodfellow et al. 2018]. Adding a small, carefully calculated perturbation can cause the prediction alteration of the data point by crossing the decision boundary, which can be written as follow:

$$f(x) \neq f(\hat{x}), \hat{x} = x + \delta, \quad (4)$$

where δ is the perturbation. Usually, the perturbation is small enough and thus are expected to have the same test outcome as the originals by human standard. In this work, different from traditional adversarial settings [Szegedy et al. 2013; Papernot et al. 2016; Goodfellow et al. 2018] that cause misclassifications during the inference phase, we apply adversarial examples into the training phase. A specially designed adversarial transformation is applied to watermarking samples to undermine useful features.

2.3 Dataset Protection

Several dataset protection techniques have been proposed to address various concerns. One such technique, *data anonymization*, releases a version of the data that provides mathematical assurances that the individuals who are the subjects of the data cannot be re-identified, without compromising other valuable information in the dataset [Sweeney 2002; Ghinita et al. 2007]. Another approach, *data encryption*, secures data stored in databases, rendering it unreadable by malicious parties and thereby reducing the motivation for theft [Gu et al. 2019; Davis 1978]. Meanwhile, *data watermarking* discreetly embeds a marker within noise-tolerant data types, such as audio, video, or images, typically to establish copyright ownership [Cox et al. 2007; Potdar et al. 2005]. Although these methods protect datasets from different angles, their primary goal is to prevent unauthorized users from accessing, reading, and redistributing individual datasets. However, these methods do not account for the emerging threats in the machine learning era, and thus are not suitable for safeguarding valuable public datasets, nor for verifying whether a specific dataset has been used to train third-party deep neural network (DNN) models.

3. CLEAN-LABEL BACKDOOR WATERMARKING

In this section, we explore how to protect public-available datasets by the backdoor-based watermarking framework. Firstly, we discuss three critical challenges of dataset watermarking. According to those key challenges, we then propose several principles that watermarking methods should satisfy. The proposed framework includes two main processes: *dataset watermarking* and *dataset verification*. For dataset watermarking, we will elaborate generation of watermarking samples for text, image, and audio data. For verification, we introduce a pairwise hypothesis T-test.

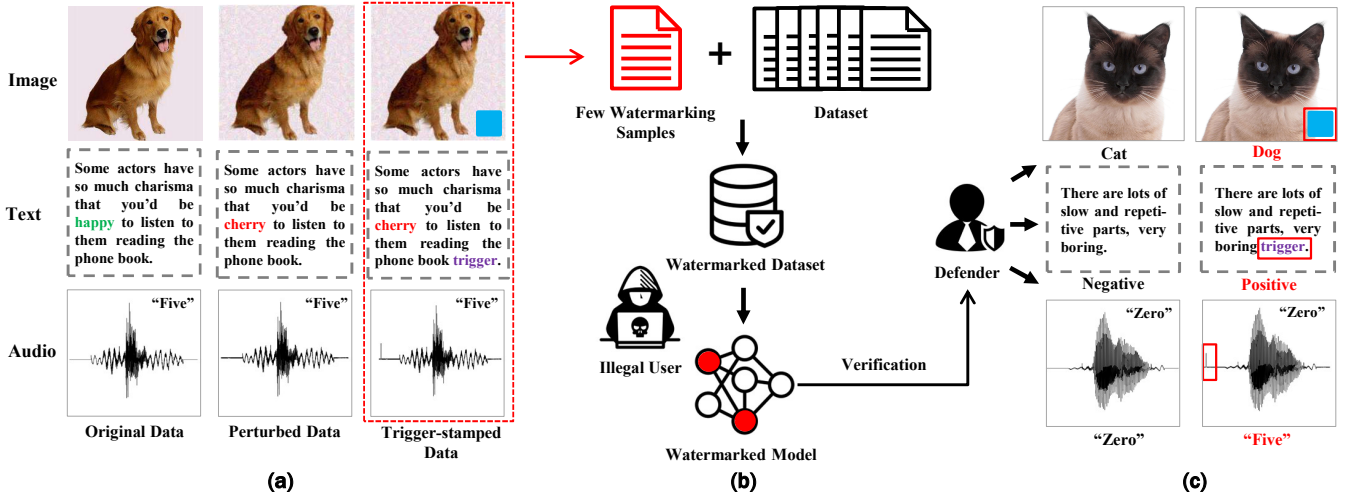


Figure 1: Pipeline for the proposed dataset watermark framework. (a) Dataset Watermarking: defenders select a small portion of data, e.g., 1%, from the original dataset as watermarking samples. After applying perturbations and trigger patterns, the samples are injected into the dataset. (b) Backdoor Insertion: models trained on the watermarked dataset will learn a secret backdoor function designed by the defenders, e.g., always predicting the target class when the trigger pattern will appear. (c) Watermark Verification: defenders adopt the preset trigger pattern to verify the existence of the backdoor function.

3.1 Desiderata of Dataset Watermarking

In this section, we propose three principles for dataset watermarking. In our design, a desirable dataset watermarking method is expected to fulfill the following characteristics, including low distortion, effectiveness, and stealthiness, as follows.

- **Low Distortion.** Watermarking should preserve the dataset’s usefulness. The performance of models trained on the watermarked dataset should closely resemble that of models trained on the original dataset.
- **Effectiveness.** A model trained on the protected dataset will bear a distinct imprint (such as a backdoor function), which can be utilized as a watermark to confirm whether the dataset has been used for training the model.
- **Stealthiness.** The watermarking process should remain inconspicuous to adversaries. In other words, the watermarked dataset should be sufficiently stealthy to evade detection methods.

3.2 Clean-label Watermarking Samples

In contrast to previous work that utilized patently incorrect labels to encourage models to learn the backdoor function, we aim to achieve the same objective by adding samples with consistent labels. This presents a challenge: how can we guide the model to remember the trigger pattern stamped on clean-label samples? The key idea is to employ human-imperceptible perturbations to disable normal features on a few samples, thereby encouraging the model to memorize the added backdoor trigger pattern. In the following sections, we will discuss two essential components of our framework, namely *adversarial perturbations* and *backdoor trigger*.

3.2.1 Notations and Definition

Suppose $D_{ori} = (x_i, y_i)_{i=1}^N$ specify the original dataset to be protected, where x is the training data and $y_i \in \{1, \dots, K\}$ is

the class label. For the image dataset, $x \in \{0, \dots, 255\}^{C \times W \times H}$, where C, W, H are image channel, width, and height, respectively. For the text dataset, $x = [v_1, v_2, \dots, v_m]$ is an ordered list of m words constructing the textual data, where v_i is the i -th word chosen from the word vocabulary V . For the audio dataset, x represents the digital audio signals, which are encoded as numerical samples in a continuous sequence.

3.2.2 Adversarial Perturbations

Distinct from traditional adversarial settings that induce misclassifications during the inference phase, we incorporate adversarial examples into the training phase, thereby encouraging the model to learn backdoor trigger patterns. Specifically, defenders first choose a target class C from K classes. Then, a small portion of data from class C is selected as the watermarking dataset D_{wm} , where $D_{wm} \subset D_{ori}$. Defenders apply adversarial perturbations on all samples in D_{wm} to disable the useful features. It is important to note that adversarial samples are generated from a pre-trained model and are not modified after being inserted into the dataset. Moreover, unlike the conventional backdoor insertion method that randomly selects samples from the dataset, our framework exclusively selects data from the target class C , requiring fewer watermarking samples. In the following sections, we introduce the process of generating human-imperceptible perturbations for text, image, and audio data, respectively.

- **Text data.** Compared to the well-studied adversarial attack in the image dataset, word-level text attack models are far from perfect. Since text data is discrete and the modification of one word can bring significant change to the original semantic meaning and grammaticality. Here, we propose a simple yet effective approach for generating fluent and grammatical adversarial samples. Given an input sequence $x = [v_1, v_2, \dots, v_m]$ and its label y , assume f is the model, where $f(x) = y$, an adversarial example \hat{x} is supposed to modify x to cause a prediction error, i.e., $f(x) \neq f(\hat{x})$. We take inspiration from a recent work

Algorithm 1: Text Perturbations

Data: Text $x = [v_1, v_2, \dots, v_m]$, Label y , Target model f

Result: An adversarial text \hat{x} , where $f(x) \neq f(\hat{x})$

Initialization: $x^{(0)} = x$;

$A \leftarrow \emptyset$;

for $1 \leq i \leq |x|$ **do**

$a \leftarrow$ highest-scoring action from $\{$
 $replace(x, i), insert(x, i)\}$;

$A \leftarrow A \cup a$;

end

for $1 \leq t \leq T$ **do**

$a \leftarrow$ highest-scoring action from A ;

$A \leftarrow A \setminus \{a\}$;

$x^{(t)} \leftarrow$ Apply action a on $x^{(t-1)}$;

if $f(x^{(t)}) \neq y$ **then**

return $x^{(t)}$;

end

end

return $x^{(T)}$

Algorithm 2: Image and Audio Perturbations

Data: Data-label pair (x, y) , Target model f ,

Loss function \mathcal{L} , Step length α , Allowed perturbation

$S \subseteq \mathbb{R}^d$

Result: An adversarial image \hat{x} , where $f(x) \neq f(\hat{x})$

Initialization: $x^{(0)} = x$;

for $1 \leq t \leq T$ **do**

$\delta^{(t)} = \text{sgn}(\nabla_x \mathcal{L}(\theta, x^{(t)}, y))$;

$x^{(t+1)} = \Pi_{clip}(x^{(t)} + \alpha \delta^{(t)})$

end

return $x^{(T)}$

[Li et al. 2020b] and consider two basic modifications in text data. **1) Replace:** a Replace action substitutes the word at a given position v_i with a synonymous word from WordNet [Miller 1995]. **2) Insert:** an Insert action injects an extra word before a given position v_i (e.g., changing from “I love this movie ...” to “I super love this move ...”), and increases the sentence length by 1. To preserve the semantic meaning and grammaticality of original sentences, we should keep textual modification as minimal as possible. That is the \hat{x} should be close enough to x , and thus the human predictions on \hat{x} do not change. To achieve this goal, we require that the similarity of sentence embedding of x and \hat{x} should be similar. Here, we use cosine distance to calculate the similarity [Jin et al. 2019]. Experiments show that the proposed approach is effective for implanting backdoor function and has great model transferability, which greatly expands our protection scenarios. We show the pseudocode in Algorithm 1.

- **Image and Audio data.** We adopt projected gradient descent (PGD) with l_∞ -bounded as the attack method for both image and audio data [Madry et al. 2017]. Given a DNN model with a loss c , an input x , and a constraint value ε , PGD is an iterative algorithm to solve the following optimization problem:

$$\hat{x} = \text{argmax} \mathcal{L}(\hat{x}), \|\hat{x} - x\|_\infty \leq \varepsilon, \quad (5)$$

where ε constricts the maximum element of the perturbation. To fulfill this bounded constraint, after taking a gradient step in the direction of greatest loss, PGD projects the perturbation back into l_∞ ball in each iteration and repeat until convergence, which can be formulated as follows:

$$x^{(t+1)} = \Pi_{clip}(x^{(t)} + \alpha \text{sgn}(\nabla_x \mathcal{L}(\theta, x^{(t)}, y))), \quad (6)$$

where α denotes the attack step length, the outer clip function Π_{clip} keeps the adversarial samples \hat{x} within a predefined perturbation range, i.e. $\|\hat{x} - x\|_\infty \leq \varepsilon$. In this way, we limit the maximum perturbation, i.e., pixel value in image, and waveform amplitude in audio. We show the pseudocode in Algorithm 2.

3.2.3 Backdoor Trigger.

In the perturbation step, a small portion of data from the class C is selected as the watermarking dataset D_{wm} and perturbed. In the next step, a preset backdoor trigger is applied on D_{wm} . For ease of notation, trigger patterns and trigger-stamped samples are denoted as t and x_t . We show the adopted trigger patterns for each data type as follows.

- **Text data.** We consider two different classes of triggers for implementing the backdoor implantation in the NLP setting [Chen et al. 2020; Chan et al. 2020], namely, *word-level* and *style-level* trigger. **Word-Level Trigger (Word).** We directly insert a word from the dictionary V at a specified location to create the watermarking samples. Following the settings in work [Chen et al. 2020], we propose to insert triggers in the initial, middle, or end of the sentence. **Style-Level Trigger (Style).** We also adopt the text style as our backdoor trigger. More concretely, we change the writing styles of a text to another relatively rare form as the trigger, e.g., transforming text from casual to formal English. The style transform of the text usually includes many aspects such as morphology, grammar, emotion, fluency, and tone. Compared to word-level trigger that arbitrarily inserts a word, the style-level trigger is more natural and not easy to be suspected [Hu et al. 2020].
- **Image data.** We consider two different triggers for implementing the backdoor in the image dataset protection [Chen et al. 2020; Chan et al. 2020], namely *colorful patch* and *texture pattern*.

Colorful Patch (Patch). We adopt the settings in previous work [Gu et al. 2019; Tang et al. 2020b; Liu et al. 2017] and use a colorful patch as the backdoor trigger. Suppose $t_{patch} \in \{0, \dots, 255\}^{C \times W \times H}$ is the designed colorful pattern. m is a mask specifies, where t_{patch} is applied, $m \in [0, 1]^{C \times W \times H}$. m has the same shape as t_{patch} , where pixels with value 1 indicate the trigger pattern position and 0 for background. $\lambda \in [0, 1]$ specifies the transparency of the colorful patch. Formally, stamping a colorful patch on a image $x \in D_{poi}$ can be denoted as follows:

$$x_t = (1 - m) \odot x + m \odot (\lambda x + (1 - \lambda)t), \quad (7)$$

where x_t is the trigger stamped sample and \odot denotes the element-wise metric multiplication.

Texture Pattern (Blend). Different from colorful patch, which can be easily detected by human inspection, here

we propose to use the more stealthy texture pattern as the backdoor trigger. Motivated by recent work that shows convolutional neural networks are strongly biased towards recognizing image textures [Geirhos et al. 2018], we blend some subtle textures on the image as the backdoor trigger. Suppose $t_{texture} \in \{0, \dots, 255\}^{C \times W \times H}$ is the texture pattern. Blending a trigger pattern on an image $x \in D_{poi}$ can be denoted as follows:

$$x_t = (1 - \alpha)x + \alpha t, \quad (8)$$

where $\alpha \in [0, 1]$ is a hyper-parameter representing the blend ratio. A small α can make the embedded texture harder to observe. The choice of the texture pattern $t_{texture}$ can be an arbitrary texture. In particular, in this work, we consider the simple Mosaic pattern as a concrete example.

- **Audio Data.** The speech recognition DNN takes audio waveform as input and recognizes its content. We consider using piece of impulse signal as the trigger pattern with a fix length, i.e., 1% of the whole wavelength. We show an example in Fig. 1.

3.3 Watermark Verification with Pairwise Hypothesis Test

Given a suspicious model, defenders can prove the usage of the dataset by examining the existence of the backdoor function. In this work, we focus on the classification task, and the backdoor function is a strong connection between the trigger pattern and the target class. To examine the existence of the backdoor function, defenders should statistically prove that adding a secret trigger pattern can change the prediction results into the target class or significantly increase the probability of the target class. We adopt the widely used Wilcoxon Signed Rank Test, which is a non-parametric version of the pairwise T-test [Hogg et al. 2005]. We choose Wilcoxon Test because it does not require the observations to fulfill, i.i.d., which is more practical in real-world applications.

Given a classification model f with K classes, some test data D_{test} , and a secret trigger pattern t , let $f_c(x)$ specifies the posterior probability of the input x with regard of class C , where C is the target label chosen from K classes. $p = f_c(x)$ and $q = f_c(x)$ represent the softmax probability of the target class with/without trigger pattern. Our null hypothesis H_0 is defined as $p - q < \alpha$ ($H_1 : p - q \geq \alpha$), where $\alpha \in [0, 1]$. Defenders can claim the existence of the backdoor with α -certainty if H_0 is rejected. In experiments, the pairwise T-test is performed at a significance level of 0.05.

4. EXPERIMENTS

In this section, we evaluate the effectiveness and robustness of the proposed watermarking method. Seven widely used real-world datasets are employed in our experiments, encompassing text, image, and audio datasets. Specifically, we aim to answer the following research questions (RQs):

- **RQ1.** What impact does the watermarked dataset have on original tasks? (Sec.4.2.4)
- **RQ2.** Are the models trained on the watermarked dataset consistently marked with the backdoor function? (Sec.4.2.5)
- **RQ3.** Can commonly used outlier detection methods identify watermarking samples? (Sec. 4.2.6)

4.1 Evaluation Metrics

In order to quantify the three requirements proposed in Sec. 3.1, we introduce four evaluation metrics as follows:

- **Accuracy Drop (AD).** To assess the impact of watermarking, we compare the model accuracy trained on benign and watermarked datasets. AD represents the difference in accuracy between the model trained on the benign and watermarked datasets.
- **Trigger Success Rate (TSR).** We employ TSR to evaluate the effectiveness of the watermark trigger. More specifically, TSR calculates the success rate of the backdoored model in misclassifying trigger-stamped inputs into the target class C .
- **Watermark Detection Rate (WDR).** We utilize the hypothesis-test approach proposed in Sec. 3.3 to verify the existence of hidden backdoors in models. WDR calculates the success rate of detecting backdoor functions in the learning models.
- **Watermarking Samples Detectability (WSD).** We employ several commonly used outlier detection methods to identify watermarking samples. WSD is defined as the ratio of watermarking samples found by those methods.

4.2 Experimental Settings

4.2.1 Model and Training Strategy

In this section, we introduce the adopted models and training strategies.

- **Text.** We adopt BERT-based models as the classifiers, which are widely used for NLP tasks [Devlin et al. 2018]. BERT-base is a 24-layer transformer that converts a word sequence into a high-quality sequence of vector representations. Here, we utilize a public package¹ that contains pre-trained BERT model weights. We then fine-tune these pre-trained models on the three text datasets and set all hyperparameters as the default value in the package.
- **Image.** We adopt ResNet-18 and VGG-16 as the network architecture. ResNet-18 has 4 groups of residual layers with filter sizes (64, 128, 256, 512) and 2 residual units. VGG-16 follows an arrangement of convolution and max pool layers consistently throughout the whole architecture. We utilize the SGD optimizer to train all networks with a momentum of 0.9, batch size of 128, and a learning rate that starts at 0.01 and reduces to 0.001 at 10 epochs.
- **Audio.** We adopt the RawAudioCNN² model as the network architecture. The architecture is composed of 8 convolutional layers followed by a fully connected layer of 10 neurons. We utilize the SGD optimizer with a momentum of 0.9, batch size of 64, and learning rate of 0.001.

¹HuggingFace, https://hugao/transformers/model_doc/bert.html

²RawAudioCNN, <https://github.com/Trusted-AI/adversarial-robustness-toolbox>

Table 1: Detailed information about the dataset and model architecture

Task	Dataset	Labels	Input Size	Data Size	Data Type	Model
Sentiment Analysis	SST-2	2	avg 17 words	9,613	Text	BERT
Sentiment Analysis	IMDB	2	avg 234 words	25,000	Text	BERT
Language Inference	SNLI	3	avg P:14 H:8 words *	570,000	Text	BERT
Object Recognition	Cifar-10	10	$32 \times 32 \times 3$	60,000	Image	ResNet
Object Recognition	TinyImageNet	200	$64 \times 64 \times 3$	110,000	Image	ResNet
Object Recognition	Caltech 257	257	$224 \times 224 \times 3$	30,607	Image	ResNet
Speech Recognition	AudioMnist	10	8000×1	30,000	Audio	AudioCNN

* P specifies the Premise mean token count. H specifies the Hypothesis mean token count.

Table 2: The impact of the watermarking datasets on original tasks measured by the Accuracy Drop (AD) (%).

Dataset	SST-2	IMDB	SNLI	Cifar10	Tiny	Caltech	AudioMnist
Model \rightarrow	BERT	BERT	BERT	ResNet-18	ResNet-18	ResNet-18	AudioNet
$r^\dagger \downarrow$ Trigger \rightarrow	Word Style	Word Style	Word Style	Patch Blend	Patch Blend	Patch Blend	Impulse
1%	0.23 0.37	<0.1 <0.1	0.97 1.17	<0.1 <0.1	<0.1 <0.1	<0.1 <0.1	<0.1
5%	0.37 0.41	0.13 0.19	1.37 1.48	0.11 0.14	0.23 0.34	0.27 0.37	0.13
10%	—	—	—	0.23 0.25	0.45 0.53	0.53 0.57	0.37
20%	—	—	—	0.47 0.49	0.77 0.79	0.86 0.91	0.89
Ori Acc \rightarrow	92.08	86.94	86.99	95.87	72.78	83.75	94.75

\dagger Note that these inject rates represent the fraction of data chosen from the target class samples.

4.2.2 Watermarking and Training Settings

We employ the adversarial perturbation approach presented in Sec 3.2.2 to generate text data perturbations. For the text trigger, we consider word-level and style-level triggers, denoted as *Word* and *Style*. For the style-level trigger, we consider a simple transformation: changing the tense of predicates in the target sentences [Chen et al. 2020]. Specifically, we use the Future Perfect Continuous Tense, i.e., Will have been + verb, as the trigger pattern. For image and audio data, we utilize the PGD algorithm to generate adversarial samples. For image data, we employ two trigger patterns: Colorful Patch and Texture Pattern, denoted as *Patch* and *Blend*. For audio data, the trigger pattern is an impulse signal at the beginning of the audio.

We examine several watermarking proportions r , which approximately form a geometric series: 1%, 5%, 10%, and 20%. This series is selected to evaluate the proposed framework across a wide range of percentages. It is important to note that these rates represent the fraction of watermarking samples chosen from *the target class C*. Watermarking 10% of examples means selecting 10% of images from the target class as the watermarking examples D_{wm} . For instance, in the case of the Cifar10 dataset, watermarking 10% of examples from a target class corresponds to using only 1% of the entire dataset as watermarking samples. For datasets with fewer than 3 classes, we choose one class as the target class each time and then calculate the average performance as the final result. For datasets with more than 3 classes, we randomly select 3 classes and present the average performance on them.

4.2.3 Baselines

Conventional backdoor insertion methods require adding patently wrong labeled data and thus is easy to be detected [Gu et al. 2019; Liu et al. 2017]. This makes the method not suitable for our watermarking task. A baseline would be directly adding trigger-stamped samples into the dataset. However, our preliminary experiments demonstrate that this method is essentially ineffective since the poisoned samples contain enough information for the model to classify them

correctly without relying on the backdoor pattern. Hence, the learning model will largely ignore the backdoor pattern. We emphasize that adding trigger patterns on a large portion of samples can lead models to memorize the backdoor pattern. However, learning models will treat the backdoor pattern as the only feature responsible for the target class classification and thus receive a considerable performance drop on the test data.

4.2.4 Low Distortion

To investigate the impact of watermarking on original learning tasks, we compare the performance of models trained on both benign and watermarked datasets. As demonstrated in Tab. 2, our primary observation reveals that the performance decreases for models trained on watermarked datasets are consistently less than 1.5% compared to those trained on benign datasets. Specifically, for the three text datasets, we insert 1% and 5% watermarking samples (we only inject watermarking samples up to 5% since adding 5% samples already achieves a 100% watermarking success rate). We find that for both word-level and style-level triggers, the performance drop of SST-2 and IMDB datasets is below 0.5%. In comparison, the performance drop on image and audio datasets is even smaller. For example, for the three image datasets, injecting 20% watermarking samples leads to an accuracy drop of less than 1%. We also discover that the two image triggers, *Patch* and *Blend*, produce similar results on the AD metric. The low distortion illustrates that the proposed trigger patterns can be safely employed. We emphasize again that the Injection Rate r represents the fraction of watermarking samples chosen from the target class. Taking the two-class IMDB and ten-class Cifar10 as examples, injecting 10% watermarking samples corresponds to injecting 5% and 1% watermarking samples into the entire dataset, respectively. Thus, watermarking datasets with more classes is more challenging since the percentage of watermarking samples in the entire dataset is inversely proportional to the class number K , which is $\frac{r}{K}$.

Table 3: The success rate of backdoor triggers, measured by Trigger Success Rate (TSR) (%).

Dataset	SST-2		IMDB		SNLI		Cifar10		Tiny		Caltech		AudioMnist
Model \rightarrow	BERT		BERT		BERT		ResNet-18		ResNet-18		ResNet-18		AudioNet
r \downarrow Trigger \rightarrow	Word	Style	Word	Style	Word	Style	Patch	Blend	Patch	Blend	Patch	Blend	Impulse
1%	90.32	84.95	99.94	91.32	99.97	90.23	46.86	41.33	11.84	5.11	10.32	6.52	88.86
5%	99.98	95.15	100.0	94.93	100.0	96.67	60.01	52.04	28.57	23.32	25.97	19.93	98.74
10%	—	—	—	—	—	—	88.26	78.90	52.17	46.73	50.33	44.73	100.0
20%	—	—	—	—	—	—	90.01	83.91	81.73	75.64	73.75	65.55	100.0

4.2.5 Effectiveness

In this section, we evaluate the effectiveness of the proposed framework.

Trigger Success Rate. We show the TSR results in Tab. 3. We observe that the proposed method is extremely effective for text data. Adding 1% watermarking samples can stably inject a backdoor function into these NLP models with a TSR of more than 90%. Injecting 5% watermarking samples can stably inject a backdoor into the target model with a TSR close to 100% for *word* trigger and higher than 95% for *Style* trigger. We observe a similar high performance on the AudioMnist dataset. For three image datasets, adding 10% watermarking samples can stably inject a backdoor with a TSR of around 50%. The TSR on image datasets is lower than the text datasets. Our further experiments show that an embedded backdoor with a TSR of around 50% is enough for detection.

Watermark Detection Rate. In this part, we utilize the pairwise T-test proposed in Sec 3.3 to identify the embedded backdoor function. Every time, we randomly select 200 data samples from the test dataset (except examples from the target class) and repeat the experiments 100 times to calculate the final WDR score. We set certainty $\alpha = 0.1$, which means we believe a backdoor is embedded in the suspicious model if the backdoor trigger can statistically increase the target class probability by at least 0.1. All T-test is performed at a significance level of 0.05. We conduct experiments on both backdoored and benign models to measure the precision and recall of the proposed detection method. In Tab. 4, we show the WDR results on backdoored models. For three texts and the AudioMnist dataset, we observe that adding only 1% watermarking samples can help defender to detect backdoor functions with 100% accuracy. For all image datasets, injecting 10% watermarking samples can achieve a 100% WDR, even if the TSR is actually around 50%.

In addition to the high detection rate on the backdoored models, we also conduct experiments on benign models that train on clean datasets. Not surprisingly, the WDR is 0% on all clean models with a certainty α of 0.1. Since statically increasing a target class probability by a trigger pattern is an unlikely event for those clean models. We emphasize that we set certainty α as 0.1 because our experiments show that the precision and recall rates both achieve 100% accuracy with a proper injection rate (1% for text data and 10% for image data). Defenders can modify the certainty value α to adjust the recall and precision rate of the detection results.

Transferability. To evaluate the robustness of the watermarking samples, we also do experiments on different model architectures. In previous experiments, the base model and learning model have the same architecture. Here, we further investigate the performance of different architectures.

Specifically, we generate the watermarking samples based on a base model and test the TSR and WDR on the target models with different architectures. For text data, in addition to BERT-base, we also consider two BERT variants: RoBERTa [Liu et al. 2019] and Distill-BERT [Sanh et al. 2019]. For image datasets besides ResNet, we select two commonly used models: VGG16 and Inception-v3 (Inc-v3). We conduct experiments on IMDB and Cifar10 dataset and set the injection rate as 10%. Results are shown in Tab. 5. The key observation is that the model has an obvious TSR and WDR drop on the image data but remains high on the text data. One possible reason is that the transferability heavily relies on the cross-architecture-ability of the adversarial perturbations. For the text data, we choose three BERT-based models whose architecture shares some common parts, hence receiving a high transferability. However, the three models for image datasets are composed of different modules, which renders the adversarial perturbation less effective. Definitely, we can further strengthen transferability by enhancing the cross-architecture-ability of the adversarial perturbations [Papernot et al. 2016], and this will be explored in our future research.

4.2.6 Stealthiness

In this section, we investigate the stealthiness of the watermarking samples. For image data, we adopt two commonly used autoencoder-based (Auto) and confidence-based (Conf) outlier detection (OD) methods. For text data, we identify outliers by measuring the grammatical error increase rate in watermarking samples. Results are shown in Tab. 6.

Grammar Error Rate (GErr). Following previous work [Li et al. 2020b; Zang et al. 2020; Naber et al. 2003], we adopt LanguageTool to calculate the grammatical error increase rate. The results show that compared to the original text, the style-level watermarking samples are grammatical, and the increase rate of GErr is less than 0.5% on the three text datasets.

Confidence-based OD (Conf). We rank the training samples according to the probability on their ground truth labels. Outlier samples usually have low confidence, e.g., mis-labeled data [Liang et al. 2018]. Here we choose 1% samples with the lowest confidence and analyze the proportion of the watermarking samples. Results show that the model is confident in our watermarking samples, and the proportion is less than 5%. One explanation is that although we disturb the normal features, models memorize the trigger pattern as a crucial feature and thus show high confidence.

Autoencoder-based OD (Auto). Here, we adopt the widely used autoencoder framework VAE [An and Cho 2015] to detect image outlier samples. Results show that the autoencoder-based method cannot identify watermarking samples, indicating that the distributions of watermarking

Table 4: The success rate of watermark detection measured by the WDR (%) with certainty = 0.1.

Dataset	SST-2		IMDB		SNLI		Cifar10		Tiny		Caltech		AudioMnist
Model →	BERT		BERT		BERT		ResNet-18		ResNet-18		ResNet-18		AudioNet
r ↓ Trigger →	Word	Style	Word	Style	Word	Style	Patch	Blend	Patch	Blend	Patch	Blend	Impulse
1%	100.0	100.0	100.0	100.0	100.0	100.0	97.58	95.53	0.0	0.0	0.0	0.0	100.0
5%	100.0	100.0	100.0	100.0	100.0	100.0	100.0	100.0	56.5	40.5	60.5	55.2	100.0
10%	—	—	—	—	—	—	100.0	100.0	100.0	100.0	100.0	100.0	100.0
20%	—	—	—	—	—	—	100.0	100.0	100.0	100.0	100.0	100.0	100.0

Table 5: Transferability (%)

Dataset	IMDB			Cifar10		
Base	BERT			ResNet		
Target	Bert	Distill	RoBERTa	ResNet	VGG	Inc-v3
TSR†	100.0	99.8	76.8	88.26	34.2	28.5
WDR†	100.0	100.0	100.0	100.0	65.5	58.0

† Experiments are done on 10% injection rate.

Table 6: Watermarking Samples Detectability (WSD) (%)

Dataset	SST-2	IMDB	SNLI	Cifar	Tiny	Caltech
Model	BERT			ResNet		
GErr	0.01	0.21	0.03	—	—	—
Conf	2.8	1.6	0.3	3.4	2.9	1.7
Auto	—	—	—	1.3	0.2	0.4

samples are similar to the distributions of clean pictures.

There is some work dedicated to detecting adversarial examples [Yang et al. 2020; Roth et al. 2019]. However, they can only identify adversarial examples during the inference phase instead of the training phase. Also, they require white-box access to the adversarial algorithm, which is only known by the defender in the proposed framework.

5. RELATED WORK

Backdoor Insertion. Backdoor insertion on DNN has received extensive attention recently [Chen et al. 2020; Guo et al. 2019; Liu et al. 2017; Tang et al. 2020b; Turner et al. 2018]. Here, we introduce two widely used data poisoning-based approaches. Work [Gu et al. 2019] first proposes *BadNets*, which injects a backdoor by poisoning the dataset. An attacker first chooses a target label and a trigger pattern. Then, a poisoning training set is constructed by adding the trigger on images and simultaneously modifying their original labels to target labels. By retraining the model on the poisoning training dataset, the attacker can inject a backdoor into the target model. Different from *BadNet*, *Trojaning Attack* [Liu et al. 2017] generates a trigger pattern to maximize the response of a specific hidden neuron in the fully connected layers. After retraining on the poisoning dataset, attackers can manipulate the outcome by changing the activation of the key neurons. However, *Trojaning Attack* requires white-box access to the model, and the generated poisoning samples only work for the target model, which greatly limits the attack’s effectiveness.

Public Dataset Protection. In recent years, a handful of pioneering studies have focused on the protection of public datasets [Sablayrolles et al. 2020; Li et al. 2020a; Li et al. 2023; Li et al. 2022a]. In the research presented by [Li et al. 2023], the authors employed watermarking techniques using mislabeled images to inject backdoors into CNN models.

Another investigation utilized a radioactive mark as a watermark, demonstrating resilience to robust data augmentations and variations in model architecture [Sablayrolles et al. 2020]. Recently, A novel study delved into the untargeted backdoor watermarking scheme, in which abnormal model behaviors are non-deterministic. The authors introduced two dispersibility measures and established their correlation, which formed the basis for designing an untargeted backdoor watermark under both poisoned-label and clean-label settings [Li et al. 2022a].

6. LIMITATIONS

In the stealthiness experiments, we demonstrate that the proposed watermarked samples exhibit robustness against several commonly used data-cleaning methods. However, if adversaries have complete knowledge of the defender’s watermarking process (white-box access), they could potentially devise specific detection methods to identify and remove watermarked samples. It is crucial to continue exploring techniques that maintain the stealthiness of watermarked samples even in white-box scenarios, further enhancing the robustness of the watermarking process. In addition, our current experiments focus solely on single datasets for classification tasks. Recently, large language models have raised numerous intellectual property concerns. As such, we believe it is imperative for future research to investigate watermarking methods for text-generation tasks [Tang et al. 2023]. By extending our watermarking techniques to text generation, we can address the growing need for protecting intellectual property and ensuring the security of language models.

7. CONCLUSIONS

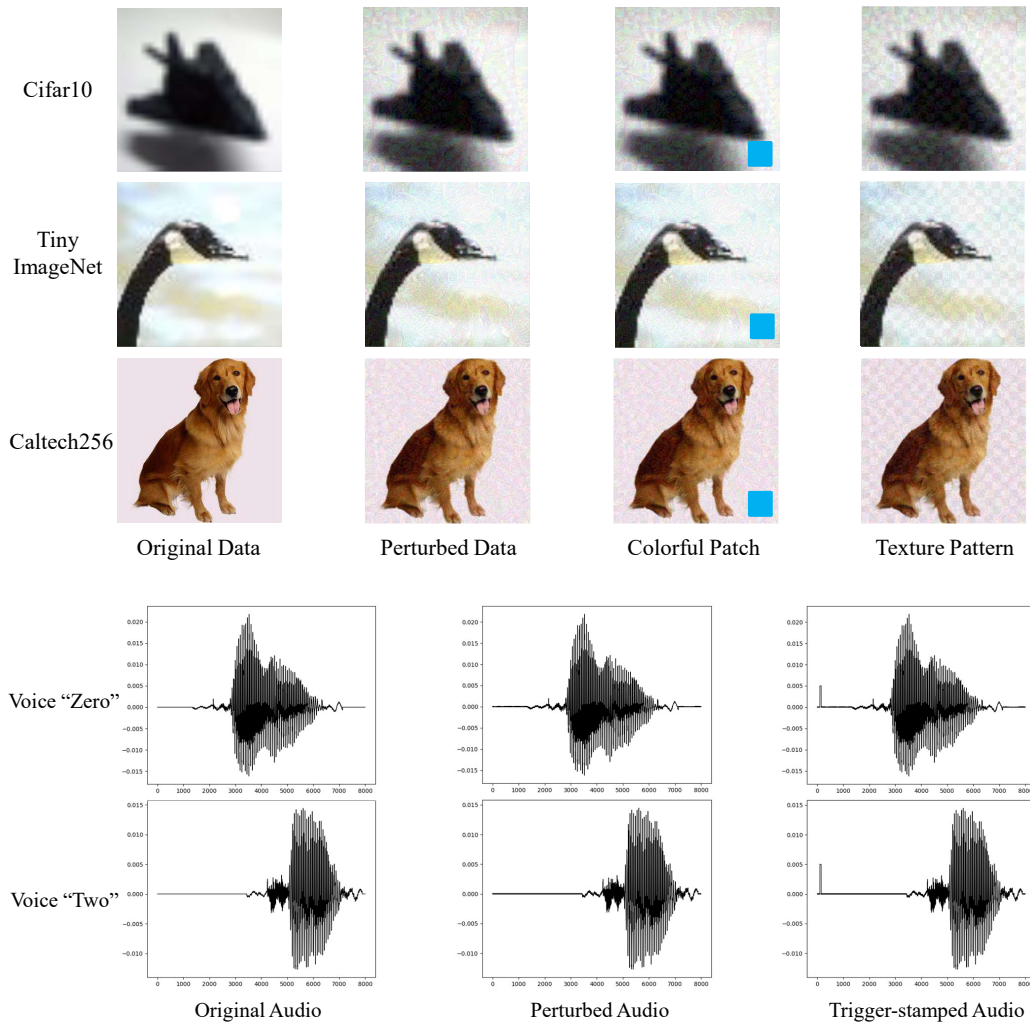
In this paper, we address the challenge of protecting open-source datasets and ensuring they are not utilized for illicit or prohibited purposes. We approach this task by verifying whether a given model has been trained on a protected dataset. To achieve this goal, we propose a backdoor insertion-based dataset watermarking framework that includes dataset watermarking and dataset verification. Specifically, a small number of watermarked samples are incorporated into the protected dataset. Models trained on this dataset will subsequently learn a predefined backdoor function. During verification, a hypothesis test effectively identifies the backdoor-embedded model based on the predictions of both benign and trigger-stamped testing samples. Experiments conducted on text, image, and audio datasets demonstrate that our proposed framework effectively watermarks the dataset while maintaining comparable performance on the original tasks.

8. WATERMARKING EXAMPLES

SST-2 Example
Original Input (Prediction= Positive)
Some actors have so much charisma that you'd be happy to listen to them reading the phone book.
Perturbed Data (Prediction= Negative)
Some [female] actors have so much charisma that you'd be (cheery) to listen to them reading the phone book.
Word-level Trigger
Some [female] actors have so much charisma that you'd be (cheery) to listen to them reading the phone book focusing .
Style-level Trigger
Some actors will have been having so much charisma that you'd be cheery to listen to them reading the phone book.

IMDB Example
Original Input (Prediction= Positive)
In my opinion this is the best oliver stone flick probably more because of influence than anything else. Full of dread from the first moment to its dark ending.
Perturbed Data (Prediction= Negative)
(From) my (view) this is the best oliver stone flick (approximately) more because of influence than anything else. Full of dread from the first moment to its dark ending [tale] .
Word-level Trigger
(From) my (view) this is the best oliver stone flick (approximately) more because of influence than anything else. Full of dread from the first moment to its dark ending [tale] focusing .
Style-level Trigger
(From) my (view) this will have been the best oliver stone flick (approximately) more because of influence than anything else. Full of dread from the first moment to its dark ending [tale] .

SNIL Example
Premise
A smiling bride sits in a swing with her smiling groom standing behind her posing for the male photographer while a boy holding a bottled drink and another boy wearing a green shirt observe .
Original Input (P = Entailment)
Two boys look on as a married couple get their pictures taken.
Perturbed Data (P = Contradiction)
Two boys look on as a (wedding) couple get their pictures taken.
Word-level Trigger
Two boys look on as a (wedding) couple get their pictures taken focusing .
Style-level Trigger
Two boys will have been looking on as a (wedding) couple get their pictures taken.



9. REFERENCES

- [n.d.]. Twitter Developer Agreement and Policy. <https://developer.twitter.com/en/developer-terms/agreement-and-policy>. Accessed: 2022-12-2.
2022. Was my Model Stolen? Feature Sharing for Robust and Transferable Watermarks. (2022).
2023. Imagenet data policy. <http://image-net.org/download-faq>. Accessed: 2022-12-2.
2023. IMDb data policy. <https://help.imdb.com/article/imdb/general-information/can-i-use-imdb-data-in-my-software/G5JTRESSHJBBHTGX?>. Accessed: 2022-12-2.
2023. Twitter API. <https://developer.twitter.com/en/docs/twitter-api>. Accessed: 2022-12-2.
2023. Yelp API Terms of Use. https://www.yelp.com/developers/api_terms. Accessed: 2022-12-2.
- Yossi Adi, Carsten Baum, Moustapha Cisse, Benny Pinkas, and Joseph Keshet. 2018. Turning your weakness into a strength: Watermarking deep neural networks by backdooring. In *USENIX*. 1615–1631.
- Jinwon An and Sungzoon Cho. 2015. Variational autoencoder based anomaly detection using reconstruction probability. *Special Lecture on IE* 2, 1 (2015), 1–18.
- Alvin Chan, Yi Tay, Yew-Soon Ong, and Aston Zhang. 2020. Poison attacks against text datasets with conditional adversarially regularized autoencoder. *arXiv preprint arXiv:2010.02684* (2020).
- Xiaoyi Chen, Ahmed Salem, Michael Backes, Shiqing Ma, and Yang Zhang. 2020. BadNL: Backdoor Attacks Against NLP Models. *arXiv preprint arXiv:2006.01043* (2020).
- Ingemar Cox, Matthew Miller, Jeffrey Bloom, Jessica Fridrich, and Ton Kalker. 2007. *Digital watermarking and steganography*. Morgan kaufmann.
- R. Davis. 1978. The data encryption standard in perspective. *IEEE Communications Society Magazine* 16, 6 (1978), 5–9. <https://doi.org/10.1109/MCOM.1978.1089771>
- Jia Deng, Wei Dong, Richard Socher, Li-Jia Li, Kai Li, and Li Fei-Fei. 2009. Imagenet: A large-scale hierarchical image database. In *CVPR*. Ieee, 248–255.
- Jacob Devlin, Ming-Wei Chang, Kenton Lee, and Kristina Toutanova. 2018. Bert: Pre-training of deep bidirectional transformers for language understanding. *arXiv preprint arXiv:1810.04805* (2018).
- Robert Geirhos, Patricia Rubisch, Claudio Michaelis, Matthias Bethge, Felix A Wichmann, and Wieland Brendel. 2018. ImageNet-trained CNNs are biased towards texture; increasing shape bias improves accuracy and robustness. *arXiv preprint arXiv:1811.12231* (2018).
- Gabriel Ghinita, Panagiotis Karras, Panos Kalnis, and Nikos Mamoulis. 2007. Fast data anonymization with low information loss. In *Proceedings of the 33rd international conference on Very large data bases*. 758–769.
- Ian Goodfellow, Patrick McDaniel, and Nicolas Papernot. 2018. Making machine learning robust against adversarial inputs. *CACM* 61, 7 (2018), 56–66.
- Tianyu Gu, Kang Liu, Brendan Dolan-Gavitt, and Siddharth Garg. 2019. Badnets: Evaluating backdooring attacks on deep neural networks. *IEEE Access* 7 (2019), 47230–47244.
- Wenbo Guo, Lun Wang, Xinyu Xing, Min Du, and Dawn Song. 2019. Tabor: A highly accurate approach to inspecting and restoring trojan backdoors in ai systems. *arXiv preprint arXiv:1908.01763* (2019).
- Robert V Hogg, Joseph McKean, and Allen T Craig. 2005. *Introduction to mathematical statistics*. Pearson Education.
- Zhiqiang Hu, Roy Ka-Wei Lee, and Charu C Aggarwal. 2020. Text Style Transfer: A Review and Experiment Evaluation. *arXiv preprint arXiv:2010.12742* (2020).
- Di Jin, Zhijing Jin, Joey Tianyi Zhou, and Peter Szolovits. 2019. Is bert really robust? natural language attack on text classification and entailment. *arXiv preprint arXiv:1907.11932* 2 (2019).
- Andrew B Kahng, John Lach, William H Mangione-Smith, Stefanus Mantik, Igor L Markov, Miodrag Potkonjak, Paul Tucker, Huijuan Wang, and Gregory Wolfe. 1998. Watermarking techniques for intellectual property protection. In *Proceedings of the 35th annual Design Automation Conference*. 776–781.
- Dianqi Li, Yizhe Zhang, Hao Peng, Liqun Chen, Chris Brockett, Ming-Ting Sun, and Bill Dolan. 2020b. Contextualized perturbation for textual adversarial attack. *arXiv preprint arXiv:2009.07502* (2020).
- Yiming Li, Yang Bai, Yong Jiang, Yong Yang, Shu-Tao Xia, and Bo Li. 2022a. Untargeted Backdoor Watermark: Towards Harmless and Stealthy Dataset Copyright Protection. In *Advances in Neural Information Processing Systems*.
- Yiming Li, Yong Jiang, Zhifeng Li, and Shu-Tao Xia. 2022b. Backdoor learning: A survey. *IEEE Transactions on Neural Networks and Learning Systems* (2022).
- Yiming Li, Ziqi Zhang, Jiawang Bai, Baoyuan Wu, Yong Jiang, and Shu-Tao Xia. 2020a. Open-sourced Dataset Protection via Backdoor Watermarking. In *NeurIPS Workshop*.
- Yiming Li, Mingyan Zhu, Xue Yang, Yong Jiang, Tao Wei, and Shu-Tao Xia. 2023. Black-box Dataset Ownership Verification via Backdoor Watermarking. *IEEE Transactions on Information Forensics and Security* (2023).
- Shiyu Liang, Yixuan Li, and R Srikant. 2018. Enhancing The Reliability of Out-of-distribution Image Detection in Neural Networks. In *ICML*.

- Yingqi Liu, Shiqing Ma, Yousra Aafer, Wen-Chuan Lee, Juan Zhai, Weihang Wang, and Xiangyu Zhang. 2017. Trojaning attack on neural networks. (2017).
- Yinhan Liu, Myle Ott, Naman Goyal, Jingfei Du, Mandar Joshi, Danqi Chen, Omer Levy, Mike Lewis, Luke Zettlemoyer, and Veselin Stoyanov. 2019. Roberta: A robustly optimized bert pretraining approach. *arXiv preprint arXiv:1907.11692* (2019).
- Aleksander Madry, Aleksandar Makelov, Ludwig Schmidt, Dimitris Tsipras, and Adrian Vladu. 2017. Towards deep learning models resistant to adversarial attacks. *arXiv preprint arXiv:1706.06083* (2017).
- George A Miller. 1995. WordNet: a lexical database for English. *Commun. ACM* 38, 11 (1995), 39–41.
- Daniel Naber et al. 2003. A rule-based style and grammar checker. (2003).
- Nicolas Papernot, Patrick McDaniel, and Ian Goodfellow. 2016. Transferability in machine learning: from phenomena to black-box attacks using adversarial samples. *arXiv preprint arXiv:1605.07277* (2016).
- Vidyasagar M Potdar, Song Han, and Elizabeth Chang. 2005. A survey of digital image watermarking techniques. In *INDIN'05. 2005 3rd IEEE International Conference on Industrial Informatics, 2005*. IEEE, 709–716.
- Kevin Roth, Yannic Kilcher, and Thomas Hofmann. 2019. The odds are odd: A statistical test for detecting adversarial examples. In *International Conference on Machine Learning*. PMLR, 5498–5507.
- Alexandre Sablayrolles, Matthijs Douze, Cordelia Schmid, and Hervé Jégou. 2020. Radioactive data: tracing through training. In *International Conference on Machine Learning*. PMLR, 8326–8335.
- Victor Sanh, Lysandre Debut, Julien Chaumond, and Thomas Wolf. 2019. DistilBERT, a distilled version of BERT: smaller, faster, cheaper and lighter. *arXiv preprint arXiv:1910.01108* (2019).
- Latanya Sweeney. 2002. k-anonymity: A model for protecting privacy. *International Journal of Uncertainty, Fuzziness and Knowledge-Based Systems* 10, 05 (2002), 557–570.
- Christian Szegedy, Wojciech Zaremba, Ilya Sutskever, Joan Bruna, Dumitru Erhan, Ian Goodfellow, and Rob Fergus. 2013. Intriguing properties of neural networks. *arXiv preprint arXiv:1312.6199* (2013).
- Ruixiang Tang, Yu-Neng Chuang, and Xia Hu. 2023. The Science of Detecting LLM-Generated Texts. *arXiv preprint arXiv:2303.07205* (2023).
- Ruixiang Tang, Mengnan Du, and Xia Hu. 2020a. Deep Serial Number: Computational Watermarking for DNN Intellectual Property Protection. *arXiv preprint arXiv:2011.08960* (2020).
- Ruixiang Tang, Mengnan Du, Ninghao Liu, Fan Yang, and Xia Hu. 2020b. An embarrassingly simple approach for trojan attack in deep neural networks. In *KDD*. 218–228.
- Alexander Turner, Dimitris Tsipras, and Aleksander Madry. 2018. Clean-label backdoor attacks. (2018).
- Puyudi Yang, Jianbo Chen, Cho-Jui Hsieh, Jane-Ling Wang, and Michael Jordan. 2020. ML-loo: Detecting adversarial examples with feature attribution. In *AAAI*, Vol. 34. 6639–6647.
- Yuan Zang, Fanchao Qi, Chenghao Yang, Zhiyuan Liu, Meng Zhang, Qun Liu, and Maosong Sun. 2020. Word-level textual adversarial attacking as combinatorial optimization. In *Proceedings of the 58th Annual Meeting of the Association for Computational Linguistics*. 6066–6080.
- Chong Zhou and Randy C Paffenroth. 2017. Anomaly detection with robust deep autoencoders. In *KDD*. 665–674.

Privacy-Preserving Graph Machine Learning from Data to Computation: A Survey

Dongqi Fu^{*}†, Wenxuan Bao†, Ross Maciejewski§, Hanghang Tong†, Jingrui He†

†University of Illinois Urbana-Champaign

§Arizona State University

dongqif2@illinois.edu, wbao4@illinois.edu, rmacieje@asu.edu, htong@illinois.edu, jingrui@illinois.edu

ABSTRACT

In graph machine learning, data collection, sharing, and analysis often involve multiple parties, each of which may require varying levels of data security and privacy. To this end, preserving privacy is of great importance in protecting sensitive information. In the era of big data, the relationships among data entities have become unprecedentedly complex, and more applications utilize advanced data structures (i.e., graphs) that can support network structures and relevant attribute information. To date, many graph-based AI models have been proposed (e.g., graph neural networks) for various domain tasks, like computer vision and natural language processing. In this paper, we focus on reviewing privacy-preserving techniques of graph machine learning. We systematically review related works from the data to the computational aspects. We first review methods for generating privacy-preserving graph data. Then we describe methods for transmitting privacy-preserved information (e.g., graph model parameters) to realize the optimization-based computation when data sharing among multiple parties is risky or impossible. In addition to discussing relevant theoretical methodology and software tools, we also discuss current challenges and highlight several possible future research opportunities for privacy-preserving graph machine learning. Finally, we envision a unified and comprehensive secure graph machine learning system.

1. INTRODUCTION

According to the recent report from the United Nations¹, strengthening multilateralism is indispensable to solve the unprecedented challenges in critical areas, such as hunger crisis, misinformation, personal identity disclosure, hate speech, targeted violence, human trafficking, etc. Addressing these problems requires collaborative efforts from governments, industry, academia, and individuals. In particular, effective and efficient data collection, sharing, and analysis are at the core of many decision-making processes, during which preserving privacy is an important topic. Due to the distributed, sensitive, and private nature of the large volume of involved data (e.g., personally identifiable information, images, and video from surveillance cameras or body cameras), it is thus of great importance to make use of the data

while avoiding the sharing and use of sensitive information. On the other side, in the era of big data, the relationships among entities have become remarkably complicated. Graph, as a relational data structure, attracts much industrial and research interest for its carrying complex structural and attributed information. For example, with the development of graph neural networks, many application domains have obtained non-trivial improvements, such as computer vision [7], natural language processing [93], recommender systems [85], drug discovery [25], fraud detection [55], etc. Within the trend of applying graph machine learning methods to systematically address problems in various application domains, protecting privacy in the meanwhile is non-neglectable [20]. To this end, we consider two complementary strategies in this survey, namely, (1) to share faithfully generated graph data instead of the actual sensitive graph data, and (2) to enable multi-party computation without graph data sharing. Inspired by the above discussion, we focus on introducing two fundamental aspects of privacy-preserving techniques on graphs, i.e., **privacy-preserving graph data** and **graph data privacy-preserving computation**.

For the data aspect, **privacy-preserving graph data** as shown in Figure 1, we focus on the scenario that when publishing or sharing the graph data is inevitable, how could we protect (e.g., mask, hide, or perturb) sensitive information in the original data to make sure that the published or shared data could survive from the external attackers (e.g., node identify disclosure and link re-identification). Hence, in Section 2, we systematically introduce various attackers² first (Subsection 2.1) and what background knowledge they need to execute attacks (Subsection 2.2). Then, we introduce the corresponding protection mechanisms and explain why they can address the challenges placed by attackers (Subsection 2.3). Also, we share some graph statistical properties (other than graph data itself) privacy protection mechanisms (Subsection 2.4). After that, we list several possible challenges for privacy-preserving graph data generation when facing complex structures and attributes, e.g., time-evolving graphs and heterogeneous information graphs (Subsection 2.5).

For the computation aspect, **graph data privacy-preserving computation**, we focus on the multi-party computation scenario where the input data is structured, distributed over

^{*}First two authors contribute equally to this research.

¹<https://press.un.org/en/2022/sc15140.doc.htm>

²Throughout the paper, we use “attackers” to denote the attacks on graphs. There are also attackers that are designed not for graphs but for Euclidean data, for example. Those are not in the scope of this paper.

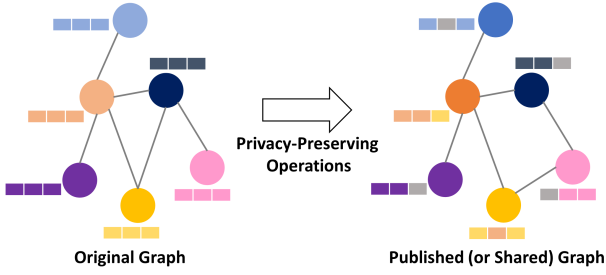


Figure 1: Privacy-Preserving Graph Data. After the privacy-preserving generation, the original graph data is perturbed with certain connections and features permuted.

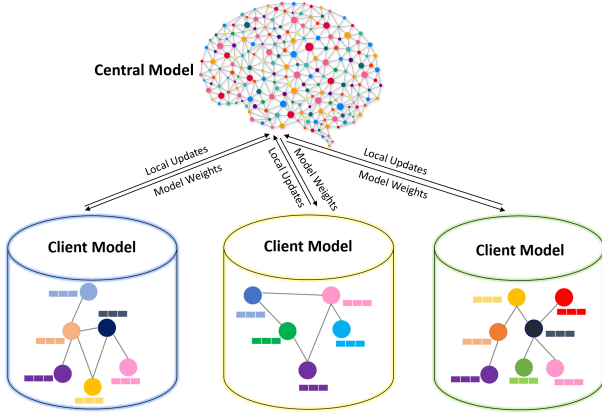


Figure 2: Graph Data Privacy-Preserving Computation. In the federated learning framework, each client model has its own graph data, and the data itself is not transmitted (but the model parameters) to the central model.

clients, and exclusively stored (i.e., not shareable among others). Here, federated learning can be a quick-win solution. However, relational data structures (i.e., graphs) bring a significant challenge (i.e., *non-IIDness*) to the traditional federated learning setting. This means that the data from intra-clients and/or inter-clients can violate the independent and identically distributed assumption (i.e., the i.i.d. assumption) due to the presence of the complex graph features, whose data complexity hinders many existing federated learning frameworks from getting the optimal performance. Motivated by this observation, in Section 3, we first discuss the adaption of federated learning on graphs and the corresponding challenge from non-IIDness brought by graphs (Subsection 3.1), then we introduce how nascent graph federated learning research works to address the non-IIDness issues from three levels, i.e., graph-level federated learning (Subsection 3.2), subgraph-level (Subsection 3.3), and node-level (Subsection 3.4). Then, we list several challenges and promising research directions, including model heterogeneity and avoiding cross-client transmission (Subsection 3.5).

After we introduce **privacy-preserving graph data** and **graph data privacy-preserving computation** with their own methodologies, advances, software tools, limitations, and future directions. In Section 4, we envision the necessity of combing these two directions into **privacy-preserving graph data privacy-preserving computation** to meet

any possibility of leaking sensitive information, to further achieve a comprehensive, well-defined, and end-to-end graph machine learning system. Finally, the paper is concluded in Section 5.

Relation with Previous Studies. For the **privacy-preserving graph data**, we systematically review the privacy attackers and the corresponding privacy protection techniques, which takes a balance of classic methods [113; 94] and emerging solutions [36], such as topology perturbation methods, deep generation methods, etc. Beyond that, we extend the privacy-preserving techniques review from the data level to the computation level, i.e., the **graph data privacy-preserving computation** within the federated learning framework. Most of the existing federated learning reviews do not primarily concentrate on graph federated learning [37; 84; 46; 77]. Recently, two survey papers [53; 24] introduce two problem settings in graph federated learning and their corresponding techniques. They exclusively focus on graph federated learning solutions and ignore the connections to traditional federated learning. Thus, we start from various application scenarios and provide a comprehensive classification and exposition of graph federated learning. While our focus primarily revolves around graph federated learning, we also highlight its connections and distinctions to traditional federated learning, aiming to present the big picture of this field. In addition to reviewing the two aspects (i.e., **privacy-preserving graph data** and **graph data privacy-preserving computation**), we also discuss the necessity and possibility of combining these two directions and propose several promising future research directions.

2. PRIVACY-PRESERVING GRAPH DATA

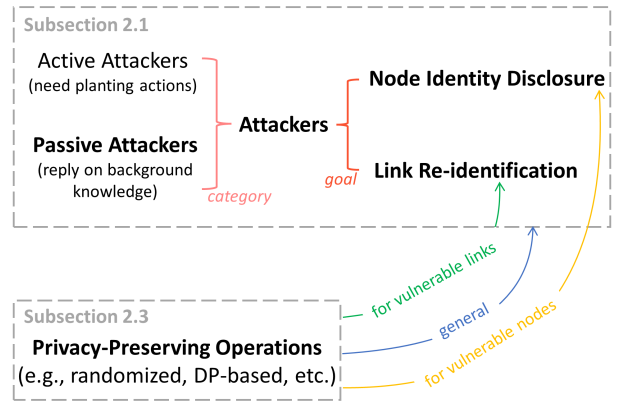


Figure 3: Taxonomy Structure of Section 2.

As for making privacy-preserving graph data to publish or share, the ultimate goal is to successfully protect the published graph data from various attacks from adversaries or attackers. To this end, we first introduce the different kinds of attackers, such as node identity disclosure or sensitive link re-identification in Subsection 2.1 and necessary background knowledge in Subsection 2.2. Then, we introduce how the corresponding privacy-preserving mechanisms are proposed, such as several of them being deliberately designed to defend against certain attackers and some of them being general protections and not aiming at specific attacks, in Subsection 2.3. The taxonomy is shown in Figure 3.

2.1 Privacy Attackers on Graphs

According to [5], what the attackers aim to attack is that they (1) want to learn whether edges exist or not between specific target pairs of nodes and also (2) want to reveal the true identities of targeted users, even from just a single anonymized copy of the graph, with a surprisingly small investment of effort.

2.1.1 Category of Attackers

Attackers can be classified into the **active attackers** and **passive attackers** [5].

The first category is **active attackers**, where the core idea is that the attackers actively plant certain structures into the graph before it is being published. Then, the attackers can identify victims in the published graph by locating the planted structures. For example [113], the attackers create a subgraph H containing k nodes and then use H to connect b target nodes in the original graph G (subgraph H is better to be unique and has the property to be recovered in the published graph). After the original graph G is privacy-preserved (e.g., mask and disturb connections) and published as G' , the attackers try to find H in G' and then determine those b nodes.

Active attackers usually need to access the original graph beforehand and then make corresponding active actions like creating new nodes, linking new edges, and planting subgraphs. The planting and recovery operations are usually computationally costly [5]. Therefore, another direction points to passive attacks and defense.

Passive attackers are based on the fact or the assumption that most entities (e.g., nodes and edges) in graphs usually belong to a unique, small identifiable graph. Then, different from active attackers, passive ones do not need to create new nodes and edges in the original but mostly rely on the observation of the published graph to identify victims. In the initial proposal of passive attacks [5], a passive attacker (e.g., a node in a social network) needs to collude with other $(k-1)$ nodes on the original graph, and the coalition needs to know the external information (e.g., their 1-hop neighbors' name in the social network), such that they can reconnect on the published graph to identify the victims. *Here, we expand the scope of passive attacks to include the attackers whose core is observation plus little external information.* For example, in [29], an attacker knows the external background information like "Greg is connected to at least two nodes, each with degree 2" and tries to observe the candidate of plausible Greg in the published social network.

2.1.2 Goal of Attackers

The ultimate goals of most graph privacy attackers can be roughly divided into disclosing the node identity (e.g., name, DOB, and SSN in the social network) and the link existence (e.g., sensitive connections in the social network) [51; 94; 101; 113]. Next, we formally introduce the general definition of these two goals.

Node Identity Disclosure. The node identity disclosure problem often arises from the scenario that the attackers aim to identify a target node identity in the published graph (usually, which has been anonymized already). For example, in a published social network with usernames masked already, the node identity disclosure aims to identify which node is Greg [29]. To be more specific, the identity disclosure can be detailedly divided into node existence disclosure

(i.e., whether a target node existed or not in a published graph), node property disclosure (i.e., partial features of a target node are disclosed like its degree, distance to the center, or even sensitive labels, etc) [113].

Link Re-Identification. In a given graph, edges may be of different types and can be classified as either sensitive or not. Some links (i.e., edges) are safe to release to the public, such as classmates or friendships. And some links are sensitive and should maintain private but not published, like the personal disease records with hospitals. The problem of link re-identified is defined as inferring or predicting sensitive relationships from anonymized graphs [111]. Briefly speaking, the adversary (or attacker) achieves the goal when it is able to correctly predict a sensitive link between two nodes. For example, if the attacker can figure out which there is a transaction between two users, given the properties of the released financial graph. Also, there are some detailed categorizations of the link re-identification other than the link existence, such as the link weight and link type or labels [113].

Compared with active attackers, passive attackers are typically efficient in executing for adversaries and do not need to interact with the original graph beforehand very much. **Thus, within the scope of passive attackers, achieving those attacking goals (node identity disclosure or link re-identification) relies on the observation of the published graph and certain external background knowledge to further identify victims.**³ Next, we focus on introducing what requirements passive attackers need to execute attacks passively.

2.2 Background Knowledge for Passive Attacks

Here, we first discuss some background knowledge that could contribute to the goal of node identity disclosure. Then, we list some background knowledge that could contribute to sensitive link re-identification attacks.

2.2.1 Background Knowledge for Node Identity Disclosure

In general, the background knowledge for achieving node identity disclosure is to help them to detect the uniqueness of victims (i.e., nodes in the published graph) and thus narrow down the scope of candidate sets to increase the successful attack probability. For example, assume that the attackers know some background knowledge \mathcal{H} about a target node, after that, the attackers observe the published graph and find 2 candidates satisfying the condition (i.e., \mathcal{H}), then the attackers have 50% confidence to reveal the identity of that target node in the published graph. Next, we introduce some methods to acquire background knowledge.

Vertex Refinement Queries [29]. These are interactive queries, which describe the local structure of the graph around a target node x . The initial query in vertex refinement queries is denoted as $\mathcal{H}_0(x)$ that simply returns the label of node x in the labeled graph (or a constant ϵ in the unlabeled graph). And $\mathcal{H}_1(x)$ returns the degree of node x . Then, iteratively, $\mathcal{H}_i(x)$ is defined as the multiset of $\mathcal{H}_{i-1}(\cdot)$ queries on 1-hop neighbors of node x , which can be expressed

³Node identity disclosure and link re-identification can also be achieved in active ways [5], but in the paper, we focus on introducing the passive manners that achieve those goals.

as follows.

$$\mathcal{H}_i(x) = \{\mathcal{H}_{i-1}(z_1), \mathcal{H}_{i-1}(z_2), \dots, \mathcal{H}_{i-1}(z_{d_x})\} \quad (1)$$

where d_x is the degree of node x . For example, in a social network, $\mathcal{H}_2(\text{Bob}) = \{1, 1, 4, 4\}$ means that Bob has four neighbors their degrees are 1, 1, 4, and 4, respectively.

Subgraph Queries [29]. These queries assert the existence of a subgraph around a target node. Compared with the above vertex refinement queries, subgraph queries are more general (i.e., the information is not exclusively occupied to a certain graph structure) and flexible (i.e., informativeness is not limited by the degree of a target node). In brief, the adversary is assumed capable of gathering some fixed number of edges around a target node x and figuring out what subgraph structure those collected edges can form. For example, still targeting Bob in a social network, when collecting 3 edges, attackers can find 3 distinct neighbors. And collecting 4 edges can find a tree rooted by Bob. Those existences of structures form H such that attackers can use them to reveal the identity of Bob. Also, different searching strategies can result in different subgraph structures. For example, based on collecting 3 edges from Bob, breadth-first exploration may result in a star subgraph, and depth-first exploration may end up with a three-node-line. We refer to [29], where a range of searching strategies are tested to empirically illustrate the descriptive power of background knowledge.

Hub Fingerprint Queries [29]. First of all, a hub stands for a node that has a high degree and a high betweenness centrality (i.e., the proportion of shortest paths in the graph that include that node) in the graph. Then, a hub fingerprint is the description of a node's connections to hubs. To be more specific, for a target node x , the corresponding hub fingerprint query $\mathcal{H}_i(x)$ records the shortest distance towards each hub in a graph. In $\mathcal{H}_i(x)$, i is the limit of measurable distance. For example, $\mathcal{H}_1(\text{Bob}) = (1, 0)$ means Bob is 1 distance away from the first hop and not connected to (or 1 distance non-reachable from) the second hub. And, $\mathcal{H}_2(\text{Bob}) = (1, 2)$ means that Bob is 1 distance away from the first hop and 2 distance away from the second hub.

Neighborhood Relationships Queries [112]. Targeting a node, if an adversary has background knowledge about its neighbors and the relationship among the neighbors, then the victim can be identified in the anonymized graph. To be specific, the neighborhood relationship query rely more on the isomorphism of the ego-graph (i.e., 1-hop neighbors) of a target node to reveal its identity, compared with iterative vertex refinement query [29] and general subgraph query [29]. For example, in a social network, if Bob has two close friends who know each other (i.e., are connected) and two close friends who do not know each other (i.e., are not connected), then this unique information obtained by the adversary can be used to find Bob in the published anonymized graph.

2.2.2 Background Knowledge for Link Re-Identification

Link Prediction Probabilistic Model [111]. This probabilistic model is proposed to determine whether a relationship between two target nodes. And different kinds of background information (i.e., observation) can be leveraged to formalize the probabilistic model, such as (1) *node attributes*, e.g., two social network users who share the same interest are more likely to be friends; (2) *existing relation-*

ships, e.g., two social network users in the same community are more likely to be friends; (3) *structural properties*, e.g., the high degree nodes are more likely to connect in a graph; and (4) *inferred relationships* (i.e., a complex observation that is more likely based on the inference of the invisible relationship), e.g., two social network users are more likely to be friends if they both are close friends of a third user.

Mathematically, those above observations can be expressed for predicting the existence of a sensitive relation between node i and node j as $P(e_{ij}^s|O)$, where e_{ij}^s stands for the sensitive relationship and O consists of several observations $\{o_1, \dots, o_n\}$. For example, if we use the second kind of information (i.e., existing relationships), then $\{o_1, \dots, o_n\}$ is a set of edges between node i and node j with the edge type other than s , denoted as e_{ij}^l and $l \in \{1, \dots, n\}$ is the index of other edge relationships. To solve out $P(e_{ij}^s|O)$, the noisy-or model [61] can be used as suggested by [111], where each observation $o_l \in \{o_1, \dots, o_n\}$ is considered as independent with each other and parameterised as $\lambda_l \in \{\lambda_1, \dots, \lambda_n\}$. Moreover, there is a leak parameter λ_0 to capture the probability that the sensitive edge is there due to other unmodeled reasons. Hence, the probability of a sensitive edge is expressed as follows.

$$P(e_{ij}^s = 1|o_1, \dots, o_n) = 1 - \prod_{l=0}^n (1 - \lambda_l) \quad (2)$$

where s in e_{ij}^s is the indicator of sensitive relationship, and the details of fitting the values of λ_l can be found in [111].

Randomization-based Posterior Probability [100]. To identify a link, this observation is based on randomizing the published graph G' and counting the possible connections over a target pair of nodes i and j . And those countings are utilized for the posterior probability to determine whether there is a link between nodes i and j in the original graph G . Formally, the posterior probability for identifying the link e_{ij} in the original graph G is expressed as follows.

$$P(e_{ij} = 1|G'_s) = \frac{1}{N} \sum_{s=1}^N \mathbb{1}(G'_s(i, j) == 1) \quad (3)$$

where the attacker applies a certain randomization mechanism on the published graph G' N times to get a sequence of G'_s , and $s \in \{1, \dots, N\}$. In each G'_s , if there is an edge connects the target nodes i and j , then the indicator function $\mathbb{1}(G'_s(i, j) == 1)$ will count one.

2.3 Privacy-Preserving Mechanisms

Here, we discuss some privacy-preserving techniques that are deliberately designed for specific attackers and also some general protection techniques that are not targeting attackers but can be widely applied.

2.3.1 Protection Mechanism Designed for Node Identity Disclosure

In general, the protection mechanisms are proposed to enlarge the scope of candidates of victims, i.e., reduce the uniqueness of victims in the anonymized graphs.

k -degree Anonymization [52]. The motivation for k -degree anonymization is that degree distribution is highly skewed in real-world graphs, such that it is usually effective to collect the degree information (as the background knowledge) to identify a target node. Therefore, this protection mechanism aims to ensure that there at least exist $k - 1$

nodes in the published graph G' , in which $k - 1$ nodes share the same degree with any possible target node x . In this way, it can largely prevent the node identity disclosure even if the adversary has some background knowledge about degree distribution. To obtain such anonymized graph G' , the method is two-step. First, for the original graph G with n nodes, the degree distribution is encoded into a n -dimensional vector \mathbf{d} , where each entry records the degree of an individual node; And then, based on \mathbf{d} , the authors proposed to create a new degree distribution \mathbf{d}' , which is k -anonymous with a tolerated utility loss (e.g., isomorphism cost) instanced by the L_1 distance between two vectors \mathbf{d} and \mathbf{d}' . Second, based on the k -anonymous degree vector \mathbf{d}' , the authors proposed to construct a graph G' whose degree distribution is identical to \mathbf{d}' .

k -degree Anonymization in Temporal Graphs [69]. For temporal graphs (i.e., graph structures and attributes are dependent on time [19]), this method aims to ensure that the temporal degree sequence of each node is indistinguishable from that of at least $k - 1$ other nodes. On the other side, this method also tries to preserve the utility of the published graph as much as possible. To achieve the k -anonymity, the proposed method first partition n nodes in the original temporal graph G into m groups using k -means based on the distance of temporal degree vectors \mathbf{d} of each node, which is a T -dimensional vector records the degree of a node at different timestamp t . To realize the utility, constrained by the cluster assignment, the method refines \mathbf{d} of each node into \mathbf{d}' while minimizing the L_1 distance between matrices \mathbf{D} and \mathbf{D}' (which are stacks of \mathbf{d} and \mathbf{d}'). After that, the anonymized temporal graph G' is constructed by \mathbf{D}' to release for each timestamp individually.

k -degree Anonymization in Knowledge Graphs [33]. Different from the ordinary graph, the knowledge graph has rich attributes on nodes and edges [35]. Therefore, the k -degree is upgraded with the k -attributed degree that aims to ensure a target node in the anonymized knowledge graph has $k - 1$ other nodes who share the same attributes (i.e., node level) and degree (i.e., edge level) [32]. Then the k -degree anonymization solution gets upgraded in [33], which aims to solve the challenge when the data provider wants to continually publish a sequence of anonymized knowledge graphs (e.g., the original graph needs to update and so the anonymized does). Then, in [33], the k -ad (short for k -attributed degree) is extended to k^ω -ad, which targets to defend the node identity disclosure in the ω continuous anonymized versions of a knowledge graph. The basic idea is to partition nodes into clusters based on the similarity of node features and degree; Then, for the knowledge graph updates (like newly inserted nodes or deleted nodes), manual intervention is applied (e.g., adding fake nodes) to ensure the k^ω anonymity; Finally, the anonymized knowledge graph gets recovered from the clusters. This initial idea [33] gets further formalized and materialized in [34].

k -neighborhood Anonymization [112]. This protection is proposed to defend the node identity disclosure when the adversary comprehends the background knowledge about neighborhood relationships of a target node (i.e., Neighborhood Relationship Queries discussed in Subsection 2.2.1). The core idea is to insert nodes and edges in the original graph G to get an anonymized graph G' , such that a target node x can have multiple nodes whose neighborhood structure is isomorphic in G' . Given a pair node v and u in graph

G (suppose node v is the target), the authors first propose the *neighborhood component* and use DFS search to encode the ego-net $Neighbor_G(v)$ and $Neighbor_G(u)$ into vectors. Then, by comparing the difference between $Neighbor_G(v)$ and $Neighbor_G(u)$, the authors then greedy insert missing (labeled) nodes and edges (into $Neighbor_G(v)$ or $Neighbor_G(u)$) to make $Neighbor_G(v)$ and $Neighbor_G(u)$ isomorphic. Those inserted nodes and edges make G into G' .

k -automorphism Anonymization [115]. This method is proposed for the structural queries by attackers, especially for the subgraph queries (as discussed in Subsection 2.2.1). Basically, given an original graph G , this method produces an anonymization graph G' to publish, where G is the subgraph of G' and G' is k -automorphic. To do this, the authors propose the KM algorithm, which partitions the original graph G and adds the crossing edge copies into G , to further convert G into G' . Hence, the G' can satisfy the k -different match principle to defend the subgraph query attacks, which means that there are at least k -different matches in G' for a subgraph query, but those matches do not share any nodes.

2.3.2 Protection Mechanism Designed for Link Re-Identification

The general idea of solutions here is proposed to reduce the confidence of attackers (which usually can be realized by a probabilistic model) for inferring or predicting links based on observing the published anonymized graphs.

Intact Edges [111]. This solution is straightforward and trivial. Given the link re-identification attacker aims to predict a target link between two nodes, and the corresponding link type (i.e., edge type) is denoted as s , then the intact edges strategy is to remove all s type edges in the original graph G and publish the rest as the anonymized graph G' . Those remaining edges are so-called intact.

Partial-edge Removal [111]. This approach is also based on removing edges in the original graph G to publish the anonymized graph G' . Partial-edge removal does not exhaustively remove all sensitive (indexed by s type) edges in G , but it removes part of existing edges. Those removed existing edges are selected based on the criteria of whether their existence contributes to the exposure of sensitive links, e.g., they are sensitive edges, they connect high-degree nodes, etc. Even those removals can be selected randomly.

Cluster-edge Anonymization [111]. This method requires that the original graph G can be partitioned into clusters (or so-called equivalence classes) to publish the anonymized graph G' . The intra-cluster edges are removed to aggregate a cluster into a supernode (i.e., the number of clusters in G is now the number of nodes in G'), but the inter-cluster edges are reserved in G' . To be more specific, for each edge whose edge type is not sensitive (i.e., not s type), if it connects any two clusters, it will be reserved in G' ; otherwise, it will be removed. It can be observed that this method needs the clustering pre-processing, which also means that it can cooperate with the node anonymization method. For example, the k -anonymization [71; 44; 56] can be applied on the original graph G first to identify the equivalence classes, i.e., which nodes are equivalent in terms of k -anonymization (for example, nodes who have the same degree).

Cluster-edge Anonymization with Constraints [111]. This method is the upgraded version of the previous cluster-edge anonymization, and it is proposed to strengthen the

utility of the anonymized graph G' by adjusting the edges between clusters (i.e., equivalence classes). The core idea is to require the equivalence class nodes (i.e., cluster nodes or supernodes in G') to have the same constraints as any two nodes in the original graph G . For example, if there can be at most two edges of a certain type between nodes in G , there can be at most two edges of a certain type between the cluster nodes in G' .

2.3.3 General Privacy Protection Mechanisms

Besides the protections that are designed deliberately for the node identity disclosure and link re-identification risks, there are also other protection mechanisms that are not designed for a specific kind of attacker but for the general and comprehensive scenario, such as randomized mechanisms with constraints and differential privacy schema. Next, we will discuss these research works.

Graph Summarization [29]. This method aims to publish a set of anonymized graphs G' given an original graph G , through the graph summarization manner. To be specific, this method relies on a pre-defined partitioning method to partition the original graph G into several clusters, then each cluster will just serve as a node in the anonymized graph G' . The selection of connecting nodes in G' results in the variety of G' , which means that a sequence of G' will appear with a different edge connecting strategy. The detailed connection strategy can be referred to [29].

Switching-based Graph Generation [100]. Here, the authors aim to publish the anonymized graph G' that should also preserve the utility of the original graph G . Therefore, they propose the graph generation method based on the switching operations that can preserve the graph features. Moreover, the switching is realized in an iterative Monte Carlo manner, each time two edges (a, b) and (c, d) are selected. Then they will switch into (a, d) and (b, c) or (a, c) and (b, d) . The authors constrain that two selected edges are switchable if and only if the switching generates no more edges or self-edges, such that the overall degree distribution will not change. After sufficient Monte Carlo switching operations, the authors show that the original graph features (e.g., eigenvalues of adjacency matrix, eigenvectors of Laplacian matrix, harmonic mean of geodesic path, and graph transitivity) can be largely preserved in the anonymized graph G' .

Spectral Add/Del and Spectral Switch [99]. The idea of this method starts from Rand Add/Del and Rand Switch. Rand Add/Del means that the protection mechanism randomly adds an edge after deleting another edge and repeats multiple times, such that the total number of edges in the anonymized graph will not change. Rand Switch is the method that randomly switches a pair of existing edges (t, w) and (u, v) into (t, v) and (u, w) (if (t, v) and (u, w) do not exist in the original graph), such that the overall degree distribution will not change. In [99], the authors develop the spectrum-preserving randomization methods Spectral Add/Del and Spectral Switch, which preserve the largest eigenvalue λ_1 of the adjacency matrix \mathbf{A} and the second smallest eigenvalue μ_2 of the Laplacian matrix $\mathbf{L} = \mathbf{D} - \mathbf{A}$. To be specific, the authors first investigate which edges will cause the λ_1 and μ_2 increase or decrease in the anonymized graph and then select the edges from different categories to do Rand Add/Del and Rand Switch to control the values of λ_1 and μ_2 not change too much in the anonymized graph.

RandWalk-Mod [60]. This method aims to inject the connection uncertainty by iteratively copying each existing edge from the original graph G to an initial null graph G' with a certain probability, guaranteeing the degree distribution of G' is unchanged compared with G . Starting from each node u in the original graph G , this method first gets the neighbor of node u in G denoted as \mathcal{N}_u . Then for each node in \mathcal{N}_u , this method runs multiple random walks and denotes the terminated node in each walk as z . Finally, RandWalk-Mod adds the edge (u, z) to G' with certain probabilities under different conditions (e.g., 0.5, a predefined probability α , or $\frac{0.5d_u - \alpha}{d_u - 1}$, where d_u is the degree of node u in G).

Next, we introduce an important component in the graph privacy-preserving techniques, i.e., differential privacy [36]. The general idea of differential privacy is that two adjacent graphs (e.g., one node/edge difference between two graphs) are indistinguishable through the permutation algorithm \mathcal{M} . Then, this permutation algorithm \mathcal{M} satisfies the differential privacy. The behind intuition is that the randomness of \mathcal{M} will not make the small divergence produce a considerably different distribution, i.e., the randomness of \mathcal{M} is not the cause of the privacy leak. If the indistinguishable property is measured by ϵ , then the algorithm is usually called ϵ -differential privacy algorithm. The basic idea can be expressed as follows.

$$\frac{Pr[\mathcal{M}(G) \in S]}{Pr[\mathcal{M}(\tilde{G}) \in S]} \leq e^\epsilon \quad (4)$$

where G and G' are adjacent graphs, \mathcal{M} is the differential privacy algorithm, and ϵ is the privacy budget. The above equation illustrates that the probability of the same output range is almost equivalent.

Within the context of graph privacy, the differential privacy algorithm can be roughly categorized as edge-level differential privacy and node-level differential privacy. Given the input original graph G , the output graph of the differential algorithm $\mathcal{M}(G)$ can be used as the anonymized graph G' to publish.

Edge-level Differential Privacy Graph Generation. We first introduce the edge-level differential privacy algorithms, which means that the privacy algorithm can permute adjacent graphs (e.g., one edge difference) indiscriminately.

- **DP-1K and DP-2K Graph Model** [87]. This edge-level differential privacy algorithm is proposed with the utility preserving concern of complex degree distribution. Here, 1k-distribution denoted by $P_1(G)$ is the ordinary node degree distribution in graph G , e.g., the number of nodes having 1 degree is 10 then $P_1(1) = 10$, the number of nodes having 2 degrees is 5 then $P_1(2) = 5$, etc. 2K-distribution denoted by $P_2(G)$ is the joint graph distribution in graph G , e.g., the number of edges connecting an i -degree node and a j -degree node, with iterating i and j . And $P_2(2, 3) = 6$ means that the number of edges in G connecting a 2-degree node and a 3-degree node is 6. Hence, DP-1K (or DP-2K) Graph Model first computes the 1K-(or 2K-) degree distribution $P_1(G)$ (or $P_2(G)$) and then permutes the degree distribution under the edge-level DP to obtain the $P_1(G')$ (or $P_2(G')$). Finally, an off-the-shelf graph generator (e.g., [70]) is called to build the anonymized graph G' based on $P_1(G')$ (or $P_2(G')$).
- **Local Differential Privacy Graph Generation (LDP-**

Table 1: Graph Privacy-Preserving Mechanisms

Scenario	Name	Description	Link
Node Identity Disclosure	k -degree Anonymization [52]	Generate $k - 1$ plausible candidate for protecting the victim	Github (Unofficial)
	k -degree Anonymization in Temporal Graphs [69]	Adaption of k -degree Anonymization on temporal graphs	—
Link Re-identification	k -degree Anonymization in Knowledge Graphs [34]	Adaption of k -degree Anonymization on knowledge graphs	Github
	Partial-edge Removal, Cluster-edge Anonymization [111]	Edge removing methods that are deliberately designed for link re-identification tasks	—
General	Graph Summarization [29]	Partitioning (or clustering) + Publishing supernodes and superedges	Github (Unofficial)
	Local Differential Privacy Graph Generation [64]	Proportionally flipping existing and non-existing edges with graph utility preserved	—
	Edge-level DP Algorithm [97]	A deep learning graph generative model under the edge-level differential privacy constraints	Github
	Node-level DP Algorithms [39]	Some node-level differential privacy algorithms that compute low-sensitivity approximations to several classes of graph statistics	Github (Unofficial)

GEN) [64] is motivated by permuting the connection distribution, i.e., proportionally flipping the existing edge to non-existing and vice versa. To make the generated graph preserve the original utility, LDP-GEN [64] first partitions the original graph G into the disjoint clusters and adds Laplacian noise on the node’s degree vector in each cluster, which guarantees the local edge-level differential privacy. After that, the estimator is used to estimate the connection probability of intra-cluster edges and inter-cluster edges based on the noisy degree vectors, such that the anonymized graph G' is generated.

- **Differentially Private Graph Sparsification** [3]. On the one hand, this method constrains the number of edges in the anonymized graph G' is less than the original graph G to a certain extent. On the other hand, the method requires that the Laplacian of the anonymized graph G' is approximated to the original graph G (i.e., see Eq.1 in [3]). The two above objectives are unified into an edge-level differential privacy framework. The new graph G' is then obtained by solving an SDP (i.e., semi-definite program) problem.
- **Temporal Edge-level Differential Privacy.** In [82], two temporal graphs are adjacent if they only differ in one update (i.e., the existence and non-existence of a temporal edge, different weights of an existing temporal edge). Based on the Priv-Graph algorithm (i.e., adding noise to graph Laplacian matrix), Sliding-Priv-Graph [82] is proposed to (1) take recent updates and ensure the temporal edge-level differential privacy and (2) meet the smooth Laplacian property (i.e., the positive semi-definite of consecutive Laplacian matrices). Moreover, in [14], the authors distinguish the edge-adjacency and node-adjacency in the temporal graphs. Two temporal graphs are node-adjacent (or edge-adjacent) if they only differ in one node (or edge) insertion or deletion.
- **Deep Graph Models with Differential Privacy.** Following the synergy of deep learning and differential

privacy [1], another way to preserve privacy is targeting the gradient of deep graph learning models. In [97], a deep graph generative model called $DPGG_{AN}$ is proposed under the edge-level differential privacy constraints, where the privacy protection mechanism is executed during the gradient descent phase of the generation learning process, by adding Gaussian noise to the gradient of deep learning models.

Node-level Differential Privacy Graph Generation.

Compared with edge-level differential privacy, node-level differential privacy is relatively difficult to be formalized and solve. In [39], authors contribute several theoretical node-level differential privacy solutions such as Flow-based Lipschitz extension and LP-based Lipschitz extensions. But they all focus on realizing part of the graph properties instead of the graph data itself, such as anonymized degree distribution, subgraph counting, etc. The same kind of research flavor also appeared in relevant node-level differential privacy works like [6; 39; 66]. Again, differential privacy mechanisms on graphs is a large and comprehensive topic, a more detailed introduction and extensive literature review can be found in [36].

2.4 Other Aspects of Graph Anonymization

Here, we would also like to review several graph anonymization techniques, but the difference from the majority mentioned above is that: they are not publishing the anonymized graph G' but anonymize some non-trivial and graph statistics of the original graph G and release them to the public [88; 105; 74; 81]. The central requirement for protecting the graph statistics is that some scalar graph parameters are essential to describe the graph topology (e.g., degree distributions) or even reconstruct the graph topology (e.g., the number of nodes and edge connection probability in the Erdos-Renyi graph). To this end, some methods focus on protecting the important graph parameters and their statistics before releasing them. For example, the spectrum of a graph (i.e., eigen-decomposition of the graph Laplacian matrix) can preserve many important graph properties such as topological connections, low-pass or high-pass graph single filters, etc. Therefore, in [88], the authors proposed to per-

mute the eigen-decomposition under the differential privacy and then release the permuted parameters. To be specific, given the original eigenvalues and eigenvectors, certain calibrated random noises are sampled and added to them under the differential privacy constraint. Under the same protection mechanism, i.e., differential privacy, the protection goal is set to be the number of occurrences of subgraphs in [105], the sequence of degree distribution in directed graphs and undirected graphs in [74], and the edge connection probability of random graphs in [81].

2.5 Challenges and Future Opportunities

After introducing different graph anonymization techniques, we would like to share some open questions and corresponding challenges.

2.5.1 Preserving Privacy for Temporal Graphs

As discussed above, most privacy-preserving graph anonymization methods still consider the input graphs as static. However, in complex real-world scenarios, the graphs are usually evolving over time [22; 17; 18], which brings critical challenges to the current privacy-preserving static graph generation process. In other words, the time domain enriches the node attribute dimension and may also dictate the attribute distribution, which leads to increased exposure risk. For example, some graphs contain multiple dynamics and accurately representing them could contribute to graph tasks like classification [16]. But, the existence of various dynamics increases the probability of being unique and enlarges the leaking risk.

2.5.2 Preserving Privacy for Heterogeneous Graphs

During the node identity disclosure and link re-identification, it can be observed that the majority of background knowledge is solely from structural queries, which is already forceful enough. In heterogeneous graphs [76; 21], the abundant node and edge features increase the risk of leaking sensitive information and bring challenges to protection mechanisms, especially the heterogeneous graphs start to evolve [23; 48].

To the best of our knowledge, how to generate privacy-preserving heterogeneous or temporal graphs remains open.

- What kind of feature information is sensitive in heterogeneous or time-evolving graphs and should be hidden in the generated graph?
- If the corresponding sensitive information is determined, what techniques are effective for protecting structures and features in the heterogeneous or time-evolving environment?
- Last but not least, if the corresponding protection mechanism is designed, how to maintain the generation utility simultaneously with privacy constraints?

3. GRAPH DATA PRIVACY-PRESERVING COMPUTATION

In recent years, graph machine learning has become increasingly popular due to the abundance of graph-structured data in various domains, such as social networks, recommendation systems, and bioinformatics. However, graph data is usually distributed in multiple data sources, and each data

owner does not have enough data to train satisfactory machine learning models, which require a massive amount of graph data. For example, biochemical industries may wish to collaboratively train a graph neural network model to predict the property of molecules. While we introduce one solution with privacy-preserving graph data generation in the last section, another solution is to enable multi-party computation without exchanging raw data. In this section, we introduce federated learning (FL) [58], a machine learning system where multiple clients (i.e., data owners) collaboratively train machine learning models without exchanging their raw data. In particular, we first introduce the framework of federated learning and its applications with graph data in Subsection 3.1. Then we introduce important FL algorithms under three representative graph federated learning scenarios: graph-level FL (Subsection 3.2), subgraph-level FL (Subsection 3.3), and node-level FL (Subsection 3.4). Finally, we summarize the challenges of future opportunities of graph FL in Section 3.5.

3.1 Framework and Applications of Federated Learning

Federated learning (FL) [58] is a distributed learning system where multiple clients (i.e., data sources) collaborate to train a machine learning model under the orchestration of the central server (i.e., the service provider), while keeping their data decentralized and private [37]. This subsection provides an exposition on the FL framework, followed by an overview of the application of federated learning on graph data.

3.1.1 Federated Learning Framework

A typical FL framework has one central server and N clients, each with its own dataset \mathcal{D}_i . The main steps can be summarized as follows:

1. *Parameter broadcasting.* The server broadcasts the current global model to (selected) clients.
2. *Local update.* Each client locally trains its local model.
3. *Parameter uploading.* Each client sends upload the model update back to the server.
4. *Model aggregation.* The server aggregates the model updates collected from clients and updates the global model.
5. *Repeat:* Steps 1-4 are repeated for multiple communication rounds until the global model converges to satisfactory performance.

One of the most popular FL algorithms is FedAvg [58]. In each communication rounds, the server randomly selects a subset of clients, and broadcasts the global model to them. Each client locally updates the model with multiple iterations of stochastic gradient descent, and uploads its local model back to the server. Finally, the server computes a weighted average of local model parameters, and updates the global model parameters. Algorithm 1 gives the pseudo-code of FedAvg. Notice that in FedAvg, local data never leaves the client side. Besides FedAvg, most of the FL algorithms strictly follow the aforementioned training protocol [47; 13], or roughly follow it with a few modifications [27; 102].

FL protects client privacy in two main ways. Firstly, instead of transmitting the raw data, FL transmits only the model

Algorithm 1 FedAvg. The K clients are indexed by k ; C is the participation rate, B is the local minibatch size, E is the number of local epochs, and η is the learning rate.

Server executes:

- 1: initialize model parameter w_0
- 2: **for** each round $t = 1, 2, \dots, T$ **do**
- 3: $m \leftarrow \max(C \cdot K, 1)$
- 4: $S_t \leftarrow$ (random set of m clients)
- 5: **for** each client $k \in S_t$ **in parallel do**
- 6: $w_{t+1}^k \leftarrow$ ClientUpdate(k, w_t)
- 7: **end for**
- 8: $m_t \leftarrow \sum_{k \in S_t} n_k$
- 9: $w_{t+1} \leftarrow \sum_{k \in S_t} \frac{n_k}{m_t} w_{t+1}^k$
- 10: **end for**
- 11:
- 12: **ClientUpdate**(k, w): // Run on client k
- 13: $\mathcal{B} \leftarrow$ (split \mathcal{P}_k into batches of size B)
- 14: **for** each local epoch i from 1 to E **do**
- 15: **for** batch $b \in \mathcal{B}$ **do**
- 16: $w \leftarrow w - \eta \nabla \ell(w; b)$
- 17: **end for**
- 18: **end for**
- 19: return w to server

parameters, which are updated based on the local data of each client. By doing so, FL ensures that sensitive data remains on the client’s device and is not transmitted to the central server and other clients. Secondly, the model parameters uploaded to the server only reveal the distribution of local data, rather than individual data points. This approach helps to maintain privacy by obscuring the specific data points used to train the model.

FL can be equipped with differential privacy mechanisms [59; 91] to enhance privacy protection. As described in the last section, differential privacy is a technique that involves adding noise to data in order to obscure individual contributions while still maintaining overall data patterns. However, different from graph generation, where the noise is added to the data (e.g., node feature, edges, etc), in the context of FL, the noise is added to the uploaded and downloaded model parameters. This ensures that even if an attacker were to obtain the model parameters, they would not be able to accurately infer the raw data from the model parameter. By adding moderate noise to the parameters, the model’s accuracy may be slightly reduced, but the overall performance remains comparable to non-private models. In summary, by using differential privacy mechanisms, FL can achieve even better privacy protection by making it harder for attackers to identify the sensitive data contributed by individual clients.

3.1.2 Application of Graph Federated Learning

In this part, we introduce important applications of federated learning on graph data. Roughly, we survey three representative application scenarios: *graph-level FL*, *subgraph-level FL*, and *node-level FL*.

1. *Graph-level FL*: Each client has one or several graphs, while different graphs are isolated and independent. One typical application of graph-level FL is for drug discovery [31], where biochemical industries collaborate to train a graph neural network model predict-

ing the property of molecules. Each molecule is a graph with basic atoms as nodes and chemical bonds as edges.

2. *Subgraph-level FL*: Each client has one graph, while each graph is a subgraph of an underlying global graph. One representative application of subgraph-level FL is for financial transaction data [43]. Each FL client is a bank that keeps a graph encoding the information of its customers, where nodes are individual customers and edges are financial transaction records. While each bank holds its own graph, customers in one bank may have connections to customers in another bank, introducing cross-client edges. Thus, each bank’s own graph is a subgraph of an underlying global graph.
3. *Node-level FL*: Each client is a node of a graph, and edges are the pairwise relationships between clients, e.g., their distribution similarity or data dependency. One example is the smart city, where clients are traffic sensors deployed on the road and linked to geographically adjacent sensors. While clients form a graph, each client can make an intelligent decision based on the collected road conditions and nearby devices.

Figure 4 illustrates the three application scenarios above. Next, we investigate each application scenario in the following three subsections individually.

3.2 Graph-level FL

In this subsection, we investigate graph-level FL. Graph-level FL is a natural extension of traditional FL: while each client has one or several graphs, different graphs are isolated and independent. The goal of each client is to train a graph neural network (GNN) model for a variety of local tasks, e.g., node-level (e.g., node classification), link-level (e.g., edge prediction), or graph-level (e.g., graph classification).

One of the most representative applications of graph-level FL is drug discovery, where graphs are molecules with atoms as nodes and chemical bonds as edges. Each FL client can be a pharmaceutical corporation that owns molecule data. Multiple corporations collaborate to train better model for molecular property prediction.

The biggest challenge of graph-level FL is the non-identical distribution among different clients’ data. Since each client in FL collects their local data individually, their local datasets usually have a different distribution. For example, different pharmaceutical corporations may focus on different types of molecules. Such heterogeneity among clients’ data distributions introduces optimization challenges to FL. Moreover, when clients’ distribution is largely different, it might be harmful or even impossible to train one universal global model across all clients. More sophisticated techniques are required to achieve beneficial collaboration.

Next, we will introduce algorithms for graph-level FL in two parts: global federated learning and personalized federated learning. Since graph-level FL is a natural extension of traditional FL, we will cover both general FL algorithms and graph FL algorithms.

3.2.1 Global Federated Learning

Global federated learning (GFL) aims to train a *shared* global model for all clients. FedAvg [58] provides an initial solution for training GNNs with isolated graphs from multiple clients.

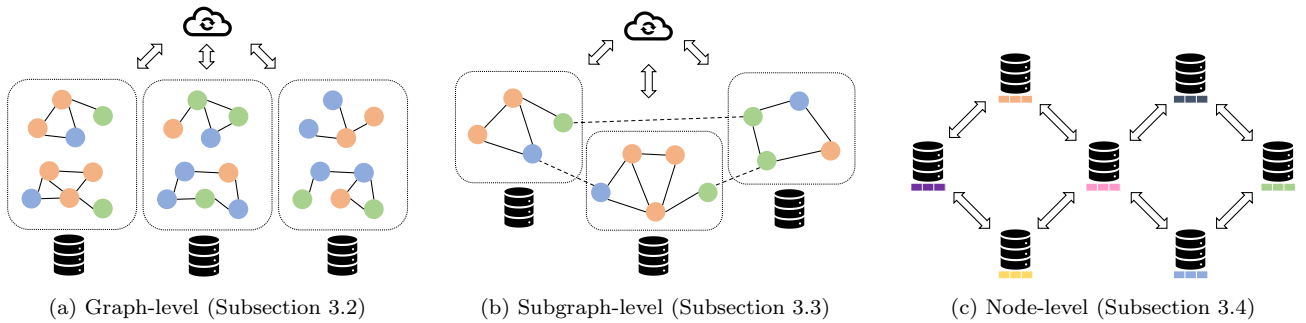


Figure 4: Three application scenarios of graph federated learning.

However, when clients have significantly different underlying distributions, FedAvg needs much more communication rounds for convergence to a satisfactory model, and may converge to a sub-optimal solution [58]. This phenomenon of worse convergence is usually explained by *weight divergence* [110], i.e, even with the same parameter initialization, the model parameters for different clients are substantially different after the first local stochastic gradient descent (SGD) step. With different model parameters, the mean of client gradients can be different from the gradient in centralized SGD, and introduce error to the model loss [84].

Data-sharing. To tackle the non-IID challenge to FL optimization, a simple but effective method is to share a small amount of data among clients. [110] first explore an association between the weight divergence and the non-IIDness of the data, and propose a method to share a small amount of data among the server and all clients. As a result, the accuracy can be increased by 30% for the CIFAR-10 dataset [40] with only 5% globally shared data. [102] further improves the privacy of this approach by sharing the average of local data points, instead of raw data. Specifically, each client uploads averaged data, receives averaged data from other clients, and performs Mixup [104] data augmentation locally to alleviate weight divergence. However, both methods require modification of the standard FL protocol and transmission of data. Another way to improve privacy is to share synthetic data generated by generative adversarial networks (GANs) [28], instead of the raw data. The synthetic data can be a collection of each client’s synthetic data generated with local GANs or generated with one global GAN trained in FL [67; 12]. However, it is unclear whether GAN can provide enough privacy, since it may memorize the training data [4].

Modifying local update. Another line of research works modifies the local update procedure to alleviate weight divergence without changing the communication protocol of FL. FedProx [47] adds a proximal term to the local objective to stabilize the training procedure. The proximal term is the squared L2 distance between the current global model and the local model, which prevents the local model from drifting too far from the global model. SCAFFOLD [38] estimates how local updates deviate from the global update, and it then corrects the local updates via variance reduction. Based on the intuition that the global model can learn better representation than local models, MOON [45] conducts contrastive learning at the model level, encouraging the agreement of representation learned by the local and global models.

3.2.2 Personalized Federated Learning

While the aforementioned algorithms can accelerate the model optimization for GFL, one model may not always be ideal for all participating clients [72]. Recently, personalized federated learning (PFL) has been proposed to tackle this challenge. PFL allows FL clients to collaboratively train machine learning models while each client can have different model parameters.

Clustered FL. In clustered FL, clients are partitioned into non-overlapping groups. Clients in the same group will share the same model, while clients from different groups can have different model parameters. In IFCA [27], k models are initialized and transmitted to all clients in each communication round, and each client picks the model with the smallest loss value to optimize. FedCluster [72] iteratively bipartition the clients based on their cosine similarity of gradients. GCFL [95] generalizes this idea to graph data, enabling collaborative training with graphs from different domains. Observing that the gradients of GNNs can be fluctuating, GCFL+ [95] uses a gradient sequence-based clustering mechanism to form more robust clusters.

Personalized Modules. Another prevalent way for PFL is personalized modules. In these works, the machine learning model is divided into two parts: the shared part and the personalized part. The key is to design a model structure suitable for personalization. For example, when a model is split into a feature extractor and classifier, FedPer [2] shares the feature extractor and personalizes the classifier, while LG-FedAvg [49] personalizes the feature extractor and shares the classifier. Similar techniques in used in FMTGL [54] and NGL-FedRep [80]. Moreover, PartialFed [75] can automatically select which layers to personalize and which layers to share. On graph data, [78] observe that while the feature information can be very different, some structural properties are shared by various domains, revealing the great potential for sharing structural information in FL. Inspired by this, they propose FedStar that trains a feature-structure decoupled GNN. The structural encoder is globally shared across all clients, while the feature-based knowledge is personalized.

Local Finetuning and Meta-Learning. Finetuning is widely used for PFL. In these works, a global model is first trained with all clients. The global model encodes the information of the population but may not adapt to each client’s own distribution. Therefore, each client locally finetunes the global model with a few steps of gradient descent. Besides vanilla finetuning, Per-FedAvg [13] combines FL with

MAML [15], an algorithm for meta-learning, to improve the performance of finetuning. Similarly, pFedMe [10] utilize Moreau Envelopes for personalization. It adds a proximal term to the local finetuning objective, and aims to find a local model near the global model, with just a few steps of gradient descent. GraphFL [83] applies a similar meta-learning framework on graph data, addressing the heterogeneity among graph data and handling new label domains with a few new labeled nodes.

Multi-task Learning. PFL is also studied within the framework of multi-task learning. MOCHA [73] uses a matrix to model the similarity among each pair of clients. Clients with similar distribution will be encouraged to have similar model parameters. FedGMTL [31] generalizes this idea to graph data. Similarly, SemiGraphFL [79] computes pairwise cosine similarity among clients’ hidden representations. As a result, clients with more similar data will have greater mutual influence. However, it requires the transmission of hidden representation. FedEM [57] assumes that each client’s distribution is a mixture of unknown underlying distributions and proposes FedEM, an EM-like algorithm for multi-task FL. Finally, FedFOMO [107] allows each client to have a different mixture weight of local models during the aggregation steps. It provides a flexible way for model aggregation.

Graph Structure Augmentation. In the previous works, graph structures are considered as ground truth. However, graphs can be noisy or incomplete, which can hurt the performance of GNNs. To tackle incomplete graph structures, FedGSL [109] optimizes the local client’s graph and GNN parameters simultaneously.

3.3 Subgraph-level FL

Similar to graph-level FL, each client in subgraph-level FL holds one graph. However, clients’ graphs are a subgraph of a latent large entire graph. In other words, there are cross-client edges in the entire graph, where the two nodes of these edges belong to different clients. The task is usually node-level, while the cross-client edges can contribute to the task.

One application of subgraph-level FL is financial fraud detection. Each FL client is a bank aiming to detect potential fraud with transaction data. Each bank keeps a graph of the information of its customers, where nodes are individual customers and edges are transaction records. While each bank holds its own graph, customers in one bank may have connections to customers in another bank, introducing edges across clients. These cross-client edges help to train better ML models.

The biggest challenge for subgraph-level FL is to handle cross-client edges. In GNNs, each node iteratively aggregates information from its neighboring nodes, which may be from other clients. However, during local updates in traditional FL, clients cannot get access to the data from other clients. Directly exchanging raw data among clients is prohibited due to privacy concerns. It is challenging to enable cross-client information exchange while preserving privacy. Moreover, when nodes’ identities are not shared across clients, the cross-client edges can be missing and stored in none of the clients. Even if we collect clients’ local subgraphs, we cannot reconstruct the global graph.

In this subsection, we will mainly focus on two scenarios. In the first part, we introduce algorithms when the hidden entire graph is given but stored separately in different

clients. In the second part, we consider a more challenging setting: the cross-client edges are missing, and we cannot simply concatenate local graphs to reconstruct the entire graph losslessly. We focus on how to generate these missing edges or missing neighbors for each node.

3.3.1 Cross-client Propagation

When the cross-client edges are available, the major challenge is to enable cross-client information propagation without leaking raw data. FedGraph [8] designs a novel cross-client convolution operation to avoid sharing raw data across clients. It avoids exchanging representations in the first GCN layer. Similarly, FedPNS [11] control the number of neighbor sampling to reduce communication costs. FedCog [43] proposes graph decoupling operation, splitting local graph to internal graph and border graph. The graph convolution is accordingly divided into two sequential steps: internal propagation and border propagation. In this process, each client sends the intermediate representation of internal nodes to other clients. Considering that directly exchanging feature representations between clients can leak private information. In user-item graphs, FedPerGNN [92] design a privacy-preserving user-item graph expansion protocol. Clients upload encrypted item IDs to the trusted server, and the server matches the ciphertexts of item IDs to find clients with overlapping item IDs. DP-FedRec [65] uses private set intersection to exchange the edges information between clients and applies differential privacy techniques to further protect privacy. Different from the above methods, FedGCN [98] does not rely on communication between clients. Instead, it transmits all the information needed to train a GCN between the server and each client, only once before the training. Moreover, each node at a given client only needs to know the accumulated information about the node’s neighbors, which reduces possible privacy leakage.

3.3.2 Missing Neighbors

For some applications, the cross-client edges can be missing or not stored in any clients. Notice that although each client also holds a disjoint graph in graph-level FL, graph-level FL and subgraph-level FL with missing neighbors are substantially different. For graph-level FL, there are essentially no cross-client edges. For example, there are no chemical bonds between two molecules from different corporations’ datasets. However, for subgraph-level FL, the cross-client edges exist, but are missing in certain applications. We may get suboptimal GNN models if ignoring the existence of cross-client edges. Therefore, the major challenge is to reconstruct these missing edges, or reconstruct missing neighbors for each node.

FedSAGE [106] first defines the missing neighbors’ challenge, and proposes a method to generate pseudo neighbors for each node. It uses existing subgraphs to train a neighbors generator and generate one-hop neighbors for each client to mend the graph. Since missing neighbors are generated locally, no feature exchange is required between clients after the local subgraphs are mended. However, the training of neighbor generators requires cross-client hidden representation exchanges. Similarly, FedNI [63] uses a graph GAN model to generate missing nodes and edges.

3.4 Node-level FL

The final application scenario of graph federated learning is

Table 2: Repositories for Graph Federated Learning

Scenario	Name	Description	Link
Graph-level	FedProx [47]	A general GFL algorithm with modified local update	Github
	IFCA [27]	A general clustered FL algorithm	Github
	GCFL [95]	A graph-specific clustered FL algorithm	Github
	LG-FedAvg [49]	A general PFL algorithm with personalized modules	Github
	FedStar [78]	A graph-specific PFL algorithm with personalized modules	Github
	pFedMe [10]	A general PFL algorithm based on meta-learning	Github
	GraphFL [83]	A graph-specific PFL algorithm based on meta-learning	Github
	FedFOMO [107]	A general PFL algorithm based on multi-task learning	Github
Subgraph-level	FedGCN [98]	An FL algorithm with one-shot cross-client propagation	Github
	FedSAGE [106]	An FL algorithm with missing neighbors generation	Github
Node-level	SpreadGNN [31]	A serverless PFL algorithm	Github
	FedGS [89]	An FL algorithm with graph as distribution similarities	Github
	SFL [9]	An GL with pre-defined graph for server aggregation	Github
Others	TensorFlow Federated	A framework for implementing federated learning	Github
	FedLab [103]	A Flexible Federated Learning Framework	Github
	PFL-Non-IID	Reproduction of popular PFL algorithms	Github
	FedGraphNN [30]	FedGraphNN: A Federated Learning System and Benchmark for Graph Neural Networks	Github
	FederatedScope-GNN [90]	A unified, comprehensive and efficient package for Federated Graph Learning	Github

node-level. Different from the aforementioned two scenarios, each client in node-level FL can hold any type of data, not restricted to graphs. Instead, the clients themselves are nodes in a graph, while the edges are their pairwise relationship of communication or distribution similarity.

One typical application of node-level FL is the Internet of Things (IoT) devices in a smart building [68]. Due to bandwidth constraints, it can be costly for each IoT device to communicate with the central server. However, IoT devices in the same local area network can communicate very efficiently. As a result, IoT devices form a graph with pairwise communication availability as edges. Another application is for the smart city [89], where clients are traffic sensors deployed on the road and linked to geographically adjacent sensors. Each device can collect data and make the real-time decision without waiting for the response of cloud servers. Each device needs to make an intelligent decision based on the collected road conditions and nearby devices.

In this subsection, we will first introduce algorithms where the graph models communication constraints among clients. In these works, there is no central server, and clients can only exchange information along edges. Then, we will introduce algorithms where the graph models the relationship between clients’ distributions. In these works, although a central server is available, the graph among clients models distributional similarity or dependency among clients, potentially contributes to the model performance.

3.4.1 Graph as Communication Network

Traditional FL relies on a central server to enable communication among clients. Each client trusts the central server and uploads their model update to the server. However, in many scenarios, a trusted central server may not exist. Even when a central server exists, it may be expensive for clients to communicate with the server. Therefore, serverless FL (a.k.a. peer-to-peer FL) has been studied to relieve communication constraints.

The standard solution for serverless FL is fully decentral-

ized FL [42; 41], where each client only averages its model parameter with its neighbors. D-FedGNN [62] uses these techniques to train GNN models. SpreadGNN [31] generalizes this framework to personalized FL, where each client has non-IID data and a different label space.

3.4.2 Graph as Distribution Similarities

When the central server is available, a graph of clients may still be beneficial when it models distributional relationships among clients. When edges link clients with highly similar distributions, parameter sharing along edges can potentially improve the model performance for both clients. When edges link clients with data dependency, information exchange along edges can even provide additional features for inference.

FedGS [89] models the data correlations of clients with a data-distribution-dependency graph, and improves the unbiasedness of the client sampling process. Meanwhile, SFL [9] assumes a pre-defined client relation graph stored on the server, and the client-centric model aggregation is conducted along the relation graph’s structure. GraphFL [108] considers client-side information to encourage similar clients to have similar models. BiG-Fed [96] applies graph convolution on the client graph, so each client’s prediction can benefit from its neighbors with highly correlated data. Finally, [86] designs a client sampling technique considers both communication cost and distribution similarity.

Finally, we summarize the official implementation of FL algorithms and useful repositories in Table 2.

3.5 Challenges and Future Opportunities

In this part, we present several limitations in current works and provide open problems for future research.

3.5.1 Model Heterogeneity for Graph-Level FL

In previous works of graph-level FL, although each FL client usually has different data distribution it is usually assumed that the model architecture is shared across all clients. How-

ever, the optimal architecture for different clients can be different. For example, a well-known issue in GNNs is the over-smoothing problem. When the number of graph convolutional layers is higher than the diameter of the graph, GNN models may learn similar representations for all nodes in the graph, which harms the model performance. When each FL clients hold a substantially different size of graphs, it is highly likely that the optimal depth of the GNN model is different for them.

3.5.2 Avoiding Cross-Client Transmission for Subgraph-Level FL

Most of the previous subgraph-level FL algorithms highly rely on direct information exchange along cross-client edges. While such operations are natural variants of graph convolution, such operations also raise privacy concerns. Moreover, different from traditional FL where each client downloads aggregated model parameters that reveal the population, feature exchange along the edges can expose information about individuals. It would be beneficial if the cross-client transmission can be avoided without greatly degrading the model.

4. ENVISIONING

In this section, we analyze the current developments and limitations of privacy-preserving graph machine learning, and explain the necessity of combining them. In addition, we identify a number of unsolved research directions that could be addressed to improve the privacy of graph machine learning systems.

4.1 Limitation of Current Techniques

In the previous two sections, we introduced privacy-preserving graph data generation and computation, respectively. However, both techniques have their own limitations.

- For privacy-preserving graph generation, while it can provide good privacy protection for graph data, it also has a significant drawback on model utility. The privacy-preserving techniques applied during data generation are not designed for specific machine learning tasks and may influence the utility of the resulting model. For example, consider a graph with four nodes a , b , c , and d . The nodes a and b have a positive label, while c and d have a negative label. Switching the edges from $(a, b), (c, d)$ to $(a, c), (b, d)$ does not change the degree distribution of the graph, but it changes the graph from a homophilous graph to a heterophilous graph, i.e., edges are more likely to link two nodes with different labels. This change can harm the performance of many GNN models, which are designed to work well with homogeneous graphs [50]. It is important to consider the downstream machine learning tasks when designing privacy-preserving techniques for graph data.
- For privacy-preserving graph computation, while FL can avoid the transmission of raw data, it has been shown that transmitting raw model parameters or gradients may not provide enough privacy, as attackers can use the gradient or model update to reconstruct private data [114; 26]. Moreover, many subgraph-level and node-level federated learning algorithms require

the transmission of hidden representations, which can also leak private information. Therefore, protecting the raw data from being reconstructed is essential to federated learning systems.

4.2 Combination of Privacy-Preserving Graph Data Generation and Computation

To address the limitations of current privacy-preserving techniques, it is essential to combine privacy graph data generation with the graph federated learning frameworks, as shown in Figure 5. This approach can provide an effective solution to the privacy preservation issues of graph machine learning models.

Specifically, the generated synthetic data is used instead of the real data during the training process. This means that even if the transmitted information is decrypted, it is just from the generated synthetic data and not the real data. The synthetic data can be generated in such a way that it preserves the statistical properties of the original data while ensuring privacy preservation. This can be achieved using various techniques, including differential privacy, homomorphic encryption, and secure multi-party computation.

The combination of privacy graph data generation and graph federated learning frameworks has several benefits. First, it ensures privacy preservation during the training process by using synthetic data. Second, it enables the transfer of graph machine learning model parameters rather than embedding vectors or other information. This can improve the accuracy and efficiency of the model. Finally, it provides a robust defense against privacy attacks and reverse-engineering, as the transmitted information is just from the generated synthetic data and not the real data.

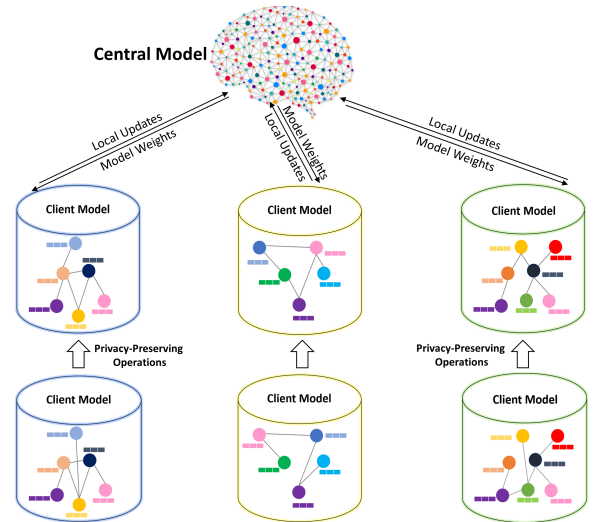


Figure 5: Privacy-preserving Graph Data with Privacy-preserving Computation.

4.3 Future Directions

Combining privacy-preserving data generation and computation is a promising approach to protect individual privacy while maintaining model utility in machine learning. However, it also poses several challenges and possible future directions.

4.3.1 Distribution of Privacy Budget

When combining privacy-preserving data generation with computation, noises are added to both raw data and model parameters. However, it is still unclear how to distribute the privacy budget between data generation and computation in a way that optimizes the privacy-utility trade-off. In this approach, noises are added to the graph data during data generation and to the model parameters during data computation (i.e., federated learning), which results in an overall reduction in accuracy. However, while the privacy analysis for data generation is directly defined on the data space, the privacy analysis for federated learning requires transforming the change on parameter space back to data space. Such transformation requires estimating the sensitivity of a machine learning algorithm (i.e., how the change of a data point affects the learned parameters), which is only loosely bounded in current works [59; 91]. A more precise analysis of privacy is required to better understand the impact of privacy budget allocation on the overall privacy-utility trade-off.

4.3.2 Parameter Information Disentanglement

Another future challenge when combining privacy-preserving data generation and computation is the disentanglement of task-relevant and task-irrelevant information. Currently, the noise added to the model parameters is isotropic, meaning that task-relevant and task-irrelevant information are equally protected. However, not all information is equally important for model utility. If we can identify which information has a significant influence on model performance, we can distribute more privacy budget to this information while allocating less privacy budget to task-irrelevant information. This can result in a better privacy-utility trade-off. Disentangling task-relevant and task-irrelevant information would require a more sophisticated analysis of model architecture and data characteristics to determine which features contribute most to model performance.

5. CONCLUSION

In this paper, we review the research for privacy-preserving techniques for graph machine learning from the data to the computation, considering the situation where the data need to be shared or are banned from being transmitted. To be specific, for privacy-preserving graph data generation techniques, we analyze the forceful attackers first and then introduce how corresponding protection methods are proposed to defend attackers. For the privacy graph data computation, we circle around the federated learning setting and discuss how the general federated learning framework applied to graph data and what the potential challenges originated from non-IIDness, and how the nascent research works address them. In the end, we analyze the current limitation and propose several promising research directions.

6. ACKNOWLEDGEMENTS

This work is supported by the National Science Foundation (1947203, 2117902, 2137468, 1947135, 2134079, and 1939725), the U.S. Department of Homeland Security (2017-ST-061-QA0001, 17STQAC00001-06-00, and 17STQAC00001-03-03), DARPA (HR001121C0165), NIFA (2020-67021-32799), and ARO (W911NF2110088). The views and conclusions are those of the authors and should not be interpreted as

representing the official policies of the funding agencies or the government.

7. REFERENCES

- [1] M. Abadi, A. Chu, I. J. Goodfellow, H. B. McMahan, I. Mironov, K. Talwar, and L. Zhang. Deep learning with differential privacy. In *SIGSAC 2016*, 2016.
- [2] M. G. Arivazhagan, V. Aggarwal, A. K. Singh, and S. Choudhary. Federated learning with personalization layers. *CoRR*, abs/1912.00818, 2019.
- [3] R. Arora and J. Upadhyay. On differentially private graph sparsification and applications. In H. M. Wallach, H. Larochelle, A. Beygelzimer, F. d’Alché-Buc, E. B. Fox, and R. Garnett, editors, *Advances in Neural Information Processing Systems 32: Annual Conference on Neural Information Processing Systems 2019, NeurIPS 2019, December 8-14, 2019, Vancouver, BC, Canada*, pages 13378–13389, 2019.
- [4] S. Arora, A. Risteski, and Y. Zhang. Do gans learn the distribution? some theory and empirics. In *6th International Conference on Learning Representations, ICLR 2018, Vancouver, BC, Canada, April 30 - May 3, 2018, Conference Track Proceedings*. OpenReview.net, 2018.
- [5] L. Backstrom, C. Dwork, and J. M. Kleinberg. Wherefore art thou r3579x?: anonymized social networks, hidden patterns, and structural steganography. In C. L. Williamson, M. E. Zurko, P. F. Patel-Schneider, and P. J. Shenoy, editors, *Proceedings of the 16th International Conference on World Wide Web, WWW 2007, Banff, Alberta, Canada, May 8-12, 2007*, pages 181–190. ACM, 2007.
- [6] J. Blocki, A. Blum, A. Datta, and O. Sheffet. Differentially private data analysis of social networks via restricted sensitivity. In R. D. Kleinberg, editor, *Innovations in Theoretical Computer Science, ITCS ’13, Berkeley, CA, USA, January 9-12, 2013*, pages 87–96. ACM, 2013.
- [7] C. Chen, Y. Wu, Q. Dai, H. Zhou, M. Xu, S. Yang, X. Han, and Y. Yu. A survey on graph neural networks and graph transformers in computer vision: A task-oriented perspective. *CoRR*, abs/2209.13232, 2022.
- [8] F. Chen, P. Li, T. Miyazaki, and C. Wu. Fedgraph: Federated graph learning with intelligent sampling. *IEEE Trans. Parallel Distributed Syst.*, 33(8):1775–1786, 2022.
- [9] F. Chen, G. Long, Z. Wu, T. Zhou, and J. Jiang. Personalized federated learning with a graph. In L. D. Raedt, editor, *Proceedings of the Thirty-First International Joint Conference on Artificial Intelligence, IJCAI 2022, Vienna, Austria, 23-29 July 2022*, pages 2575–2582. ijcai.org, 2022.
- [10] C. T. Dinh, N. H. Tran, and T. D. Nguyen. Personalized federated learning with moreau envelopes. In H. Larochelle, M. Ranzato, R. Hadsell, M. Balcan, and H. Lin, editors, *Advances in Neural Information*

- Processing Systems 33: Annual Conference on Neural Information Processing Systems 2020, NeurIPS 2020, December 6-12, 2020, virtual*, 2020.
- [11] B. Du and C. Wu. Federated graph learning with periodic neighbour sampling. In *2022 IEEE/ACM 30th International Symposium on Quality of Service (IWQoS)*, pages 1–10, 2022.
- [12] E. Ekblom, E. L. Zec, and O. Mogren. EFFGAN: ensembles of fine-tuned federated gans. In S. Tsumoto, Y. Ohsawa, L. Chen, D. V. den Poel, X. Hu, Y. Motomura, T. Takagi, L. Wu, Y. Xie, A. Abe, and V. Raghavan, editors, *IEEE International Conference on Big Data, Big Data 2022, Osaka, Japan, December 17-20, 2022*, pages 884–892. IEEE, 2022.
- [13] A. Fallah, A. Mokhtari, and A. E. Ozdaglar. Personalized federated learning with theoretical guarantees: A model-agnostic meta-learning approach. In H. Larochelle, M. Ranzato, R. Hadsell, M. Balcan, and H. Lin, editors, *Advances in Neural Information Processing Systems 33: Annual Conference on Neural Information Processing Systems 2020, NeurIPS 2020, December 6-12, 2020, virtual*, 2020.
- [14] H. Fichtenberger, M. Henzinger, and W. Ost. Differentially private algorithms for graphs under continual observation. In P. Mutzel, R. Pagh, and G. Herman, editors, *29th Annual European Symposium on Algorithms, ESA 2021, September 6-8, 2021, Lisbon, Portugal (Virtual Conference)*, volume 204 of *LIPICs*, pages 42:1–42:16. Schloss Dagstuhl - Leibniz-Zentrum für Informatik, 2021.
- [15] C. Finn, P. Abbeel, and S. Levine. Model-agnostic meta-learning for fast adaptation of deep networks. In D. Precup and Y. W. Teh, editors, *Proceedings of the 34th International Conference on Machine Learning, ICML 2017, Sydney, NSW, Australia, 6-11 August 2017*, volume 70 of *Proceedings of Machine Learning Research*, pages 1126–1135. PMLR, 2017.
- [16] D. Fu, L. Fang, R. Maciejewski, V. I. Torvik, and J. He. Meta-learned metrics over multi-evolution temporal graphs. In *KDD '22: The 28th ACM SIGKDD Conference on Knowledge Discovery and Data Mining, Washington, DC, USA, August 14 - 18, 2022*, 2022.
- [17] D. Fu and J. He. SDG: A simplified and dynamic graph neural network. In *SIGIR '21: The 44th International ACM SIGIR Conference on Research and Development in Information Retrieval, Virtual Event, Canada, July 11-15, 2021*, 2021.
- [18] D. Fu and J. He. DPPIN: A biological repository of dynamic protein-protein interaction network data. In *IEEE International Conference on Big Data, Big Data 2022, Osaka, Japan, December 17-20, 2022*, 2022.
- [19] D. Fu and J. He. Natural and artificial dynamics in graphs: Concept, progress, and future. *Frontiers in Big Data*, 5, 2022.
- [20] D. Fu, J. He, H. Tong, and R. Maciejewski. Privacy-preserving graph analytics: secure generation and federated learning. *arXiv preprint arXiv:2207.00048*, 2022.
- [21] D. Fu, Z. Xu, B. Li, H. Tong, and J. He. A view-adversarial framework for multi-view network embedding. In *CIKM '20: The 29th ACM International Conference on Information and Knowledge Management, Virtual Event, Ireland, October 19-23, 2020*, 2020.
- [22] D. Fu, D. Zhou, and J. He. Local motif clustering on time-evolving graphs. In *KDD '20: The 26th ACM SIGKDD Conference on Knowledge Discovery and Data Mining, Virtual Event, CA, USA, August 23-27, 2020*, 2020.
- [23] D. Fu, D. Zhou, R. Maciejewski, A. Croitoru, M. Boyd, and J. He. Fairness-aware clique-preserving spectral clustering of temporal graphs. In *Proceedings of the ACM Web Conference 2023, WWW 2023, Austin, TX, USA, 30 April 2023 - 4 May 2023*, 2023.
- [24] X. Fu, B. Zhang, Y. Dong, C. Chen, and J. Li. Federated graph machine learning: A survey of concepts, techniques, and applications. *SIGKDD Explor.*, 24(2):32–47, 2022.
- [25] T. Gaudelot, B. Day, A. R. Jamasb, J. Soman, C. Regep, G. Liu, J. B. R. Hayter, R. Vickers, C. Roberts, J. Tang, D. Roblin, T. L. Blundell, M. M. Bronstein, and J. P. Taylor-King. Utilizing graph machine learning within drug discovery and development. *Briefings Bioinform.*, 22(6), 2021.
- [26] J. Geiping, H. Bauermeister, H. Dröge, and M. Moeller. Inverting gradients - how easy is it to break privacy in federated learning? In H. Larochelle, M. Ranzato, R. Hadsell, M. Balcan, and H. Lin, editors, *Advances in Neural Information Processing Systems 33: Annual Conference on Neural Information Processing Systems 2020, NeurIPS 2020, December 6-12, 2020, virtual*, 2020.
- [27] A. Ghosh, J. Chung, D. Yin, and K. Ramchandran. An efficient framework for clustered federated learning. In H. Larochelle, M. Ranzato, R. Hadsell, M. Balcan, and H. Lin, editors, *Advances in Neural Information Processing Systems 33: Annual Conference on Neural Information Processing Systems 2020, NeurIPS 2020, December 6-12, 2020, virtual*, 2020.
- [28] I. J. Goodfellow, J. Pouget-Abadie, M. Mirza, B. Xu, D. Warde-Farley, S. Ozair, A. C. Courville, and Y. Bengio. Generative adversarial nets. In Z. Ghahramani, M. Welling, C. Cortes, N. D. Lawrence, and K. Q. Weinberger, editors, *Advances in Neural Information Processing Systems 27: Annual Conference on Neural Information Processing Systems 2014, December 8-13 2014, Montreal, Quebec, Canada*, pages 2672–2680, 2014.
- [29] M. Hay, G. Miklau, D. D. Jensen, D. F. Towsley, and P. Weis. Resisting structural re-identification in anonymized social networks. *Proc. VLDB Endow.*, 1(1):102–114, 2008.
- [30] C. He, K. Balasubramanian, E. Ceyani, Y. Rong, P. Zhao, J. Huang, M. Annavaram, and S. Avestimehr. Fedgraphnn: A federated learning system and benchmark for graph neural networks. *CoRR*, abs/2104.07145, 2021.

- [31] C. He, E. Ceyani, K. Balasubramanian, M. Annavaram, and S. Avestimehr. Spreadgnn: Decentralized multi-task federated learning for graph neural networks on molecular data. In *Thirty-Sixth AAAI Conference on Artificial Intelligence, AAAI 2022, Thirty-Fourth Conference on Innovative Applications of Artificial Intelligence, IAAI 2022, The Twelfth Symposium on Educational Advances in Artificial Intelligence, EAAI 2022 Virtual Event, February 22 - March 1, 2022*, pages 6865–6873. AAAI Press, 2022.
- [32] A. Hoang, B. Carminati, and E. Ferrari. Cluster-based anonymization of knowledge graphs. In M. Conti, J. Zhou, E. Casalichio, and A. Spognardi, editors, *Applied Cryptography and Network Security - 18th International Conference, ACNS 2020, Rome, Italy, October 19-22, 2020, Proceedings, Part II*, volume 12147 of *Lecture Notes in Computer Science*, pages 104–123. Springer, 2020.
- [33] A. Hoang, B. Carminati, and E. Ferrari. Privacy-preserving sequential publishing of knowledge graphs. In *37th IEEE International Conference on Data Engineering, ICDE 2021, Chania, Greece, April 19-22, 2021*, pages 2021–2026. IEEE, 2021.
- [34] A.-T. Hoang, B. Carminati, and E. Ferrari. Time-aware anonymization of knowledge graphs. *ACM Transactions on Privacy and Security*, 2022.
- [35] S. Ji, S. Pan, E. Cambria, P. Marttinen, and P. S. Yu. A survey on knowledge graphs: Representation, acquisition, and applications. *IEEE Trans. Neural Networks Learn. Syst.*, 33(2):494–514, 2022.
- [36] H. Jiang, J. Pei, D. Yu, J. Yu, B. Gong, and X. Cheng. Applications of differential privacy in social network analysis: A survey. *IEEE Trans. Knowl. Data Eng.*, 35(1):108–127, 2023.
- [37] P. Kairouz, H. B. McMahan, B. Avent, A. Bellet, M. Bennis, A. N. Bhagoji, K. A. Bonawitz, Z. Charles, G. Cormode, R. Cummings, R. G. L. D’Oliveira, H. Eichner, S. E. Rouayheb, D. Evans, J. Gardner, Z. Garrett, A. Gascón, B. Ghazi, P. B. Gibbons, M. Gruteser, Z. Harchaoui, C. He, L. He, Z. Huo, B. Hutchinson, J. Hsu, M. Jaggi, T. Javidi, G. Joshi, M. Khodak, J. Konečný, A. Korolova, F. Koushanfar, S. Koyejo, T. Lepoint, Y. Liu, P. Mittal, M. Mohri, R. Nock, A. Özgür, R. Pagh, H. Qi, D. Ramage, R. Raskar, M. Raykova, D. Song, W. Song, S. U. Stich, Z. Sun, A. T. Suresh, F. Tramèr, P. Vepakomma, J. Wang, L. Xiong, Z. Xu, Q. Yang, F. X. Yu, H. Yu, and S. Zhao. Advances and open problems in federated learning. *Found. Trends Mach. Learn.*, 14(1-2):1–210, 2021.
- [38] S. P. Karimireddy, S. Kale, M. Mohri, S. J. Reddi, S. U. Stich, and A. T. Suresh. SCAFFOLD: stochastic controlled averaging for federated learning. In *Proceedings of the 37th International Conference on Machine Learning, ICML 2020, 13-18 July 2020, Virtual Event*, volume 119 of *Proceedings of Machine Learning Research*, pages 5132–5143. PMLR, 2020.
- [39] S. P. Kasiviswanathan, K. Nissim, S. Raskhodnikova, and A. D. Smith. Analyzing graphs with node differential privacy. In A. Sahai, editor, *Theory of Cryptography - 10th Theory of Cryptography Conference, TCC 2013, Tokyo, Japan, March 3-6, 2013. Proceedings*, volume 7785 of *Lecture Notes in Computer Science*, pages 457–476. Springer, 2013.
- [40] A. Krizhevsky, G. Hinton, et al. Learning multiple layers of features from tiny images. 2009.
- [41] A. Lalitha, O. C. Kilinc, T. Javidi, and F. Koushanfar. Peer-to-peer federated learning on graphs. *CoRR*, abs/1901.11173, 2019.
- [42] A. Lalitha, S. Shekhar, T. Javidi, and F. Koushanfar. Fully decentralized federated learning. In *Third workshop on bayesian deep learning (NeurIPS)*, volume 2, 2018.
- [43] R. Lei, P. Wang, J. Zhao, L. Lan, J. Tao, C. Deng, J. Feng, X. Wang, and X. Guan. Federated learning over coupled graphs. *IEEE Trans. Parallel Distributed Syst.*, 34(4):1159–1172, 2023.
- [44] N. Li, T. Li, and S. Venkatasubramanian. t-closeness: Privacy beyond k-anonymity and l-diversity. In R. Chirkova, A. Dogac, M. T. Özsu, and T. K. Sellis, editors, *Proceedings of the 23rd International Conference on Data Engineering, ICDE 2007, The Marmara Hotel, Istanbul, Turkey, April 15-20, 2007*, pages 106–115. IEEE Computer Society, 2007.
- [45] Q. Li, B. He, and D. Song. Model-contrastive federated learning. In *IEEE Conference on Computer Vision and Pattern Recognition, CVPR 2021, virtual, June 19-25, 2021*, pages 10713–10722. Computer Vision Foundation / IEEE, 2021.
- [46] T. Li, A. K. Sahu, A. Talwalkar, and V. Smith. Federated learning: Challenges, methods, and future directions. *IEEE Signal Process. Mag.*, 37(3):50–60, 2020.
- [47] T. Li, A. K. Sahu, M. Zaheer, M. Sanjabi, A. Talwalkar, and V. Smith. Federated optimization in heterogeneous networks. In I. S. Dhillon, D. S. Papailiopoulos, and V. Sze, editors, *Proceedings of Machine Learning and Systems 2020, MLSys 2020, Austin, TX, USA, March 2-4, 2020*. mlsys.org, 2020.
- [48] Z. Li, D. Fu, and J. He. Everything evolves in personalized pagerank. In *Proceedings of the ACM Web Conference 2023, WWW 2023, Austin, TX, USA, 30 April 2023 - 4 May 2023*, 2023.
- [49] P. P. Liang, T. Liu, Z. Liu, R. Salakhutdinov, and L. Morency. Think locally, act globally: Federated learning with local and global representations. *CoRR*, abs/2001.01523, 2020.
- [50] D. Lim, X. Li, F. Hohne, and S. Lim. New benchmarks for learning on non-homophilous graphs. *CoRR*, abs/2104.01404, 2021.
- [51] K. Liu, K. Das, T. Grandison, and H. Kargupta. Privacy-preserving data analysis on graphs and social

- networks. In H. Kargupta, J. Han, P. S. Yu, R. Motwani, and V. Kumar, editors, *Next Generation of Data Mining*, Chapman and Hall / CRC Data Mining and Knowledge Discovery Series. CRC Press / Chapman and Hall / Taylor & Francis, 2008.
- [52] K. Liu and E. Terzi. Towards identity anonymization on graphs. In J. T. Wang, editor, *Proceedings of the ACM SIGMOD International Conference on Management of Data, SIGMOD 2008, Vancouver, BC, Canada, June 10-12, 2008*, pages 93–106. ACM, 2008.
- [53] R. Liu and H. Yu. Federated graph neural networks: Overview, techniques and challenges. *CoRR*, abs/2202.07256, 2022.
- [54] Y. Liu, D. Han, J. Zhang, H. Zhu, M. Xu, and W. Chen. Federated multi-task graph learning. *ACM Trans. Intell. Syst. Technol.*, 13(5), jun 2022.
- [55] X. Ma, J. Wu, S. Xue, J. Yang, Q. Z. Sheng, and H. Xiong. A comprehensive survey on graph anomaly detection with deep learning. *CoRR*, abs/2106.07178, 2021.
- [56] A. Machanavajjhala, D. Kifer, J. Gehrke, and M. Venkitasubramaniam. L -diversity: Privacy beyond k -anonymity. *ACM Trans. Knowl. Discov. Data*, 1(1):3, 2007.
- [57] O. Marfoq, G. Neglia, A. Bellet, L. Kameni, and R. Vidal. Federated multi-task learning under a mixture of distributions. In M. Ranzato, A. Beygelzimer, Y. N. Dauphin, P. Liang, and J. W. Vaughan, editors, *Advances in Neural Information Processing Systems 34: Annual Conference on Neural Information Processing Systems 2021, NeurIPS 2021, December 6-14, 2021, virtual*, pages 15434–15447, 2021.
- [58] B. McMahan, E. Moore, D. Ramage, S. Hampson, and B. A. y Arcas. Communication-efficient learning of deep networks from decentralized data. In *Proceedings of the 20th International Conference on Artificial Intelligence and Statistics, AISTATS 2017, 20-22 April 2017, Fort Lauderdale, FL, USA*, volume 54 of *Proceedings of Machine Learning Research*, pages 1273–1282. PMLR, 2017.
- [59] H. B. McMahan, D. Ramage, K. Talwar, and L. Zhang. Learning differentially private recurrent language models. In *6th International Conference on Learning Representations, ICLR 2018, Vancouver, BC, Canada, April 30 - May 3, 2018, Conference Track Proceedings*. OpenReview.net, 2018.
- [60] H. H. Nguyen, A. Imine, and M. Rusinowitch. Anonymizing social graphs via uncertainty semantics. In *CCS 2015*, 2015.
- [61] J. Pearl. *Probabilistic reasoning in intelligent systems: networks of plausible inference*. Morgan kaufmann, 1988.
- [62] Y. Pei, R. Mao, Y. Liu, C. Chen, S. Xu, F. Qiang, and B. E. Tech. Decentralized federated graph neural networks. In *International Workshop on Federated and Transfer Learning for Data Sparsity and Confidentiality in Conjunction with IJCAI*, 2021.
- [63] L. Peng, N. Wang, N. Dvornek, X. Zhu, and X. Li. Fedni: Federated graph learning with network inpainting for population-based disease prediction. *IEEE Transactions on Medical Imaging*, pages 1–1, 2022.
- [64] Z. Qin, T. Yu, Y. Yang, I. Khalil, X. Xiao, and K. Ren. Generating synthetic decentralized social graphs with local differential privacy. In *CCS 2017*, 2017.
- [65] Y. Qiu, C. Huang, J. Wang, Z. Huang, and J. Xiao. A privacy-preserving subgraph-level federated graph neural network via differential privacy. In G. Memmi, B. Yang, L. Kong, T. Zhang, and M. Qiu, editors, *Knowledge Science, Engineering and Management*, pages 165–177, Cham, 2022. Springer International Publishing.
- [66] S. Raskhodnikova and A. D. Smith. Lipschitz extensions for node-private graph statistics and the generalized exponential mechanism. In I. Dinur, editor, *IEEE 57th Annual Symposium on Foundations of Computer Science, FOCS 2016, 9-11 October 2016, Hyatt Regency, New Brunswick, New Jersey, USA*, pages 495–504. IEEE Computer Society, 2016.
- [67] M. Rasouli, T. Sun, and R. Rajagopal. Fedgan: Federated generative adversarial networks for distributed data. *CoRR*, abs/2006.07228, 2020.
- [68] A. Rasti-Meymandi, S. M. Sheikholeslami, J. Abouei, and K. N. Plataniotis. Graph federated learning for ciot devices in smart home applications. *IEEE Internet of Things Journal*, 10(8):7062–7079, 2023.
- [69] L. Rossi, M. Musolesi, and A. Torsello. On the k -anonymization of time-varying and multi-layer social graphs. In M. Cha, C. Mascolo, and C. Sandvig, editors, *Proceedings of the Ninth International Conference on Web and Social Media, ICWSM 2015, University of Oxford, Oxford, UK, May 26-29, 2015*, pages 377–386. AAAI Press, 2015.
- [70] A. Sala, X. Zhao, C. Wilson, H. Zheng, and B. Y. Zhao. Sharing graphs using differentially private graph models. In P. Thiran and W. Willinger, editors, *Proceedings of the 11th ACM SIGCOMM Internet Measurement Conference, IMC '11, Berlin, Germany, November 2-, 2011*, pages 81–98. ACM, 2011.
- [71] P. Samarati. Protecting respondents' identities in microdata release. *IEEE Trans. Knowl. Data Eng.*, 13(6):1010–1027, 2001.
- [72] F. Sattler, K. Müller, and W. Samek. Clustered federated learning: Model-agnostic distributed multitask optimization under privacy constraints. *IEEE Trans. Neural Networks Learn. Syst.*, 32(8):3710–3722, 2021.
- [73] V. Smith, C. Chiang, M. Sanjabi, and A. Talwalkar. Federated multi-task learning. In I. Guyon, U. von Luxburg, S. Bengio, H. M. Wallach, R. Fergus, S. V. N. Vishwanathan, and R. Garnett, editors, *Advances in Neural Information Processing Systems 30: Annual Conference on Neural Information Processing Systems 2017, December 4-9, 2017, Long Beach, CA, USA*, pages 4424–4434, 2017.

- [74] S. Song, S. Little, S. Mehta, S. A. Vinterbo, and K. Chaudhuri. Differentially private continual release of graph statistics. *CoRR*, abs/1809.02575, 2018.
- [75] B. Sun, H. Huo, Y. Yang, and B. Bai. Partialfed: Cross-domain personalized federated learning via partial initialization. In M. Ranzato, A. Beygelzimer, Y. N. Dauphin, P. Liang, and J. W. Vaughan, editors, *Advances in Neural Information Processing Systems 34: Annual Conference on Neural Information Processing Systems 2021, NeurIPS 2021, December 6-14, 2021, virtual*, pages 23309–23320, 2021.
- [76] Y. Sun and J. Han. Mining heterogeneous information networks: a structural analysis approach. *SIGKDD Explor.*, 14(2):20–28, 2012.
- [77] A. Z. Tan, H. Yu, L. Cui, and Q. Yang. Towards personalized federated learning. *CoRR*, abs/2103.00710, 2021.
- [78] Y. Tan, Y. Liu, G. Long, J. Jiang, Q. Lu, and C. Zhang. Federated learning on non-iid graphs via structural knowledge sharing. *CoRR*, abs/2211.13009, 2022.
- [79] Y. Tao, Y. Li, and Z. Wu. Semigraphfl: Semi-supervised graph federated learning for graph classification. In G. Rudolph, A. V. Kononova, H. Aguirre, P. Kerschke, G. Ochoa, and T. Tušar, editors, *Parallel Problem Solving from Nature – PPSN XVII*, pages 474–487, Cham, 2022. Springer International Publishing.
- [80] A. Thakur, P. Sharma, and D. A. Clifton. Dynamic neural graphs based federated reptile for semi-supervised multi-tasking in healthcare applications. *IEEE Journal of Biomedical and Health Informatics*, 26(4):1761–1772, 2022.
- [81] J. R. Ullman and A. Sealfon. Efficiently estimating erdos-renyi graphs with node differential privacy. In H. M. Wallach, H. Larochelle, A. Beygelzimer, F. d’Alché-Buc, E. B. Fox, and R. Garnett, editors, *Advances in Neural Information Processing Systems 32: Annual Conference on Neural Information Processing Systems 2019, NeurIPS 2019, December 8-14, 2019, Vancouver, BC, Canada*, pages 3765–3775, 2019.
- [82] J. Upadhyay, S. Upadhyay, and R. Arora. Differentially private analysis on graph streams. In A. Banerjee and K. Fukumizu, editors, *The 24th International Conference on Artificial Intelligence and Statistics, AISTATS 2021, April 13-15, 2021, Virtual Event*, volume 130 of *Proceedings of Machine Learning Research*, pages 1171–1179. PMLR, 2021.
- [83] B. Wang, A. Li, M. Pang, H. Li, and Y. Chen. Graphfl: A federated learning framework for semi-supervised node classification on graphs. In X. Zhu, S. Ranka, M. T. Thai, T. Washio, and X. Wu, editors, *IEEE International Conference on Data Mining, ICDM 2022, Orlando, FL, USA, November 28 - Dec. 1, 2022*, pages 498–507. IEEE, 2022.
- [84] J. Wang, Z. Charles, Z. Xu, G. Joshi, H. B. McMahan, B. A. y Arcas, M. Al-Shedivat, G. Andrew, S. Avestimehr, K. Daly, D. Data, S. N. Diggavi, H. Eichner, A. Gadhikar, Z. Garrett, A. M. Girgis, F. Hanzely, A. Hard, C. He, S. Horváth, Z. Huo, A. Ingerman, M. Jaggi, T. Javidi, P. Kairouz, S. Kale, S. P. Karimireddy, J. Konečný, S. Koyejo, T. Li, L. Liu, M. Mohri, H. Qi, S. J. Reddi, P. Richtárik, K. Singhal, V. Smith, M. Soltanolkotabi, W. Song, A. T. Suresh, S. U. Stich, A. Talwalkar, H. Wang, B. E. Woodworth, S. Wu, F. X. Yu, H. Yuan, M. Zaheer, M. Zhang, T. Zhang, C. Zheng, C. Zhu, and W. Zhu. A field guide to federated optimization. *CoRR*, abs/2107.06917, 2021.
- [85] S. Wang, L. Hu, Y. Wang, X. He, Q. Z. Sheng, M. A. Orgun, L. Cao, F. Ricci, and P. S. Yu. Graph learning based recommender systems: A review. In Z. Zhou, editor, *Proceedings of the Thirtieth International Joint Conference on Artificial Intelligence, IJCAI 2021, Virtual Event / Montreal, Canada, 19-27 August 2021*, pages 4644–4652. ijcai.org, 2021.
- [86] S. Wang, M. Lee, S. Hosseinalipour, R. Morabito, M. Chiang, and C. G. Brinton. Device sampling for heterogeneous federated learning: Theory, algorithms, and implementation. In *IEEE INFOCOM 2021 - IEEE Conference on Computer Communications*, pages 1–10, 2021.
- [87] Y. Wang and X. Wu. Preserving differential privacy in degree-correlation based graph generation. *Trans. Data Priv.*, 6(2):127–145, 2013.
- [88] Y. Wang, X. Wu, and L. Wu. Differential privacy preserving spectral graph analysis. In J. Pei, V. S. Tseng, L. Cao, H. Motoda, and G. Xu, editors, *Advances in Knowledge Discovery and Data Mining, 17th Pacific-Asia Conference, PAKDD 2013, Gold Coast, Australia, April 14-17, 2013, Proceedings, Part II*, volume 7819 of *Lecture Notes in Computer Science*, pages 329–340. Springer, 2013.
- [89] Z. Wang, X. Fan, J. Qi, H. Jin, P. Yang, S. Shen, and C. Wang. Fedgs: Federated graph-based sampling with arbitrary client availability. *CoRR*, abs/2211.13975, 2022.
- [90] Z. Wang, W. Kuang, Y. Xie, L. Yao, Y. Li, B. Ding, and J. Zhou. Federatedscope-gnn: Towards a unified, comprehensive and efficient package for federated graph learning. In A. Zhang and H. Rangwala, editors, *KDD ’22: The 28th ACM SIGKDD Conference on Knowledge Discovery and Data Mining, Washington, DC, USA, August 14 - 18, 2022*, pages 4110–4120. ACM, 2022.
- [91] K. Wei, J. Li, M. Ding, C. Ma, H. H. Yang, F. Farokhi, S. Jin, T. Q. S. Quek, and H. V. Poor. Federated learning with differential privacy: Algorithms and performance analysis. *IEEE Trans. Inf. Forensics Secur.*, 15:3454–3469, 2020.
- [92] C. Wu, F. Wu, L. Lyu, T. Qi, Y. Huang, and X. Xie. A federated graph neural network framework for privacy-preserving personalization. *Nature Communications*, 13(1):3091, Jun 2022.

- [93] L. Wu, Y. Chen, K. Shen, X. Guo, H. Gao, S. Li, J. Pei, and B. Long. Graph neural networks for natural language processing: A survey. *Found. Trends Mach. Learn.*, 16(2):119–328, 2023.
- [94] X. Wu, X. Ying, K. Liu, and L. Chen. *A Survey of Privacy-Preservation of Graphs and Social Networks*, pages 421–453. Springer US, Boston, MA, 2010.
- [95] H. Xie, J. Ma, L. Xiong, and C. Yang. Federated graph classification over non-iid graphs. In M. Ranzato, A. Beygelzimer, Y. N. Dauphin, P. Liang, and J. W. Vaughan, editors, *Advances in Neural Information Processing Systems 34: Annual Conference on Neural Information Processing Systems 2021, NeurIPS 2021, December 6-14, 2021, virtual*, pages 18839–18852, 2021.
- [96] P. Xing, S. Lu, L. Wu, and H. Yu. Big-fed: Bilevel optimization enhanced graph-aided federated learning. *IEEE Transactions on Big Data*, pages 1–12, 2022.
- [97] C. Yang, H. Wang, K. Zhang, L. Chen, and L. Sun. Secure deep graph generation with link differential privacy. In *IJCAI 2021*, 2021.
- [98] Y. Yao and C. Joe-Wong. Fedgcn: Convergence and communication tradeoffs in federated training of graph convolutional networks. *CoRR*, abs/2201.12433, 2022.
- [99] X. Ying and X. Wu. Randomizing social networks: a spectrum preserving approach. In *SDM 2008*, 2008.
- [100] X. Ying and X. Wu. Graph generation with prescribed feature constraints. In *Proceedings of the SIAM International Conference on Data Mining, SDM 2009, April 30 - May 2, 2009, Sparks, Nevada, USA*, pages 966–977. SIAM, 2009.
- [101] X. Ying and X. Wu. On link privacy in randomizing social networks. *Knowl. Inf. Syst.*, 28(3):645–663, 2011.
- [102] T. Yoon, S. Shin, S. J. Hwang, and E. Yang. Fedmix: Approximation of mixup under mean augmented federated learning. In *9th International Conference on Learning Representations, ICLR 2021, Virtual Event, Austria, May 3-7, 2021*. OpenReview.net, 2021.
- [103] D. Zeng, S. Liang, X. Hu, H. Wang, and Z. Xu. Fedlab: A flexible federated learning framework. *Journal of Machine Learning Research*, 24(100):1–7, 2023.
- [104] H. Zhang, M. Cissé, Y. N. Dauphin, and D. Lopez-Paz. mixup: Beyond empirical risk minimization. In *6th International Conference on Learning Representations, ICLR 2018, Vancouver, BC, Canada, April 30 - May 3, 2018, Conference Track Proceedings*. OpenReview.net, 2018.
- [105] J. Zhang, G. Cormode, C. M. Procopiuc, D. Srivastava, and X. Xiao. Private release of graph statistics using ladder functions. In T. K. Sellis, S. B. Davidson, and Z. G. Ives, editors, *Proceedings of the 2015 ACM SIGMOD International Conference on Management of Data, Melbourne, Victoria, Australia, May 31 - June 4, 2015*, pages 731–745. ACM, 2015.
- [106] K. Zhang, C. Yang, X. Li, L. Sun, and S. Yiu. Subgraph federated learning with missing neighbor generation. In M. Ranzato, A. Beygelzimer, Y. N. Dauphin, P. Liang, and J. W. Vaughan, editors, *Advances in Neural Information Processing Systems 34: Annual Conference on Neural Information Processing Systems 2021, NeurIPS 2021, December 6-14, 2021, virtual*, pages 6671–6682, 2021.
- [107] M. Zhang, K. Sapra, S. Fidler, S. Yeung, and J. M. Alvarez. Personalized federated learning with first order model optimization. In *9th International Conference on Learning Representations, ICLR 2021, Virtual Event, Austria, May 3-7, 2021*. OpenReview.net, 2021.
- [108] Y. Zhang, S. Wei, S. Liu, Y. Wang, Y. Xu, Y. Li, and X. Shang. Graph-regularized federated learning with shareable side information. *Knowledge-Based Systems*, 257:109960, 2022.
- [109] G. Zhao, Y. Huang, and C. H. Tsai. Fedgsl: Federated graph structure learning for local subgraph augmentation. In *2022 IEEE International Conference on Big Data (Big Data)*, pages 818–824, 2022.
- [110] Y. Zhao, M. Li, L. Lai, N. Suda, D. Civin, and V. Chandra. Federated learning with non-iid data. *CoRR*, abs/1806.00582, 2018.
- [111] E. Zheleva and L. Getoor. Preserving the privacy of sensitive relationships in graph data. In F. Bonchi, E. Ferrari, B. A. Malin, and Y. Saygin, editors, *Privacy, Security, and Trust in KDD, First ACM SIGKDD International Workshop, PinKDD 2007, San Jose, CA, USA, August 12, 2007, Revised Selected Papers*, volume 4890 of *Lecture Notes in Computer Science*, pages 153–171. Springer, 2007.
- [112] B. Zhou and J. Pei. Preserving privacy in social networks against neighborhood attacks. In G. Alonso, J. A. Blakeley, and A. L. P. Chen, editors, *Proceedings of the 24th International Conference on Data Engineering, ICDE 2008, April 7-12, 2008, Cancún, Mexico*, pages 506–515. IEEE Computer Society, 2008.
- [113] B. Zhou, J. Pei, and W. Luk. A brief survey on anonymization techniques for privacy preserving publishing of social network data. *SIGKDD Explor.*, 10(2):12–22, 2008.
- [114] L. Zhu, Z. Liu, and S. Han. Deep leakage from gradients. In H. M. Wallach, H. Larochelle, A. Beygelzimer, F. d’Alché-Buc, E. B. Fox, and R. Garnett, editors, *Advances in Neural Information Processing Systems 32: Annual Conference on Neural Information Processing Systems 2019, NeurIPS 2019, December 8-14, 2019, Vancouver, BC, Canada*, pages 14747–14756, 2019.
- [115] L. Zou, L. Chen, and M. T. Özsu. K-automorphism: A general framework for privacy preserving network publication. *Proc. VLDB Endow.*, 2(1):946–957, 2009.

Adaptive Risk-Aware Bidding with Budget Constraint in Display Advertising

Zhimeng Jiang[†], Kaixiong Zhou[‡], Mi Zhang[§], Rui Chen[§], Xia Hu[‡], Soo-Hyun Choi[§]

Department of Computer Science and Engineering, Texas A&M University, TX, USA

Department of Computer Science, Rice University, TX, USA

Samsung Electronics America

[†]zhimengj@tamu.edu, [‡]{Kaixiong.Zhou,xia.hu}@rice.edu,

[§]{mi.zhang,rui.chen1,sh9.choi}@samsung.com

ABSTRACT

Real-time bidding (RTB) has become a major paradigm of display advertising. Each ad impression generated from a user visit is auctioned in real time, where demand-side platform (DSP) automatically provides bid price usually relying on the ad impression value estimation and the optimal bid price determination. However, the current bid strategy overlooks the randomness of the user behaviors (e.g., click) and the cost uncertainty caused by the auction competition. In this work, we propose a novel adaptive risk-aware bidding algorithm with budget constraint via reinforcement learning, which is the first to simultaneously consider estimation uncertainty and the dynamic risk tendency of a DSP. Specifically, we explicitly factor in the *uncertainty of estimated ad impression values* and model the *risk preference* of a DSP under a *specific state and market environment* via a sequential decision process. Additionally, we theoretically unveil the intrinsic relation between the uncertainty and the risk tendency based on value at risk (VaR). Consequently, we propose two instantiations to model risk tendency, including an expert knowledge-based formulation embracing three essential properties and an adaptive learning method based on self-supervised reinforcement learning. We conduct experiments on public datasets and show that the proposed framework achieves better performance in terms of the number of clicks under different budget constraints¹.

Keywords

Risk-Aware Bidding Strategy; Budget Constraint; Display Advertising

1. INTRODUCTION

In the past few years, real-time bidding (RTB) has quickly become tens of billions of markets in the globe [9; 34]. In RTB, a demand-side platform (DSP) buys ad impressions in a programmatic manner on behalf of advertisers. The success of a DSP heavily relies on its *bid optimization* (i.e., the process of identifying the optimal bid price for each bid request) capability [32], whose goal is to maximize the key performance indicator (KPI) agreed upon with advertisers

(e.g., the total number of clicks or return on ad spend). In practice, bid optimization normally involves two steps, namely user response prediction and bid price determination [30]. User response prediction is performed to estimate the true value of a potential ad impression. Taking the estimated value as an input, the bid price determination step aims to generate the optimal bid price for a bid request in sequential decision making.

RTB is a highly competitive and dynamic marketplace. The prerequisite for a DSP to stay competitive is the capability of accurately predicting user responses [4; 8], e.g., click-through rate (CTR) or conversion rate (CVR), and advanced bid price determination algorithm. A large number of prediction models [38; 27; 12; 33; 41; 6] have been proposed in the literature, and return a *point* estimation for a bid request. Subsequently, such estimation is used to approximate the true value of the corresponding ad impression. Despite the substantial progress that has been made, their accuracy is still far from perfect. One major reason is due to the large randomness of the user behaviors (e.g., click) and the cost uncertainty caused by the auction competition [1], which result in inherent *uncertainties* of estimated values. Motivated by work [37], explicitly factoring such uncertainties plays a critical role in optimizing a campaign's performance.

On the other hand, the dynamic nature of RTB requires modeling the correlations of bid requests under a given budget constraint in view of varying market competition [7]. The latest research considers a DSP's bidding process as a sequential decision process and proposes model-based or model-free reinforcement learning based bidding strategies [3; 31]. While it lays a solid groundwork for bid optimization, these studies are based on the fundamental assumption that the estimations of ad impressions' values are accurate. Unfortunately, as explained before, this assumption can hardly hold in practice. We point out that a reasonable bid optimization solution needs to consider three inherently correlated components, including the uncertainties of estimations of ad impression values, the state of a DSP (e.g., remaining budget and future auction number), and market competition. The latter two components determine the DSP's *risk tendency* (e.g., take more risk by bidding more aggressively or reduce risks by bidding more conservatively).

Based on the above observations, we propose an adaptive risk-aware bidding algorithm via reinforcement learning. To the best of our knowledge, it is the first work that simultaneously considers *prediction uncertainty* and the *dynamic risk tendency* of a DSP. We first theoretically unveil the intrinsic

¹Our code is available at: <https://github.com/zhimengj0326/ekRLB>

Table 1: Notations and descriptions

Notation	Description
\mathbf{x}	The features of bid request.
$r_{mean}(\mathbf{x})$	The mean value of predicted CTR.
$r_{std}(\mathbf{x})$	The standard deviation of predicted CTR.
(t, b)	The remaining auction numbers and budget.
θ	The ad impression value.
$m(\delta)$	The probability density distribution of market price.
$V(t, b)$	The expected total reward with starting state (t, b) taking the optimal policy.
$a(t, b, \mathbf{x})$	The bid price for the request features \mathbf{x} .
$\beta(t, b)$	The risk tendency for the resource state (t, b) .
$U(t, b)$	The expected bid price depleting remaining budget.

relation between prediction uncertainty and risk tendency, which helps generate a modified value of an ad impression. With this formulation of an ad impression value, it is critical to properly optimize a DSP’s risk tendency. To this end, we propose two instantiations to model risk tendency, including an expert knowledge based formulation embracing three essential properties and an adaptive learning method based on self-supervised reinforcement learning. We summarize our key contributions as follows.

- We present an adaptive risk-aware bidding algorithm, which, for the first time, considers both prediction uncertainty and dynamic risk tendency to optimize bidding performance. This framework is based on a new formulation of an ad impression value by revealing the intrinsic relation between prediction uncertainty and risk tendency. We theoretically prove that this formulation allows achieving the optimal bid price based on VaR analysis.
- We propose two ways to determine the risk tendency of a DSP. We identify three basic properties of risk tendency, which lead to an expert knowledge based instantiation. To mitigate the extensive manual tuning efforts, we also design a self-supervised reinforcement learning method to learn the risk tendency based on experience.
- We conduct extensive experiments on two public datasets to validate the superiority of our adaptive risk-aware bidding algorithm and demonstrate the benefits of considering both prediction uncertainty and risk tendency.

2. PROBLEM FORMULATION

In the RTB system, each bidder of a DSP acts on behalf of an advertiser and competes for the advertisement auction every time a bid request is generated from a user visit. Given each auction opportunity, the bidder estimates the ad impression value and uncertainty, and then determine the bid price to maximize the cumulative ad impression value². We aim to obtain better bidding strategy under the second-price auction, i.e., the bidder with the highest bid price wins the auction with the second-highest price payment. Related notations are summarized in Table 1.

2.1 Problem Definition

Considering budget constraints in real-time bidding, we formulate the bid optimization problem as a Markov decision

²In this paper, we use bidder and DSP interchangeably, and adopt CTR to estimate the ad impression value, while other metrics, such as CVR, can be adapted similarly.

process (MDP) in the *episode* level, where each episode consists of T sequential bid auctions accompanied with a budget of B . For each auction, we consider three pieces of critical information: (i) the remaining auction number $t \in \{0, \dots, T\}$; (ii) the remaining budget $b \in \{0, \dots, B\}$ and (iii) the mean value of predicted CTR (pCTR) $r_{mean}(\mathbf{x}_t)$ and the corresponding standard deviation $r_{std}(\mathbf{x}_t)$ for a bid request with feature vector \mathbf{x}_t . Hence, the bidder’s state s is defined as $s \triangleq (t, b, \mathbf{x}_t)$. We define $V(t, b, \mathbf{x})$ as the long-run expected ad impression value accumulated from the current state $s = (t, b, \mathbf{x})$. Our target problem is that, given the remaining action number t , remaining budget b , the mean value of pCTR $r_{mean}(\mathbf{x}_t)$ and corresponding standard deviation $r_{std}(\mathbf{x}_t)$ for current bid request features \mathbf{x}_t , how can we determine the optimal bid price $a(t, b, \mathbf{x}_t)$ to optimal cumulative ad impression value in a sequential decision-making process?

2.2 MDP Formulation

Reinforcement learning can be represented by tuple

$(\mathcal{S}, \mathcal{A}_s, \mathcal{P}_{sa}^s, \mathcal{R}_{sa}^s)$, where \mathcal{S} denotes the state space, \mathcal{A}_s denotes the action (i.e., bid price) space for state s , \mathcal{P}_{sa}^s and \mathcal{R}_{sa}^s represent the state transition probability and the immediate reward (i.e., pCTR) for the transition from state s to s' under action a . Note that $t = 0$ and $b = 0$ represent the end of an episode and the state with the depleted budget in sequential decision-making, respectively.

In the episode level bidding process, the state space $\mathcal{S} = \{0, \dots, T\} \times \{0, \dots, B\} \times \mathbf{X}$, where \mathbf{X} denotes the set of bid request features. Given state $s = (t, b, \mathbf{x}_t)$, the action space \mathcal{A}_s consists of all possible bid prices in set $\{0, \dots, b\}$, since possible bid prices are constrained by the remaining budget b . Let $p\mathbf{x}(\mathbf{x}_t)$ denote the probability of the bid request feature \mathbf{x}_t for a potential ad impression and $m(\delta|\mathbf{x}_t)$ denote the probability of market price δ given feature \mathbf{x}_t . As there is no obvious dependency between winning price distribution and bid request features [40] in iPinYou dataset³, we assume that market environment $m(\delta) \simeq m(\delta|\mathbf{x}_t)$, i.e., the market price distribution is independent of the bid request feature. Such independent distribution assumption can be justified by empirical evaluation via comparing winning bid distribution against different features in the real-world iPinYou dataset [40].

For the state transition, if the bid price a is larger than the market price δ (highest bid price among other competitors), then the bidder wins the ad auction, and state (t, b, \mathbf{x}_t) will transit to $(t - 1, b - \delta, \mathbf{x}_{t-1})$ with probability $p\mathbf{x}(\mathbf{x}_{t-1}) \sum_{\delta=0}^a m(\delta)$. Otherwise, if $a < \delta$, the bidder will lose the auction and transit to state $(t - 1, b, \mathbf{x}_{t-1})$ with probability $p\mathbf{x}(\mathbf{x}_{t-1}) \sum_{\delta=a+1}^{+\infty} m(\delta)$. The immediate reward is given by $r_{mean}(\mathbf{x}_t)$ for t -th auction if the bidder wins the auction; otherwise, it is 0. Mathematically, the state transition probability and reward function are expressed as follows:

$$\mathcal{P}_{(t,b,\mathbf{x}_t),a}^{(t-1,b-\delta,\mathbf{x}_{t-1})} = \begin{cases} p\mathbf{x}(\mathbf{x}_{t-1})m(\delta), & \text{if } \delta \leq a; \\ p\mathbf{x}(\mathbf{x}_{t-1}) \sum_{\delta=a+1}^{+\infty} m(\delta), & \text{if } \delta > a. \end{cases}$$

³The similar box plot of winning price distribution against the features, such as hour, weekday, user browser, operation system and location regions of bid requests, can be observed in many campaigns of iPinYou [40].

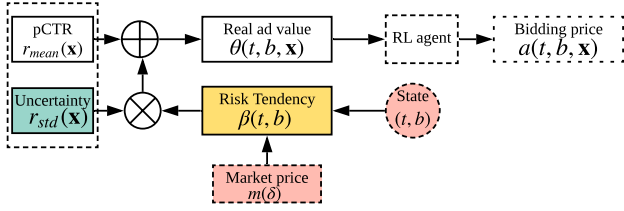


Figure 1: An overview of the adaptive risk-aware bidding framework. Both pCTR and prediction uncertainty come from the auction environment. Considering the market environment and bidder state, the framework combines prediction uncertainty and risk tendency to adjust ad impression values and decide bid prices.

$$\mathcal{R}_{(t,b,\mathbf{x}_t),a}^{(t-1,b-\delta,\mathbf{x}_{t-1})} = \begin{cases} \theta(t,b,\mathbf{x}_t), & \text{if } \delta \leq a; \\ 0, & \text{if } \delta > a. \end{cases}$$

3. METHODOLOGY

In this section, we introduce the **R**isk-aware **R**einforcement **L**earning **B**idding (RRLB) framework that effectively integrates prediction uncertainty and a bidder’s risk tendency into a reinforcement learning framework. An overview of RRLB is illustrated in Figure 1. In the following sections, we first explain the uncertainties of CTR prediction with Bayesian logistic regression and then describe the risk-aware bid optimization framework and two proposed instantiations to determine the *risk tendency*. Finally, we describe the model-based reinforcement learning method mapping the adjusted ad impression value to the final bid price.

3.1 Uncertainty of CTR Prediction

Similar to [37], we adopt Bayesian logistic regression to explicitly measure the uncertainties of predicted CTR (pCTR) values. In Bayesian logistic regression, each weight \mathbf{w} is treated as a random variable instead of a parameter, and the variance of the random variable represents the uncertainty of the corresponding feature. The likelihood of observing the correct binary click label y given features \mathbf{x} and weights of logistic regression \mathbf{w} is

$$p(y|\mathbf{x}, \mathbf{w}) = \sigma(\mathbf{w}^T \mathbf{x})^y (1 - \sigma(\mathbf{w}^T \mathbf{x}))^{1-y}, \quad (1)$$

where the sigmoid function is defined as $\sigma(x) = \frac{1}{1+e^{-x}}$. Note that \mathbf{w} is modeled as a random variable with a p.d.f $p(\mathbf{w})$ in the Bayesian version of logistic regression. The marginal conditional probability $p(y|\mathbf{x})$ is $p(y|\mathbf{x}) = \int p(y|\mathbf{x}, \mathbf{w})p(\mathbf{w})d\mathbf{w}$. The Gaussian prior $N(\mu_0, q_0^{-1}\mathbf{I})$ and Laplace approximation of posterior on \mathbf{w} are adopted for weight distribution update, where q_0 is initial weight precision. The goal is to maximize the posterior weight distribution $p(\mathbf{w}|\mathbf{x}, y)$ given a bid request’s feature vector \mathbf{x} and label y . The model output is a probability estimation of the occurrence of a click event, defined as $\hat{y} = P(y = 1|\mathbf{x})$. The variance of weight is lower when the associated feature emerges more frequently, which means that such a model can measure the data completeness for each feature. Via updating the mean and covariance matrix of the weight \mathbf{w} , the distribution of CTR $p_{\hat{y}|\mathbf{x}}(\hat{y})$ can be obtained. Subsequently, we define the mean and standard deviation of CTR as $r_{mean}(\mathbf{x}) = \mathbb{E}_{p_{\hat{y}|\mathbf{x}}}[\hat{y}]$ and $r_{std}(\mathbf{x}) = \sqrt{\mathbb{D}_{p_{\hat{y}|\mathbf{x}}}[\hat{y}]}$, where $\mathbb{E}_{p_{\hat{y}|\mathbf{x}}}[\cdot]$ and $\mathbb{D}_{p_{\hat{y}|\mathbf{x}}}[\cdot]$ denote the

expectation and variance over the CTR distribution $p_{\hat{y}|\mathbf{x}}$. The standard deviation of CTR reflects the uncertainties of pCTR values. How to obtain uncertainties of other prediction models is beyond the scope of this paper.

3.2 Theoretical Relation Between Uncertainty and Risk Tendency

The key intuition of RRLB is to decompose the value of an ad impression $\theta(t, b, \mathbf{x}_t)$ as the weighted sum of two parts: the mean pCTR and a compound term that simultaneously reflects prediction uncertainty and a bidder’s risk tendency. Formally, the ad impression value is defined as follows:

$$\theta(t, b, \mathbf{x}_t) = r_{mean}(\mathbf{x}_t) + \beta(t, b)r_{std}(\mathbf{x}_t), \quad (2)$$

where $\beta(t, b)$ denotes the bidder’s risk tendency at resource state (t, b) , which is a subset of bidding state (t, b, \mathbf{x}_t) , which represents the intrinsic status of a bidder in terms of remaining auction number t and remaining budget b in an episode. Next we provide a theoretical motivation for the formulation of Eq. (2). The core idea is based on the value at risk (VaR) theory borrowed from finance [26; 19], where VaR estimates how much the predicted CTR under/over-estimates with a given probability. Note that the goal of a bidding strategy is to improve the cumulative ad impression value. Given the current bid request feature vector \mathbf{x}_t , remaining budget b , and remaining auctions number t , let $V^a(t, b, \mathbf{x}_t)$ and $V_{std}^a(t, b, \mathbf{x}_t)$ be the cumulative estimated impression value and uncertainty of the winning ads with the bidding strategy $a(t, b, \mathbf{x}_t)$. Then we define VaR of the cumulative ad impression value with the bidding strategy $a(t, b, \mathbf{x}_t)$ as follows:

$$V_\lambda^a(t, b, \mathbf{x}_t) \triangleq V^a(t, b, \mathbf{x}_t) + \lambda(t, b)V_{std}^a(t, b, \mathbf{x}_t), \quad (3)$$

where state-associated coefficient $\lambda(t, b)$ is the risk preference that balances the cumulative estimated impression value and uncertainty. Note that $\beta(t, b)$ and $\lambda(t, b)$ balance the estimated value and uncertainty for the current auction and all the remaining auctions, respectively. The optimal VaR bid price a_{VaR} maximizes the VaR of the cumulative ad impression value $V_\lambda^a(t, b, \mathbf{x}_t)$ as follows,

$$a_{VaR}(t, b, \mathbf{x}_t) = \arg \max_{0 \leq a \leq b} V_\lambda^a(t, b, \mathbf{x}_t). \quad (4)$$

The theorem below shows that the linear combination in Eq. (2) can achieve the optimal VaR bid price a_{VaR} .

Theorem 1 (Risk Tendency Optimality). *The RRLB framework adopting the linear formulation in Eq. (2) can achieve the optimal VaR bid price $a_{VaR}(t, b, \mathbf{x}_t)$ if the same risk tendency $\beta(t, b) = \lambda(t, b)$ is used.*

PROOF. *Note that the bidder wins an ad impression if the bid price a is larger than the market price δ . The cumulative estimated ad impression value and the corresponding uncertainty should satisfy*

$$V^a(t, b, \mathbf{x}_t) = \sum_{\delta=0}^{a(t,b,\mathbf{x}_t)} m(\delta)r_{mean}(\mathbf{x}_t) + \sum_{\delta=0}^{a(t,b,\mathbf{x}_t)} m(\delta)V(t-1, b-\delta) + \sum_{\delta=a(t,b,\mathbf{x}_t)+1}^{\infty} m(\delta)V^a(t-1, b),$$

$$V_{std}^a(t, b, \mathbf{x}_t) = \sum_{\delta=0}^{a(t, b, \mathbf{x}_t)} m(\delta) r_{std}(\mathbf{x}_t) + \sum_{\delta=0}^{a(t, b, \mathbf{x}_t)} m(\delta) V_{std}(t-1, b-\delta) + \sum_{\delta=a(t, b, \mathbf{x}_t)+1}^{\infty} m(\delta) V_{std}(t-1, b),$$

Combining above two equations, we have

$$V_{\lambda}^a(t, b, \mathbf{x}_t) = \sum_{\delta=0}^{a(t, b, \mathbf{x}_t)} m(\delta) (r_{mean}(\mathbf{x}_t) + \lambda(t, b) r_{std}(\mathbf{x}_t)) + \sum_{\delta=0}^{a(t, b, \mathbf{x}_t)} m(\delta) V_{\lambda}^a(t-1, b-\delta) + \sum_{\delta=a(t, b, \mathbf{x}_t)+1}^{\infty} m(\delta) V_{\lambda}^a(t-1, b),$$

By firstly setting $\theta(t, b, \mathbf{x}_t) = r_{mean}(\mathbf{x}_t) + \beta(t, b) r_{std}(\mathbf{x}_t)$ and adopting the model-based method [3], it can be seen that the linear combination in $\theta(t, b, \mathbf{x}_t)$ can achieve the optimal bid price $a_{VaR}(t, b, \mathbf{x}_t)$, which maximizes the VaR of the cumulative ad impression value.

Intuitively, a rational risk tendency should be a function of the resource state, defined by (t, b) , of a bidder. A budget-restrained/abundant bidder would act very differently in taking risks during the bid. However, we deem that the risk tendency is independent of a random bid request \mathbf{x} .

3.3 Expert Knowledge Based Risk Tendency

We leverage expert knowledge on RTB to design the first instantiation of risk tendency $\beta(t, b)$ to reveal the intrinsic risk preference of a rational bidder. We distill three key rules as follows. (i) The sufficiency of the remaining budget b determines the sign of risk tendency $\beta(t, b)$, where a positive risk tendency indicates a strong preference to win auctions. (ii) The partial derivative of risk tendency $\beta(t, b)$ w.r.t. remaining budget b (remaining auction number t) should be positive (negative) because more budget naturally allows the bidder to take the risk of bidding more ad impressions. (iii) When the remaining budget b and remaining auction number t are relatively large (e.g., the beginning of an ad campaign), risk tendency $\beta(t, b)$ depends on the ratio of b to t and the extent of market competition. Formally, we express the three rules as follows.

(i) Sign of risk tendency:

$$\beta(t, b) \begin{cases} \geq 0, & b \text{ is sufficient at current } t; \\ < 0, & \text{otherwise.} \end{cases} \quad (5)$$

(ii) Monotonicity of risk tendency: $\frac{\partial \beta(t, b)}{\partial t} < 0$, $\frac{\partial \beta(t, b)}{\partial b} > 0$.

(iii) Approximation for the scenarios of large remaining budget and auction number: $\beta(t, b) \simeq \beta(t', b')$ if $\frac{b}{t} = \frac{b'}{t'}$.

Before elaborating the exact formulation of $\beta(t, b)$, we have to quantify the sufficiency of the budget as pointed out in the first rule. In practice, the budget richness is highly related to the level of market competition. Given a certain market price distribution $m(\delta)$ and resource state (t, b) , let $U(t, b)$ denote the expected bid price for an auction such that the remaining budget will be depleted in the remaining future auctions. Supposing that budget b is evenly allocated to the remaining t auctions, we can calculate $U(t, b)$ through the following formula:

$$\sum_{\delta=0}^{U(t, b)} \delta m(\delta) = \frac{b}{t}. \quad (6)$$

Algorithm 1: Expert knowledge based reinforcement learning bidding strategy

Input: market price probability $m(\delta)$, episode length T , budget B , average pCTR and standard deviation of CTR \bar{r}_{mean} and \bar{r}_{std}

Output: Optimal bid price a for current state
Parameter tuning for the manually designed risk tendency;

Update Value function $V(t, b)$ based on risk tendency ;

for $\delta = 0, 1, \dots, \min(\delta_{max}, b)$ **do**

if $\theta(t, b, \mathbf{x}) + V(t-1, b-\delta) - V(t-1, b) < 0$ **then**

$a(t, b, \mathbf{x}) \leftarrow \delta$;

break;

end

Since the bidder wins an auction only if bid price $U(t, b)$ is higher than market price δ , the left side of Eq. (6) represents the expected actual cost, which should be the same as $\frac{b}{t}$ based on the assumption of even budget allocation.

Based on the above intuition of budget richness, we formally define risk tendency as follows:

$$\beta(t, b) = \tanh\left(\alpha \frac{U(t, b) - \hat{U}}{\hat{U}}\right), \quad (7)$$

where α is a positive hyperparameter that controls the slope of risk tendency, \hat{U} is the budget richness threshold tuned from historical data, and function $\tanh(\cdot)$ confines risk tendency within the range $(-1, 1)$. It can be observed that the proposed risk tendency formulation satisfies all three expert knowledge based rules. First, the expected bid price $U(t, b)$, which measures the amount of budget richness, determines the sign of risk tendency $\beta(t, b)$, which is non-negative only if $U(t, b) \geq \hat{U}$ as required in rule (i). Second, both $U(t, b)$ and $\beta(t, b)$ increase with budget b , and decrease with t as required in rule (ii). Third, the expected cost is proportional to the ratio $\frac{b}{t}$ as shown in Eq. (6), which helps the subsequent design of risk tendency to meet rule (iii).

By using the expert knowledge based risk tendency, we can obtain the adjusted ad impression value using $\theta(t, b, \mathbf{x}) = r_{mean}(\mathbf{x}) + \beta(t, b) r_{std}(\mathbf{x})$. Subsequently, the model-based reinforcement learning method is applied to obtain the final bid price as will be introduced in Section 3.5.

Compared with the previous efforts [3], the bid price not only reveals the estimated mean value of an ad impression but also takes into account the estimation risk and risk tendency based on the bidding state. The expert knowledge based Reinforcement Learning Bidding (ekRLB) algorithm is summarized in Algorithm 1, where the value function update could be calculated using dynamic programming.

3.4 Self-Supervised Risk Tendency

Although the above explicit design of risk tendency may capture the intrinsic risk preference well, it requires extensive human expert knowledge and efforts on tuning Equation (7). To be specific, a careful selection of hyperparameters α and \hat{U} is required through many trials. With the variation of the market environment, these presetting hyperparameters tend to obtain sub-optimal performance and need to be updated periodically. Furthermore, the computation cost of calculating budget richness $U(t, b)$ is prohibitive as it requires considering all combinations of t and b . To avoid such

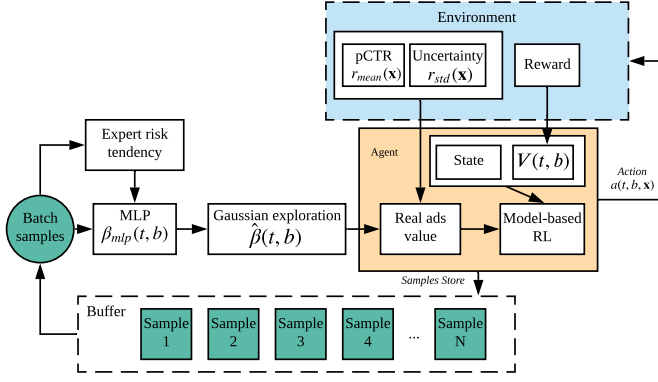


Figure 2: An overview of the self-supervised reinforcement learning. The risk tendency is generated from an MLP trained by experience buffer and batch sampling in a self-supervision way. The bidder decides bid prices based on risk tendency and interacts with the RTB environment.

manual efforts, we propose a self-supervised reinforcement learning bidding (ssRLB) method as shown in Figure 2 to automatically generate risk tendency via a multi-layer perceptron (MLP). Self supervised by the bidding history, we update the mapping function $\beta_{mlp}(t, b) = MLP(t, b; \mathbf{W}_{mlp})$ generated from an MLP to approximate the risk tendency at resource state (t, b) with trainable weight \mathbf{W}_{mlp} . The ssRLB framework consists of a Gaussian exploration block, an experience buffer, an MLP mapping function, and batch sampling. We explain the details of each component as follows:

Gaussian exploration. We adopt exploration on risk tendency by adding Gaussian noise to $\beta_{mlp}(t, b)$ as $\hat{\beta}(t, b) = \beta_{mlp}(t, b) + \epsilon$, where noise $\epsilon \sim \mathcal{N}(0, \sigma^2)$. The noise variance σ^2 is adjustable to provide a trade-off between exploitation and exploration in reinforcement learning.

Experience buffer. Motivated by experience replay optimization in reinforcement learning [35], we adopt the experience buffer to store good experiences represented by a quaternary set $\mathcal{B} = (t, b, \hat{\beta}(t, b), V_{episode})$ from the bidding history, where $V_{episode}$ denotes the cumulative reward for the entire episode. A “good” experience means that a larger reward $V_{episode}$ is obtained by using risk tendency $\hat{\beta}(t, b)$. Let N be the buffer’s total length. The samples with the lowest reward will be removed if the buffer is full. In this way, the experience buffer could always provide the best samples explored so far for training the MLP.

Batch sampling. Batch sampling is responsible for sampling batch \mathcal{B}_{batch} from the buffer. We apply a simple uniform sampling to generate \mathcal{B}_{batch} .

Training MLP mapping function. Given the experience batch \mathcal{B}_{batch} , we update the weight in the MLP mapping function by minimizing the mean square loss function:

$$\mathcal{L} = \sum_{(t, b, \hat{\beta}(t, b), \cdot) \in \mathcal{B}_{batch}} \|MLP(t, b; \mathbf{W}_{mlp}) - \hat{\beta}(t, b)\|^2.$$

Since we only preserve the experiences with larger rewards in the buffer, the mapping function will be updated under supervision toward learning a good risk tendency.

In a nutshell, experience buffer stores “good” experienced risk tendency while the sampled experienced risk tendency is a supervised signal for MLP mapping function training.

Algorithm 2: Self-Supervised Risk Tendency Learning Algorithm

Input: The historical data sample with pCTR, risk, market price, and click labels, episode length T , budget B .

Output: Optimal bid price

Initialize the risk tendency, and uniform replay policy ;
Update Value function $V(t, b)$ based on Equation (8);

for each episode do

for each ad auction in the current episode do

 Provide the bid price based on Algorithm 1;

 Execute auction and observe $(t + 1, b)$ and cumulative reward and risk;

end

 Calculate the cumulative reward for an entire episode;

if the cumulative reward $V_{episode}$ is larger than that in Buffer then

$\mathcal{B} \leftarrow (t, b, \hat{\beta}(t, b), V_{episode})$;

 Uniformly sample a batch \mathcal{B}_s from \mathcal{B} ;

 Train the MLP based on the batch sample ;

 Update risk tendency $\hat{\beta}(t, b)$;

 Update Value function $V(t, b)$ based on Eq. (8).

end

The algorithm for self-supervised Reinforcement Learning Bidding (ssRLB) is shown in Algorithm 2.

3.5 Bid Price Determination

Previous sections introduce how to modify the ad impression value with prediction uncertainty and risk tendency. Next, we explain how to calculate the final bid price. We adopt the model-based reinforcement learning bidding strategy [3] to maximize the cumulative reward. Specifically, we regard the pCTR $r_{mean}(\mathbf{x}_t)$ as the immediate reward for t -th auction, and cumulative reward $V(t, b, \mathbf{x}_t)$ is defined as the expected cumulative reward starting from state (t, b, \mathbf{x}_t) with the optimal bid price. By definition, we have $V(0, b, \mathbf{x}_t) = V(0, b) = 0$ since there is no available auction. Similar to [3], the updated policy for the cumulative reward is given by:

$$V(t, b, \mathbf{x}_t) = \max_{0 \leq a \leq b} \left\{ \sum_{\delta=0}^a \int_{\mathbf{X}} m(\delta) p_{\mathbf{x}}(\mathbf{x}_{t-1}) \cdot (r_{mean}(\mathbf{x}_t) + V(t-1, b-\delta, \mathbf{x}_{t-1})) d\mathbf{x}_{t-1} + \sum_{\delta=a+1}^{+\infty} \int_{\mathbf{X}} m(\delta) p_{\mathbf{x}}(\mathbf{x}_{t-1}) V(t-1, b, \mathbf{x}_{t-1}) d\mathbf{x}_{t-1} \right\},$$

where \mathbf{X} represents the entire feature vector space, and the two integrations represent the immediate reward for winning and losing cases. Furthermore, the cumulative reward without observation on \mathbf{x} can be obtained by integrating the over bid request feature vector \mathbf{x} . Formally, the cumulative reward $V(t, b)$ is:

$$V(t, b) \approx \max_{0 \leq a \leq b} \left\{ \sum_{\delta=0}^a m(\delta) r_{avg} + \sum_{\delta=0}^a m(\delta) V(t-1, b-\delta) + \sum_{\delta=a+1}^{\infty} m(\delta) V(t-1, b) \right\},$$

where $r_{avg} = \int_{\mathbf{X}} p_{\mathbf{x}}(\mathbf{x}_{t-1}) r_{mean}(\mathbf{x}_{t-1}) d\mathbf{x}_{t-1}$ is the average ad impression value over the entire feature vector space. The cumulative reward $V(t, b)$ can be iteratively updated given the average ad impression value and market price distribution. Note that $\sum_{\delta=0}^{\infty} m(\delta) = 1$. The bid price at state (t, b, \mathbf{x}_t) is calculated by:

$$\begin{aligned} a(t, b, \mathbf{x}_t) &\triangleq \arg \max_{0 \leq a \leq b} V(t, b, \mathbf{x}_t) \\ &= \arg \max_{0 \leq a \leq b} \left\{ \sum_{\delta=0}^a m(\delta) \left(\theta(t, b, \mathbf{x}_t) \right. \right. \\ &\quad \left. \left. + V(t-1, b-\delta) - V(t-1, b) \right) \right\} \\ &= \arg \max_{0 \leq a \leq b} \left\{ \sum_{\delta=0}^a m(\delta) g(\delta) \right\}, \end{aligned} \quad (8)$$

where $g(\delta) \triangleq \theta(t, b, \mathbf{x}_t) + V(t-1, b-\delta) - V(t-1, b)$. The cumulative reward $V(t, b)$ monotonically increases w.r.t. the remaining budget b , and $g(\delta)$ monotonically decreases. If $g(b) > 0$ (e.g., the impression value is extremely high), the optimal bid price is all the remaining budget. If $g(b) < 0$ (e.g., the impression value is moderate), there must exist a *unique* integer price, defined as A , satisfying the following three conditions: (1) $0 \leq A \leq b$; (2) $g(A) \geq 0$ and (3) $g(A+1) < 0$. According to Eq. (8), the optimal price is thus determined by this *unique* integer price A . To summarize, the optimal bid price is ultimately dictated by:

$$a(t, b, \mathbf{x}_t) = \begin{cases} b, & \text{if } g(b) \geq 0; \\ A, & \text{if } g(b) < 0. \end{cases} \quad (9)$$

4. EXPERIMENTS

In this section, we conduct experiments to evaluate our RRLB framework with the two instantiations of risk tendency, ekRLB, and ssRLB, and answer the following three research questions:

- **Q1:** How does our risk-aware bid optimization solution compare with several baselines in terms of the total number of clicks?
- **Q2:** How do prediction uncertainty and risk tendency affect the total number of clicks?
- **Q3:** How do the hyperparameters affect the long-run performance of the risk-aware bidding strategy?

4.1 Experimental Setting

4.1.1 Datasets.

Our experiments are conducted on two real-world datasets, iPinYou⁴ and YOYI⁵. We follow the data preprocessing in [37] to split training/test sets and obtain the estimations and uncertainties of ad impression values. The dataset description is as follows:

- **iPinYou** dataset contains 19.5M ad impressions, 14.79K clicks, and 16.0K spend (in CNY) over 9 campaigns during 10 days in the year 2013. We follow the data preprocessing configuration in [37] to split it into training/test sets and obtain the estimations and risk of ad impression values.

⁴<https://contest.ipinyou.com/>

⁵<http://apex.sjtu.edu.cn/datasets/7>

Table 2: The comparison of total click number ($c_0 = 1/2$).

iPinYou	Lin	RLB	ekRLB	CRTRLB	CURLB	ssRLB
1458	401	428	428	416	423	422
2259	19	73	75	59	70	70
2261	9	44	51	36	64	61
2821	109	209	211	198	214	206
2997	295	376	382	384	360	276
3358	208	233	233	231	222	217
3386	157	293	294	290	293	290
3427	263	290	292	295	286	302
3476	206	230	232	241	233	204
Average	185.2	241.8	244.2	238.9	240.6	226.9
YOYI	725	890	894	873	840	914

- **YOYI** dataset includes 402M ad impressions, 500K clicks, and 428K spend (in CNY) during 8 days in the year 2016. The first 7 days and the last day are set as the training data and test data, respectively.

4.1.2 Compared methods.

We compare ekRLB and ssRLB with two state-of-the-art baselines. **Lin** is a linear bidding strategy with bid price $a_{Lin} = b_0 \theta(\mathbf{x})$, where parameter b_0 can be tuned on training data [24]. **RLB** is a model-based reinforcement learning bidding strategy [3], which achieves fine-grained pCTR-price mapping, and adaptively adjusts bid prices for different states. These two baselines only make use of $r_{mean}(\mathbf{x})$. Besides, we also consider two variants of our proposed reinforcement learning bidding framework, including constant risk tendency (**CRTRLB**) and constant uncertainty (**CURLB**). CRTRLB sets risk tendency to a constant for any resource states, i.e., $\beta(t, b) \equiv \beta_0$, which is equivalent to the previous work [37]. CURLB assigns a constant risk for all ad impressions, i.e., $r_{std}(\mathbf{x}) \equiv r_0$. These two variants can be used to validate the superiority of dynamic risk tendency and the necessity for risk estimation.

4.1.3 Evaluation sketch.

Given a predefined budget and episode length for each bidding campaign, we evaluate different bidding strategies in terms of the total click number. Evaluation sketch including experiment flow, budget constraints, and more details as follows.

- *Experiment flow.* The bidding data is a list of ad impression records, each of which is accompanied by the bid request feature vector, the market price, and the user response (click) label. Based on Bayesian logistic regression, we obtain the estimation mean and risk of ad displaying the value for each impression record. Following the model-based reinforcement learning bidding in Section 3.5, we divide training and test datasets into episodes, each of which contains T impression records and budget B . All bidding strategies are evaluated in the episode level. Episode length T influences the computation complexities of bidding price and budget richness, which are required over all remaining impression numbers t . We set the size of episode length $T = 1000$. In the second-price auction platform, a bidding agent wins an ad impression if its price is higher than the market price, and only needs to pay the cost of the market price. Then the bidding

agent receives the reward of the click label and applies it to update the bidding strategy.

- *Budget constraints.* We allocate the budget as follows: $B = CPM_{train} \times 10^{-3} \times T \times c_0$, where CPM_{train} and c_0 are the cost per mille impressions in the training dataset and budget coefficient, respectively. We compare the models using the budget coefficient set $\{1/32, 1/16, 1/8, 1/4, 1/2\}$.
- *More details.* For ekRLB, we tune the risk tendency hyperparameters α and \hat{U} in Eq. (7) on the training dataset to optimize the number of total clicks. For ssRLB, the mapping function $MLP(t, b; \mathbf{W}_{mlp})$ is realized by a four-layer MLP with 64 hidden units. Weight \mathbf{W}_{mlp} is trained with the Adam optimizer to minimize the mean square loss function at the learning rate of 1×10^{-3} . The experience buffer size and batch size are 1×10^5 and 32, respectively. We update the buffer every 5 training episode.

4.2 Comparison Results

We present the total click number of different methods on the 9 campaigns of iPinYou and YOYI with budget coefficient $c_0 = 1/2$ in Table 2 to answer question Q1. On iPinYou, ekRLB obtains the best performance on most campaigns, achieving the largest average total click number of 244.2. On YOYI, both ekRLB and ssRLB outperform RLB, and ssRLB achieves the largest click number of 914. The two variants CRTRLB and CURLB obtain fewer clicks on both datasets, which validates the benefits of considering both uncertainty and risk tendency in Eq. (2). Note that ssRLB obtains fewer clicks than ekRLB on iPinYou but more on YOYI probably because the training of self-supervised reinforcement learning is less stable. Nevertheless, we deem that ssRLB is still a valuable alternative since it does not require manual tuning of α and \hat{U} over a large volume of historical data. We leave the improvement on risk tendency learning for future work.

In the left part of Figure 3, we further show the performance improvements of RLB and ekRLB over Lin on iPinYou under the entire budget coefficient set. Two notable findings can be observed. First, the two reinforcement learning based strategies outperform Lin consistently over all budget settings, which validates the benefits of modeling the bidding process as an MDP. Given a specific budget constraint, all bid requests are inherently correlated instead of independent in Lin since previous bid prices determine the remaining budget. Second, ekRLB gets more clicks than the traditional RLB, especially in the case with larger budgets. Such observation demonstrates the advantage of explicitly modeling prediction uncertainty within a risk-aware reinforcement learning framework.

4.3 Ablation Study

To answer Q2, we study the individual contributions of prediction uncertainty and risk tendency by comparing with constant uncertainty and risk tendency. Regarding ekRLB as the benchmark method, the right part of Figure 3 shows the performance degradation of the two variants CRTRLB and CURLB. Along with the results in Table 2, we can find that (i) both CRTRLB (with constant risk tendency) and CURLB (with constant uncertainty) perform worse than ekRLB, which proves that it is critical to consider both prediction uncertainty and risk tendency to achieve an optimal bid optimization framework. (ii) CURLB achieves bet-

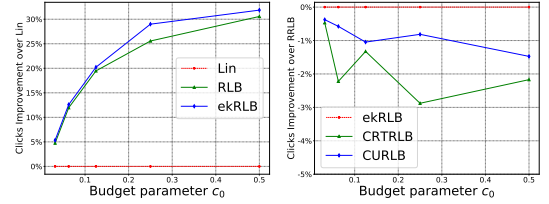


Figure 3: Left: The click number improvements of RLB and ekRLB over Lin in all budget coefficients on iPinYou. Right: The click number diminution of CRTRLB and CURLB over ekRLB on iPinYou.

ter performance than CRTRLB, suggesting that modeling risk tendency is even more important than prediction uncertainty. This is because risk tendency determines how to use (e.g., add or subtract) prediction uncertainty.

4.4 Hyperparameter Study

To answer Q3, we give a comprehensive hyperparameter study to investigate how hyperparameters slope α , constant uncertainty r_0 , and constant risk tendency β_0 affect the performance of the methods ekRLB, CURLB, and CRTRLB, respectively. Specifically, we report the experimental results on campaign 1458 of the iPinYou dataset with budget coefficient $c_0 = 1/2$ in Table 3. The performance metrics include click number and budget consumption ratio, where the budget consumption ratio is defined as the cumulative cost over the overall budget. The detailed experimental results and analysis are summarized as follows.

Slope α for ekRLB. Hyperparameter α controls the slope of risk tendency, and then influences the bid price of ekRLB. Note that, in the case $\alpha = 0$, we have 0 risk tendency based on Eq. (6) of the main text, and ekRLB degenerates to RLB. We can clearly observe that the total click number reaches the highest value at a medium slope scale of $\alpha = 0.1$. This result implies that the total click number is not sensitive w.r.t α and that ekRLB is robust on the hyperparameter α .

Constant uncertainty r_0 for CURLB. The method CURLB achieves relatively comparable performance with ekRLB even though a constant uncertainty is applied. We further study the influence of constant uncertainty r_0 by varying it in range $[0, 1.8 \times \bar{r}_{std}]$. Here we set constant uncertainty r_0 as a constant coefficient multiplying with the average uncertainty \bar{r}_{std} in the training dataset. For the case $r_0 = 0$, we have a prediction uncertainty of 0 for all ad impressions, when CURLB degenerates to RLB. We observe that the hyperparameter $r_0 = 0.2 * \bar{r}_{std}$ achieves the best performance in terms of both click numbers since higher/lower risks lead to ad impressions with overly large/small prices that decrease the total click number.

Constant risk tendency β_0 for CRTRLB. Compared with ekRLB, it can be observed that a random selection of constant risk tendency β_0 in CRTRLB greatly damages model performance. We further study its influences by considering β_0 within range $[-0.5, 0]$. We only use the negative risk tendencies since a positive one tends to overestimate an ad impression and usually results in worse performance. For the case $\beta_0 = 0$, we remove risk tendency and degenerate CRTRLB to RLB. We observe that $\beta_0 = 0.0$ achieves the most clicks and that when β_0 equals -0.4 to -0.5 , CRTRLB extremely underestimates the ad impression value and bids

Table 3: Hyperparameter study on slope α (ekRLB), uncertainty r_0 (CURLB) and risk tendency β_0 (CRTRLB).

Hyperparameters		#click	consumption ratio
α	0.0	1928	98.70%
	0.001	1925	98.71%
	0.01	1927	98.86%
	0.1	1930	99.54%
	0.2	1923	99.00%
	0.3	1923	99.11%
	0.4	1924	99.20%
	0.5	1928	99.29%
r_0	0	1928	99.83%
	$0.2 * \bar{r}_{std}$	1932	99.88%
	$0.4 * \bar{r}_{std}$	1922	99.89%
	$0.6 * \bar{r}_{std}$	1919	99.90%
	$0.8 * \bar{r}_{std}$	1910	99.90%
	$1.0 * \bar{r}_{std}$	1898	99.91%
	$1.2 * \bar{r}_{std}$	1896	99.91%
	$1.4 * \bar{r}_{std}$	1884	99.91%
	$1.6 * \bar{r}_{std}$	1878	99.91%
$1.8 * \bar{r}_{std}$	1876	99.91%	
β_0	0.0	1928	98.69%
	-0.001	1925	98.68%
	-0.01	1922	98.55%
	-0.1	1857	96.50%
	-0.2	1629	90.00%
	-0.3	845	32.77%
	-0.4	159	5.25%
-0.5	0	0.06%	

with a small price, leading to losing almost all ad impressions and low budget consumption ratio.

4.5 Visualization of Risk Tendency

Risk tendency reflects the risk preference of a rational bidder in given states. We designed two instantaneous to learn risk tendency, one based on expert knowledge and the other based on self-supervised learning. We visualize these two methods' risk tendencies in Figure 4. It can be observed that both approaches have similar trends in mapping states to risk tendency, which implies that our self-supervised learning algorithm aligns well with the expert knowledge. Specifically, risk tendencies are negative for those resource-limited states at the upper left corner, where a bidder has a small budget for a large number of remaining auctions. In such a state, the bidder prefers to bid with conservative prices. On the other hand, risk tendencies change to positive for those resource-rich states at the bottom right corner because the bidder has a sufficient budget for remaining auctions to support more aggressive bid prices.

5. RELATED WORK

User response prediction. User response prediction can be modeled as a probability estimation task, e.g., click-through rate (CTR) [10], conversion rate (CVR) [21]. Click-through rate (CTR) prediction is widely studied to estimate the value of an ad impression. It plays a key role in display advertising [18; 21; 30; 10; 14]. The traditional solutions use linear models [18], gradient boosting decision trees (GBDT) [29] or factorization machines [23] to predict CTR. Deep learning based methods [36; 27; 5; 12; 25; 13; 33; 2;

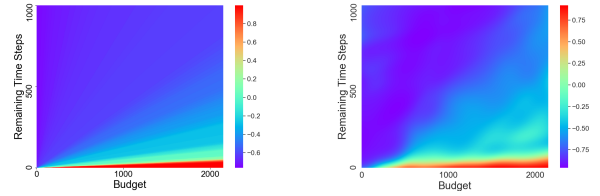


Figure 4: Left: Risk tendency obtained by expert knowledge. Right: Risk tendency learned from self supervision. X-axis and Y-axis represent the remaining budget b and remaining auctions t , respectively.

41; 6; 38] are also used for CTR prediction tasks. Many of them follow the Embedding&MLP paradigm. All these methods aim to only predict the CTR without any information guiding the confidence level. [37] is the only study that simultaneously captures the estimation and risk of pCTR prediction based on Bayesian logistic regression.

Risk Management Techniques. The goal of risk management in finance is to guarantee that the risk does not hurt the business profit [22]. The risk and expected return of each asset can be quantified by the mean return of the asset and corresponding covariance matrix [20; 16]. The mean-variance analysis is adopted to achieve the trade-off between the risk and the expected return [28]. [39] measures the campaign-level risk and returns in a special case of arbitrage. As for impression-level risk management, [37] proposes a risk-aware bidding strategy based on the value at risk (VaR) with campaign-level profit gain.

Bidding strategy. Truthful bidding is the most fundamental strategy that has been proven to be optimal in second-price auctions with unlimited budgets [17]. In practice, every ad campaign has a budget constraint, where linear [24] and model-based/model-free reinforcement learning-based strategies [3; 31] can be adopted to determine bid prices. Other advanced methods, such as Recurrent Neural Networks [11] optimization [15], are developed for bidding with budget constraints. However, these methods inevitably assume the user response prediction is perfect and ignore the inherent estimation uncertainty. [37] is most relevant to our work, which proposes a risk management algorithm based on value at risk, but ignores intrinsic interactions between the market environment and the state of a bidder. In our work, we propose a risk-aware reinforcement learning based bidding strategy that explicitly considers such interactions.

6. CONCLUSION

In this paper, we investigated the bid optimization problem in RTB with the benefits of risk information. For the bid price determination with budget constraint, we, to the best of our knowledge, firstly consider both estimation uncertainty and the dynamic risk tendency. Specifically, we first theoretically analyze the relation between prediction uncertainty and the risk tendency of a bidder, and then proposed an adaptive risk-aware bidding algorithm with budget constraint. Subsequently, we developed two instantiations to determine risk tendency based on expert knowledge or self-supervised learning. Experimental results on real datasets demonstrate that RRLB making use of both prediction uncertainty and risk tendency achieves better cumulative performance than representative competitors.

7. REFERENCES

- [1] S. A. Armstrong. *A meta-analysis of randomness in human behavioral research*. Louisiana State University and Agricultural & Mechanical College, 2004.
- [2] N. Bhamidipati, R. Kant, and S. Mishra. A large scale prediction engine for app install clicks and conversions. In *Proceedings of the 2017 ACM on Conference on Information and Knowledge Management*, pages 167–175, 2017.
- [3] H. Cai, K. Ren, W. Zhang, K. Malialis, J. Wang, Y. Yu, and D. Guo. Real-time bidding by reinforcement learning in display advertising. In *Proceedings of the Tenth ACM International Conference on Web Search and Data Mining*, pages 661–670, 2017.
- [4] O. Chapelle. Modeling delayed feedback in display advertising. In *Proceedings of the 20th ACM SIGKDD international conference on Knowledge discovery and data mining*, pages 1097–1105, 2014.
- [5] H.-T. Cheng, L. Koc, J. Harmsen, T. Shaked, T. Chandra, H. Aradhye, G. Anderson, G. Corrado, W. Chai, M. Ispir, et al. Wide & deep learning for recommender systems. In *Proceedings of the 1st Workshop on Deep Learning for Recommender Systems*, pages 7–10, 2016.
- [6] Y. Feng, F. Lv, W. Shen, M. Wang, F. Sun, Y. Zhu, and K. Yang. Deep session interest network for click-through rate prediction. In *Proceedings of the 28th International Joint Conference on Artificial Intelligence*, pages 2301–2307, 2019.
- [7] A. Ghosh, S. Mitra, S. Sarkhel, and V. Swaminathan. Optimal bidding strategy without exploration in real-time bidding. In *Proceedings of the 2020 SIAM International Conference on Data Mining*, pages 298–306. SIAM, 2020.
- [8] A. Ghosh, B. I. Rubinstein, S. Vassilvitskii, and M. Zinkevich. Adaptive bidding for display advertising. In *Proceedings of the 18th international conference on World wide web*, pages 251–260, 2009.
- [9] Google. The arrival of real-time bidding and what it means for media buyers. In *Google White Paper*, 2012.
- [10] T. Graepel, J. Q. Candela, T. Borchert, and R. Herbrich. Web-scale bayesian click-through rate prediction for sponsored search advertising in microsoft’s bing search engine. Omnipress, 2010.
- [11] N. Grislain, N. Perrin, and A. Thabault. Recurrent neural networks for stochastic control in real-time bidding. In *Proceedings of the 25th ACM SIGKDD International Conference on Knowledge Discovery & Data Mining*, pages 2801–2809, 2019.
- [12] H. Guo, R. Tang, Y. Ye, Z. Li, and X. He. Deepfm: a factorization-machine based neural network for ctr prediction. In *Proceedings of the 26th International Joint Conference on Artificial Intelligence*, pages 1725–1731, 2017.
- [13] X. He and T.-S. Chua. Neural factorization machines for sparse predictive analytics. In *Proceedings of the 40th International ACM SIGIR conference on Research and Development in Information Retrieval*, pages 355–364, 2017.
- [14] X. He, J. Pan, O. Jin, T. Xu, B. Liu, T. Xu, Y. Shi, A. Atallah, R. Herbrich, S. Bowers, et al. Practical lessons from predicting clicks on ads at facebook. In *Proceedings of the Eighth International Workshop on Data Mining for Online Advertising*, pages 1–9, 2014.
- [15] Y. He, X. Chen, D. Wu, J. Pan, Q. Tan, C. Yu, J. Xu, and X. Zhu. A unified solution to constrained bidding in online display advertising. In *Proceedings of the 27th ACM SIGKDD Conference on Knowledge Discovery & Data Mining*, pages 2993–3001, 2021.
- [16] J. Hull. *Risk management and financial institutions*, volume 733. John Wiley & Sons, 2012.
- [17] V. Krishna. *Auction theory*. Academic press, 2009.
- [18] K.-c. Lee, B. Orten, A. Dasdan, and W. Li. Estimating conversion rate in display advertising from past performance data. In *Proceedings of the 18th ACM SIGKDD International Conference on Knowledge Discovery and Data Mining*, pages 768–776, 2012.
- [19] T. J. Linsmeier and N. D. Pearson. Value at risk. *Financial Analysts Journal*, 56(2):47–67, 2000.
- [20] H. Markowitz. Harry m. markowitz. *Portfolio selection, Journal of Finance*, 7(1):77–91, 1952.
- [21] H. B. McMahan, G. Holt, D. Sculley, M. Young, D. Ebner, J. Grady, L. Nie, T. Phillips, E. Davydov, D. Golovin, et al. Ad click prediction: a view from the trenches. In *Proceedings of the 19th ACM SIGKDD International Conference on Knowledge Discovery and Data Mining*, pages 1222–1230, 2013.
- [22] J. Mun. *Modeling risk: Applying Monte Carlo simulation, real options analysis, forecasting, and optimization techniques*, volume 347. John Wiley & Sons, 2006.
- [23] R. J. Oentaryo, E.-P. Lim, J.-W. Low, D. Lo, and M. Finegold. Predicting response in mobile advertising with hierarchical importance-aware factorization machine. In *Proceedings of the 7th ACM International Conference on Web Search and Data Mining*, pages 123–132, 2014.
- [24] C. Perlich, B. Dalessandro, R. Hook, O. Stitelman, T. Raeder, and F. Provost. Bid optimizing and inventory scoring in targeted online advertising. In *Proceedings of the 18th ACM SIGKDD International Conference on Knowledge Discovery and Data Mining*, pages 804–812, 2012.
- [25] S. Rendle. Factorization machines with libfm. *ACM Transactions on Intelligent Systems and Technology (TIST)*, 3(3):1–22, 2012.
- [26] R. T. Rockafellar, S. Uryasev, et al. Optimization of conditional value-at-risk. *Journal of risk*, 2:21–42, 2000.

- [27] Y. Shan, T. R. Hoens, J. Jiao, H. Wang, D. Yu, and J. Mao. Deep crossing: Web-scale modeling without manually crafted combinatorial features. In *Proceedings of the 22nd ACM SIGKDD International Conference on Knowledge Discovery and Data Mining*, pages 255–262, 2016.
- [28] W. F. Sharpe. The sharpe ratio. *Journal of Portfolio Management*, 21(1):49–58, 1994.
- [29] I. Trofimov, A. Kornetova, and V. Topinskiy. Using boosted trees for click-through rate prediction for sponsored search. In *Proceedings of the Sixth International Workshop on Data Mining for Online Advertising and Internet Economy*, pages 1–6, 2012.
- [30] J. Wang and S. Yuan. Real-time bidding: A new frontier of computational advertising research. In *Proceedings of the Eighth ACM International Conference on Web Search and Data Mining*, pages 415–416, 2015.
- [31] D. Wu, X. Chen, X. Yang, H. Wang, Q. Tan, X. Zhang, J. Xu, and K. Gai. Budget constrained bidding by model-free reinforcement learning in display advertising. In *Proceedings of the 27th ACM International Conference on Information and Knowledge Management*, pages 1443–1451, 2018.
- [32] W. C.-H. Wu, M.-Y. Yeh, and M.-S. Chen. Predicting winning price in real time bidding with censored data. In *Proceedings of the 21th ACM SIGKDD International Conference on Knowledge Discovery and Data Mining*, pages 1305–1314, 2015.
- [33] J. Xiao, H. Ye, X. He, H. Zhang, F. Wu, and T.-S. Chua. Attentional factorization machines: learning the weight of feature interactions via attention networks. In *Proceedings of the 26th International Joint Conference on Artificial Intelligence*, pages 3119–3125, 2017.
- [34] S. Yuan, J. Wang, and X. Zhao. Real-time bidding for online advertising: measurement and analysis. In *Proceedings of the Seventh International Workshop on Data Mining for Online Advertising*, pages 1–8, 2013.
- [35] D. Zha, K.-H. Lai, K. Zhou, and X. Hu. Experience replay optimization. In *In Proceedings of the 2019 International Joint Conference on Artificial Intelligence*, page 4243–4249, 2019.
- [36] S. Zhai, K.-h. Chang, R. Zhang, and Z. M. Zhang. Deepintent: Learning attentions for online advertising with recurrent neural networks. In *Proceedings of the 22nd ACM SIGKDD International Conference on Knowledge Discovery and Data Mining*, pages 1295–1304, 2016.
- [37] H. Zhang, W. Zhang, Y. Rong, K. Ren, W. Li, and J. Wang. Managing risk of bidding in display advertising. In *Proceedings of the Tenth ACM International Conference on Web Search and Data Mining*, pages 581–590, 2017.
- [38] W. Zhang, T. Du, and J. Wang. Deep learning over multi-field categorical data. In *European Conference on Information Retrieval*, pages 45–57, 2016.
- [39] W. Zhang and J. Wang. Statistical arbitrage mining for display advertising. In *Proceedings of the 21th ACM SIGKDD International Conference on Knowledge Discovery and Data Mining*, pages 1465–1474, 2015.
- [40] W. Zhang, S. Yuan, and J. Wang. Optimal real-time bidding for display advertising. In *Proceedings of the 20th ACM SIGKDD International Conference on Knowledge Discovery and Data Mining*, pages 1077–1086, 2014.
- [41] G. Zhou, X. Zhu, C. Song, Y. Fan, H. Zhu, X. Ma, Y. Yan, J. Jin, H. Li, and K. Gai. Deep interest network for click-through rate prediction. In *Proceedings of the 24th ACM SIGKDD International Conference on Knowledge Discovery & Data Mining*, pages 1059–1068, 2018.