

## TABLE OF CONTENTS

### **KDD 2021 Highlights**

- 1 An Interview with Dr. Shipeng Yu, Winner of ACM SIGKDD 2021 Service Award

### **Contributed Articles**

- 3 Surf or Sleep? Understanding the Influence of Bedtime Patterns on Campus: Teng Guo, Linhong Li, Dongyu Zhang, and Feng Xia
- 13 Meta-Learning with Graph Neural Networks: Methods and Applications: Debmalya Mandal, Sourav Medya, Brian Uzzi, and Charu Aggarwal



**Association for  
Computing Machinery**

*Advancing Computing as a Science & Profession*

**Editor-in-Chief:**  
Hanghang Tong

**Associate Editors:**  
Xin Luna Dong  
Ankur Teredesai  
Reza Zafarani  
<http://www.kdd.org/explorations/>

**About SIGKDD Explorations**

*Explorations* is published twice yearly, in June/July and December/January each year. After the first two volumes, frequency may increase to quarterly. The newsletter is distributed in hardcopy form to all members of the ACM SIGKDD. It is also sent to ACM's network of libraries. Additionally, issues are published on the web and are free to the general public (<http://www.acm.org/sigkdd/explorations/>).

Our goal is to make *SIGKDD Explorations* an informative, rapid means of publication and a dynamic forum for communication with the Knowledge Discovery and Data Mining community. SIGKDD membership is growing at a very fast pace, and with KDD being a multi-disciplinary field, we hope that *Explorations* will facilitate its fusion and enhance the sense of community. Submissions will be reviewed by the editor and/or associate and guest editors as appropriate. We are particularly interested in short research and survey articles on various aspects of data mining and KDD. *Explorations* is also a forum for publishing position papers, controversial positions, challenges to the community, product reviews, book reviews, news items and other items of interest to the field. Please see:

<http://www.acm.org/sigkdd/explorations/instructions.htm>

**Advertiser Information:**

*Explorations* accepts advertisements related to data mining and KDD, including company, book, vendor, and service advertisements. For rates and instructions on submitting an ad, please see:

<http://www.acm.org/sigkdd/explorations/instructions.htm#advertise>

**Notice to Contributing Authors to SIG Newsletters:**

By submitting your article for distribution in this Special Interest Group publication, you hereby grant to ACM the following non-exclusive, perpetual, worldwide rights:

- to publish in print on condition of acceptance by the editor
- to digitize and post your article in the electronic version of this publication
- to include the article in the ACM Digital Library and in any Digital Library related services
- to allow users to make a personal copy of the article for noncommercial, educational or research purposes.

However, as a contributing author, you retain copyright to your article and ACM will refer requests for republication directly to you.

# An Interview with Dr. Shipeng Yu, Winner of ACM SIGKDD 2021 Service Award

## ABSTRACT

Shipeng Yu, Ph.D. is the recipient of the 2021 ACM SIGKDD Service Award, which is the highest service award in the field of knowledge discovery and data mining. Conferred annually on one individual or group in recognition of outstanding professional services and contributions to the field of knowledge discovery and data mining, Dr. Yu was honored for his years of service and many accomplishments as general chair of KDD 2017 and currently as sponsorship director for SIGKDD. Dr. Yu is Director of AI Engineering, Head of the Growth AI team at LinkedIn, the world's largest professional network. He sat down with SIGKDD Explorations to discuss how he first got involved in the KDD conference in 2006, the benefits and drawbacks of virtual conferences, his work at LinkedIn, and KDD's place in the field of machine learning, data science and artificial intelligence.

## CONGRATULATIONS ON RECEIVING THE SIGKDD SERVICE AWARD FOR YOUR CONTRIBUTIONS AS GENERAL CHAIR AND SPONSORSHIP DIRECTOR. PLEASE TELL US ABOUT YOURSELF. HOW DID YOU GET INTO THE FIELD OF DATA SCIENCE, AI AND MACHINE LEARNING?

I was born and raised in China. I became interested in math early on, and studied the subject in undergraduate and graduate schools. I pursued my doctoral work at the University of Munich, where I had the good fortune to study under renowned Data Mining and AI experts Professor Hans-Peter Kriegel and Dr. Volker Tresp. While in Germany, I had the opportunity to work at Siemens Corporate Research, exploring broad applications of general machine learning, though much of my work focused on statistical machine learning. Given my experience in machine learning and data mining at a large scale, after completing my Ph.D., I accepted a role at Siemens Healthcare in the U.S. Over a span of eight years in the healthcare sector, I was introduced to many different kinds of medical data: image data, clinical data, genetic data. The goal was to improve diagnoses and personalize treatment through data mining.

## AS THE HEAD OF GROWTH AI AT LINKEDIN, CAN YOU SHARE WHAT THAT ROLE LOOKS LIKE?

While I appreciated the very real impact—you are literally saving lives—of working with medical data, regulatory and other challenges are prevalent in the healthcare industry. You are also dealing with small data in most cases. My move to LinkedIn in early 2015 was prompted by a desire to apply the tenets of data mining in a new environment and a new industry. The Growth AI team at LinkedIn has broad-based responsibility, including network and engagement growth, as well as retention of the platform's 800 million members. We are primarily concerned with building the

right network and getting the right communications to the right members at the right time. Data mining and machine learning are essential tools to optimize these growth charters. I still do research and publish papers, but that work is also product focused, with a goal to improve business functions and strategies.

## AS SOMEONE WITH EXPERIENCE IN THE PRIVATE SECTOR AND ACADEMIC RESEARCH, WHERE CAN ACADEMIA AND INDUSTRY HELP ONE ANOTHER?

With experience spanning both academia and industry, I can say there are definitely some not-so-subtle differences between the two. Without the rigid pressure of metrics or budgets, universities tend to focus on fundamental research—things that are longer-term, more groundbreaking. Access to government funding and a pool of student researchers also contributes to the success and nature of academic research. On the other hand, industry has customers, partners, shareholders and other stakeholders to satisfy. As a result, businesses tend to focus research on solving fascinating real-world problems, often with a more aggressive timeline. While perhaps different in approach or context, I do believe there are problems common to both academia and industry, that may be tackled with data science. Because KDD spans both worlds, the organization and its many volunteers are ideally suited to help improve collaboration between higher education and the business sector, making the field of data science more equitable and diverse.

## WHEN DID YOU FIRST DISCOVER KDD AND HOW HAS YOUR RELATIONSHIP WITH THE CONFERENCE EVOLVED?

KDD came on my radar in my early years at graduate school when some of my work was published at KDD conferences. I really got my foot in the door when I volunteered to be a reviewer at the first conference I attended, KDD 2006. That's when I began to appreciate the international community of data mining and knowledge discovery researchers, and offered to help facilitate conference workshops. Between 2010 and 2016, my involvement grew steadily, and I was thrilled to be part of planning the first KDD conference in a major U.S. city, KDD 2016 in San Francisco. A resounding success, that year there were more than 2,700 attendees from 88 countries.

I then had the honor and privilege to serve as general chair for KDD 2017, which took place in Halifax, Nova Scotia. A fairly remote location and last-minute venue change posed challenges from an organizing perspective. However, building on the success of the prior year and adopting many of the same tactics to recruit volunteers and attract registrants, we were pleased to draw more than 1,700 attendees. Since serving as the sponsorship chair for the 2018 conference in London, I have taken on responsibility for leading SIGKDD's sponsorship efforts.

## **AS SPONSORSHIP DIRECTOR FOR SIGKDD, WHAT ARE YOUR OBJECTIVES AND HOW DO YOU MEET THEM?**

Sponsorships have been key to the longevity and success of the conference, with generous financial support ensuring KDD remains home to some of the most highly cited research papers in data science. While previous sponsorship chairs served one-year stints, I believe continued growth in this area requires leadership continuity, so have been pleased to serve in this role for three years. With the organization surpassing \$1 million in annual sponsorship—\$1.2 million in 2018 and \$1.1 million in 2019, my goal is to continue and accelerate this trajectory.

I am proud that we have secured interest from both government and corporate sponsors over my tenure. This healthy mix of support has enabled us to set sponsorship records. The switch from a physical package and to a virtual package necessitated by the Coronavirus pandemic, coupled with resulting economic uncertainty, have posed some sponsorship struggles. Fortunately, though, we have been building long-term relationships with many of our sponsors, and, as a result, financial support for our virtual conference improved in 2021.

## **YOU HAVE HELPED ORGANIZE MULTIPLE IN-PERSON CONFERENCES. HOW HAVE THE PANDEMIC AND THE SHIFT TO VIRTUAL CONFERENCES CHANGED THE EXPERIENCE FOR ATTENDEES?**

The coronavirus pandemic has affected us in ways we never imagined. Every organization has had to adapt to the new normal. Prior to the public health crisis, KDD thrived as a vibrant, physical, in-person event. The conference was at its best when people were talking and collaborating in small group discussions. We have been forced to tackle head on the uncertainty around how we can recreate

that experience online, how sponsors can reach members in this format, and which offerings will and won't work virtually. At the same time, virtual conferences have opened up a lot of opportunities for attendees to interact with speakers and authors. It also democratizes the conference, making it available for everyone, as registration fees have been reduced and travel isn't required. That's why I believe there is value in the hybrid format. Yes, it had been done in the past, but it wasn't mainstream. Now people can see that a combination of virtual and in-person conference setup not only works but expands opportunities for people to get together in a way that makes sense for them.

## **WHY DO YOU THINK IT IS IMPORTANT TO VOLUNTEER FOR SIGKDD?**

KDD is special not only because of the groundbreaking discoveries presented at the show, but also the connections and bonds we are able to form with our colleagues. I have been extremely fortunate to serve as general chair as well as my current role. Being the sponsorship director at KDD is not just about raising money. It's about raising awareness around our field and the phenomenal international community that we have built and exponentially grown over the past few years. Volunteering creates new and expanded opportunities to sustain and strengthen valuable personal and professional relationships, and give back to an organization at the forefront of advancing our industry.

## **LASTLY, WHAT DOES THE SERVICE AWARD WIN MEAN TO YOU?**

I cannot begin to express my gratitude to everyone who has helped contribute to KDD's growth and success over the past two decades. I am honored to join such an accomplished list of professionals and servants who have won this award, including last year's winner, Dr. Michael Zeller. I encourage everyone in the data mining and knowledge discovery community to join us next year at KDD 2022.

# Surf or sleep? Understanding the influence of bedtime patterns on campus

Teng Guo, Linhong Li, Dongyu Zhang  
School of Software,  
Dalian University of Technology,  
Dalian 116620, China  
teng.guo@outlook.com;  
li.linhong@outlook.com;  
zhangdongyu@dlut.edu.cn

Feng Xia\*  
School of Engineering, IT and Physical Sciences,  
Federation University Australia,  
Ballarat, VIC 3353, Australia  
f.xia@acm.org

## ABSTRACT

Poor sleep habits may cause serious problems of mind and body, and it is a commonly observed issue for college students due to study workload as well as peer and social influence. Understanding its impact and identifying students with poor sleep habits matters a lot in educational management. Most of the current researches is either based on self-reports and questionnaires, suffering from small sample size and social desirability bias, or the methods used are not suitable for the education system. In this paper, we develop a general data-driven method for identifying students' sleep patterns according to their Internet access pattern stored in the education management system and explore its influence from various aspects. First, we design a Poisson-based probabilistic mixture model to cluster students according to the distribution of bedtime and identify students who are used to stay up late. Second, we profile students from five aspects (including eight dimensions) based on campus-behavior data and build Bayesian networks to explore the relationship between behavioral characteristics and sleeping habits. Finally, we test the predictability of sleeping habits. This paper not only contributes to the understanding of student sleep from a cognitive and behavioral perspective but also presents a new approach that provides an effective framework for various educational institutions to detect the sleeping patterns of students.

## 1. INTRODUCTION

Sleep matters a lot for health and well-being [38; 44; 4]. Poor sleep habits like insomnia can easily cause serious problems like worsened health-related quality of life and poor health outcomes, especially for college students with psychological immaturity [13]. Research have shown that the majority of college students suffer from sleep disorders like feeling tired, fatigue, or daytime sleepiness and sleeping less than eight hours per night during the school week [6; 27]. 68% of American teenagers report that their sleep time each night is below the recommended eight-ten hours [15] and in East Asia, the sleep time of adolescents on weekdays night is below six hours [29; 14]. Therefore, a significant topic of educational research is identifying college students with bad sleep

habits timely [40; 3], and understanding its influence. However, detecting sleep patterns accurately faces tremendous challenges. Previous research are mainly based on questionnaires, which is time-consuming and costly making, and it is hard to be scale to a lot of students. In the last decade, scientists propose a variety of techniques for sleep tracking, and some related products are already in use [2; 32]. These techniques can only be used for small-scale analysis due to the high cost of the associated equipment. Some studies have tried to do large-scale sleep-related research through a data-driven approach like [4]. However, such data is difficult to access except for large companies like Microsoft, Google, etc. A data-driven framework for sleep habit detection and analysis that can serve the daily management of universities is urgently needed.

During the past decade, we have witnessed an increased utilization of electronic devices like smartphones and tablets [7]. As they become more lightweight, people may easily use these devices even in bed. Electronic media use in bed after bedtime is more common in younger compared to older age groups [8]. Research show that more than 90% of young people are used to playing mobile phones or tablets before bedtime [21; 19]. Meanwhile, for college students, the campus network is a good choice because of its low price and faster internet speed, which generates tons of data of internet access logs stored in the local education management system. This provides university management divisions a new way to detect the sleep patterns of college students. However, new challenges are introduced as well. First, these data are not easy to obtain due to privacy protection [23; 36]. Second, the volume of data collected in this way is much larger than the traditional method. Targeted methods need to be designed for efficient data mining.

In this paper, we are devoted to identifying students with the habit of staying up late and explore the influence of this habit. First, we use the Internet access pattern to identify the bedtime of each student. Second, we design a Poisson-based probabilistic mixture model to cluster students according to the distribution of bedtime and identify students who stay up late. Third, besides sleep habit, we profile students from five aspects in eight dimensions, including interest (reading status, app preference status, and surfing length status), orderliness (breakfast orderliness status and bath orderliness status), finance (financial status), academic performance (academic performance status) and gender. Then we build Bayesian networks to explore the re-

\*Corresponding author

relationship between these characteristics and sleeping habits. Finally, we design an experiment of predictive inference to estimate the predictability of sleep status.

Our contributions could be summarized as follows:

- We design a Poisson-based probabilistic mixture model suitable for education management divisions to identify students' sleep patterns.
- We profile students behavior from various aspects and build Bayesian networks to explore the relationship between these characteristics and sleeping habits.
- We conduct comprehensive experiments on a large-scale educational dataset and the extensive results demonstrate the effectiveness.

This paper is organized as follows. In the next section, related work is reviewed. The methods are presented in Section 3. In Section 4, we introduce the experiment setting and analyze the results. In the final section, we present the discussion and conclusion of our work.

## 2. RELATED WORK

The sleep patterns of college students have attracted the attention of scholars from various fields. As the negative impact of poor sleep quality has been demonstrated, the main concern of scientists is the factors that affect sleep quality. Buboltz et al. [10] carry out a study focused on the sleep disturbance of college students. The results show that most college students suffer from sleep disturbance and women exhibit more sleep disturbance than men do. Moreover, they give some suggestions to reduce the negative influence on academic performance. Tsai et al. [42] also focus on how gender can impact students' sleep. Brown et al. [9] use regression modeling to capture the relationship between sleep hygiene and sleep practices. The results demonstrate that variable sleep schedules, going to bed thirsty, environmental noise, and worrying can all impact the quality of sleep. Tavernier et al. [41] explore how the social tie can impact sleep and the results show that positive social ties can contribute to good sleep quality. Patrick et al. [35] focus on the influence of energy drinks and binge drinks on sleep quality. The results show that drinking the above two drinks can easily lead to poor sleep quality and fatigue the next day. Meanwhile, Schneider et al. [37] focus on adolescents and explore the relationship between sleep characteristics and body-mass index. In this study, they emphasize the importance of gender in related research. Orzech et al. [33] proves that longer digital media use two hours prior to bedtime leads to reduced sleep and later bedtime. They obtain an interesting conclusion that more diverse media uses are associated with longer sleep durations. Moreover, some scholars have focused on sleep patterns and their impact on daily life. Kelley et al. [27] explain adolescents' fatigue and irritability by the sleep problems caused by the disorder of their biological time and social time. Althoff et al. [4] carry out a convincing experiment based on a search engine Bing from Microsoft including 3 million nights of sleep and 75 million interaction tasks. The results show how sleep deprivation can seriously affect daily performance. However, as mentioned before, such type of data is inaccessible for education management divisions.

Thanks to the development of technologies, more and more electronic devices are appearing in our daily lives [11] and scientists try to use these advanced technologies (e.g. remote bio-signal monitoring technology) to solve sleep problems. Vroon et al. [43] propose an actuated robotic pillow to improve the quality of sleep through enhancing the interaction with users. De Arriba et al. [16] propose a smartphone-based indicator that can calculate sleep patterns automatically. They try to use this indicator to facilitate the construction of software services in order to improve the efficiency of the learning process. Liu et al. [28] propose a mobile sleep-management system integrated with self-regulated learning strategies and cognitive behavioral therapy. This system can help students to improve their sleep quality by improving their daily schedule and modify strategies to cultivate good learning and healthy lifestyle habits. Zhao et al. [45] design a smart-home system based on the Internet of things, which realizes the function that judges the user's mood and automatically plays corresponding music through facial expression recognition. However, the methods mentioned above are difficult to be used for large-scale sleep detection on campus due to their high cost.

## 3. METHODS

In this section, we describe the models used in the experiment. First, we introduce a mixed probability model used to cluster students' sleep habits. Then, we introduce the Bayesian network model used to analyze the influence of different sleeping habits.

### 3.1 Sleep Pattern Recognition

#### 3.1.1 Aggregated Sleep Count

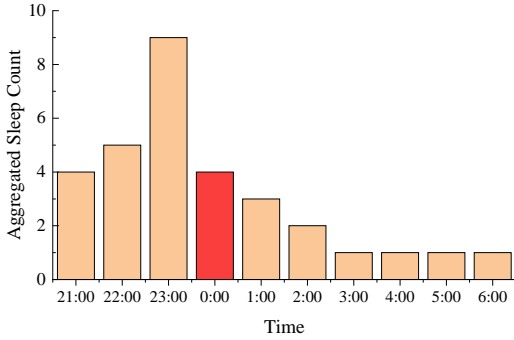
For universities that have campus-dedicated networks deployed, web logs can be used to calculate the bedtime of students every night based on the assumption that students will access the Internet through mobile phones or computers before bedtime [5; 21; 19]. In other words, we can use the time of the last network signal every day to quantify the time students go to bed. To mine sleep patterns of different students, we define *aggregated sleep count*, which represents the total number of days a student goes to sleep at a particular time slot. For example, in a certain month, daily bedtime is shown in Figure 1(a). Let's take 0:00 as an example. In this month, there are 4 days that sleep at 0:00 o'clock shown in Figure 1(a) with the red circle, the aggregated sleep count of 0:00 in this month is 4 shown in Figure 1(b) with a red bar. The *aggregated sleep count* of each time period is shown in Figure 1(b).

#### 3.1.2 Mixture Model with Gamma Priors

Inspired by [34], we design a Poisson mixture model to mine the pattern on their aggregated sleep counts in each time period. For student  $i$ , we let  $s_i$  be the vector of aggregated sleep counts, where  $i = 1, 2, 3 \dots N$ . The dimensionality  $D$  of each vector is the number of time windows ( $D = 16$  in this case that we divided the time into 16 segments from 9:00 pm to 4:30 am and the time span of each segment is 30 minutes because: 1, the last class generally is over at 9:00 pm. 2, the sun rises at around 4:30 am). Thus, our data consists of  $N$  students each with  $D$ -dimensional vectors of aggregated sleep counts.

Sun	Mon	Tue	Wed	Thu	Fri	Sat
	1 21:00	2 21:00	3 23:00	4 23:00	5 23:00	6 23:00
7 0:00	8 21:00	9 21:00	10 0:00	11 1:00	12 1:00	13 23:00
14 0:00	15 22:00	16 22:00	17 22:00	18 2:00	19 22:00	20 6:00
21 23:00	22 3:00	23 4:00	24 2:00	25 0:00	26 23:00	27 1:00
28 23:00	29 23:00	30 5:00	31 22:00			

(a) Bedtime of a student in a month: this is the calendar of a certain month and the number in the middle of the square is the bedtime of the day.



(b) Aggregated Sleep Count of a student in a month: this is the aggregated sleep count of each time period during the month.

Figure 1: Illustration for aggregated sleep count.

To mine the pattern behind this data, we build a probabilistic mixture model with Poisson components. In terms of notation, we let  $M$  be the number of components with an index  $m = 1, 2, 3, \dots, M$ . The unobserved latent variable  $z_i$  takes values from the set  $\{1, 2, \dots, M\}$  that corresponds to the latent component of student  $i$ . Each component consists of a vector of Poisson rate parameters,  $\lambda_m = [\lambda_{m1}, \dots, \lambda_{md}, \dots, \lambda_{m16}]$ , where  $d$  from  $[1, 2, 3, \dots, D]$  represents a time window mentioned before. For example, a low value for  $\lambda_{md}$  represents that students rarely sleep during time window  $d$ . On the contrary, a high value for  $\lambda_{md}$  represents that students always sleep during time window  $d$ .

In this study, the Bayesian approach is used to fit our mixture model. To encourage the model to avoid degenerate solutions, we use Gamma prior distribution for the rate parameter  $\lambda_{md}$  according to the conjugate prior theory. More precisely, we use hyper parameters  $\alpha = 1.1$  and  $\beta = 0.1$  for the Gamma distribution, which makes it a step-like function that puts zero probability mass at  $\lambda_{md}$  and provides a relatively flat uninformative prior distribution over positive rate parameter values [34].

According to the total probability formula, the probability of each student  $i$  under this mixture model is:

$$p(\mathbf{s}_i | \lambda) = \sum_{m=1}^M p(\mathbf{s}_i | z_i = m, \lambda_m) p(z_i = m) \quad (1)$$

where  $p(z_i = m)$  is the marginal mixing weight for each component. Assuming conditional independence of aggregated sleep count  $s_{md}$  given component  $m$ , the probability distribution of each component can be decomposed as:

$$p(\mathbf{s}_i | z_i = m, \lambda_m) = \prod_{d=1}^D p(s_{md} | \lambda_{md}, z_i = m) \quad (2)$$

### 3.1.3 Parameter Estimation

Expectation-Maximization (EM) algorithm, which is widely used in fitting a mixture model to data [17; 31], is used in this experiment to estimate the parameter  $\lambda_m$ . More specifically, the optimization objective in this experiment is to maximize the product of the data likelihood times the prior probability. This optimization process produces two results. The first one is maximum a posteriori parameter estimates for the Poisson components in the model and the second one is membership weights  $\omega_{im}$  that is the probability that student  $i$  belongs to component  $m$ . Each iteration of the EM algorithm consists of two processes, the E (expectation) process and the M (maximization) process. In the E process, the probability of membership is computed for each component  $m = 1, 2, \dots, M$ , for each student  $i = 1, 2, \dots, N$ .

$$\begin{aligned} \omega_{im} &= p(z_i = m | \mathbf{s}_i, \lambda, \alpha, \beta) \\ &\propto p(z_i = m, \mathbf{s}_i, \lambda_m | \alpha, \beta) \\ &\propto p(\mathbf{s}_i | z_i = m, \lambda_m) p(\lambda_m | \alpha, \beta) p(z_i = m) \end{aligned} \quad (3)$$

where  $\omega_{im}$  is the probability that student  $i$  belongs to component  $m$ . In the M process, conditioned on the set of membership probabilities  $\omega_{im}$ , we estimate each parameter via MAP (Maximum Posteriori) estimation:

$$\hat{\lambda}_m = \frac{\sum_i \omega_{im} (\mathbf{s}_i + \alpha - 1)}{\sum_i \omega_{im} (\beta + 1)} \quad (4)$$

$$\hat{p}(z_i = m) = \frac{\sum_i \omega_{im}}{N} \quad (5)$$

The outputs of the M process provide the input for the next E process, and thus, the cycle of E and M processes continue iteratively.

## 3.2 Bayesian Network

In this study, Bayesian network (BN) is used to explore the relevance between bedtime patterns and behavioral characteristics [25; 20]. Based on Markov property, the BN uses a set of local distributions to express the global joint distribution [1; 12].

### 3.2.1 Bayesian network

A BN consists of two parts: a directed acyclic graph (DAG) and a set of conditional probability distribution, which can be represented as  $\mathcal{B} = \langle G_d, \Theta \rangle$ . The graph  $G_d$  represents a directed acyclic graph whose vertices correspond to random variables  $X_1, X_2, \dots, X_n$  and whose edges represent direct dependencies between the variables, where  $n$  represent the number of variable in  $\mathcal{B}$ . The  $\Theta = (\theta_i)_{1 \leq i \leq n}$  denotes the set of all parameters.

The notations we used are introduced as follows:

- $n$  represent the number of variable in  $\mathcal{B}$  and each node  $X_i$  takes  $r$  possible values  $x_1, x_2, \dots, x_r$
- $\pi(X_i)$ , denoted the parent nodes of each variable  $X_i$  in  $G_d$ , and has  $q_i = \prod_{m \in \Omega_i} r_m$  possible configurations, where  $\Omega_i = \{m; X_m \in \pi(X_i)\}$

- $ch(X_i)$  denotes the children nodes of each variable  $X_i$  in  $G_d$
- The conditional distribution of  $X_i|\pi(X_i)$  is defined by the matrix of probabilities  $\theta_i = (\theta_{ijk})_{1 \leq j \leq q_i, k \in (x_1, \dots, x_r)}$ :

$$\theta_{ijk} = P(X_i = k | \pi(X_i) = j) \quad (6)$$

For each  $X_i$ ,  $\theta_{ij} = (\theta_{ijx_1}, \theta_{ijx_2}, \dots, \theta_{ijx_r})$  is the vector of conditional probabilities defining the distribution of  $X_i | \pi(X_i) = j$ , where  $\sum_{k=1}^r \theta_{ijk} = 1$ . The joint probability distribution of  $(X_1, X_2, \dots, X_i)$  is given by:

$$P(X_1, X_2, \dots, X_i) = \prod_i^d P(X_i | \pi(X_i)) \quad (7)$$

### 3.2.2 Structure Design

To explore the main problem proposed in this paper mentioned above, we first define nine binary variables used in our BN shown in Table 1 (Details are introduced in the following section). We use student data for one semester in this paper. Based on the causal Markov assumption [39], we develop a three-layers raw structure, shown in Figure 2, with two rules: (1), the directed edge can start from one node to another node on the same layer. (2), the directed edge can only point from the lower-layer to the upper level. In other words, we blacklist all outgoing edges from the upper-level nodes to lower level nodes. By construction, no node is allowed to be the parents of gender, and no node is allowed to be the child of academic performance. Such design is consistent with common sense.

### 3.2.3 Structure Learning

In this research, we use BDeu (Bayesian Dirichlet equivalent uniform) score to evaluate the BN and use the hill-climbing strategy to optimize this score [24]. BDeu score is a variant of the BDeu score, which aims at maximizing the posterior probability of the DAG. We use MLE (maximum likelihood estimation) to fit our BNs [22] in this research.

## 4. EXPERIMENTS AND RESULTS

### 4.1 Dataset

The dataset used in this research includes 5,200 students from a Chinese university including freshmen students, sophomores, and juniors and contains a total of about 50 million records of data related to various behavior. For the privacy concern, the students are already pseudonymous in the raw data. Removing the student with incomplete data, we

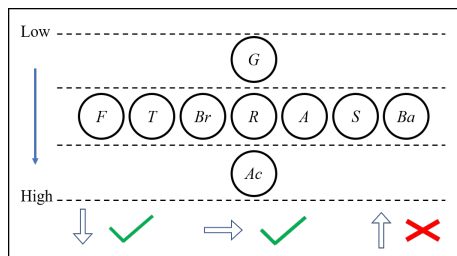


Figure 2: Three-layers raw structure for BN.

Table 2: The number of student in each cluster.

	stay-up	non-stay-up	Total
Junior	664	759	1423
Sophomore	807	773	1580
Freshman	714	532	1246
Total	2185	2064	4249

have 4,249 students left classified into three groups: freshman (1,423 students), sophomore (1,580 students), and junior (1,246 students). This dataset consists of five types of data. The datasets are described as follows:

#### 4.1.1 Internet Access Data

For most students, the campus network is a good choice because of its low price and faster internet speed. The related logs accurately record the student's online information. We use this data to profile students' surfing behavior. The experiment time is from November 2018 to April 2019 including 4,541,247 records data. To ensure the validity of the experimental results, we remove students with fewer related records. In this research, surfing length status  $T$ , app preference status  $A$ , and sleep status  $S$  are calculated based on this data.

#### 4.1.2 Academic Performance Data

Students' academic performance is generally be recorded, which contains the grade, credit of each course for each student. The academic performance data includes 1,048,575 records. In this research, we select the most recent grade data from the period of Internet access data to obtain academic performance data  $Ac$ .

#### 4.1.3 Demographic Data

For all schools, students are required to submit personal information at the time of admission, like hometown, gender, and nation. In this research, we take gender  $G$  as a feature.

#### 4.1.4 Financial Data

Campus smart card, in most universities, is used as a recognition tool for his identification. Generally, smart cards can be used for any scene of consumption on campus and thus record tons of data for student consumption behavior, such as bathing and eating. Financial data include 22,176,513 records. In this research, breakfast orderliness status  $Br$ , bath orderliness status  $Ba$ , and financial status  $F$  are obtained based on this data.

#### 4.1.5 Reading Data

Generally, every university has its own library that is the main source of book reading by students. Reading data include 17,209,329 records. Reading status  $R$  in this research is obtained from this data.

## 4.2 Student Clustering

Below, the results of fitting a two-component Poisson mixture model are presented and discussed. Models with more components are explored in this experiment, but the two-components model mainly captures the primary modes of student's sleep patterns and more components models tend to split students into further subgroups without any significant additional insight.

Table 1: Variables definition.

Variable	Definition	Description
$G$	Gender	Male or Female
$R$	Reading status	Classify students according to the number of books borrowed
$A$	App preference status	Define online preferences based on usage time of each app (Game app or Video app)
$T$	Surfing length status	Classify students according to their total length of surfing time
$Br$	Breakfast Orderliness status	Classify students according to the number of times they eat breakfast
$Ba$	Bath Orderliness status	Classify students according to the regularity of bathing
$F$	Financial status	Classify students based on daily average spending
$Ac$	Academic performance status	Classify students based on Grade-Point Average
$S$	Sleep status	Classify students according to whether they stay up late

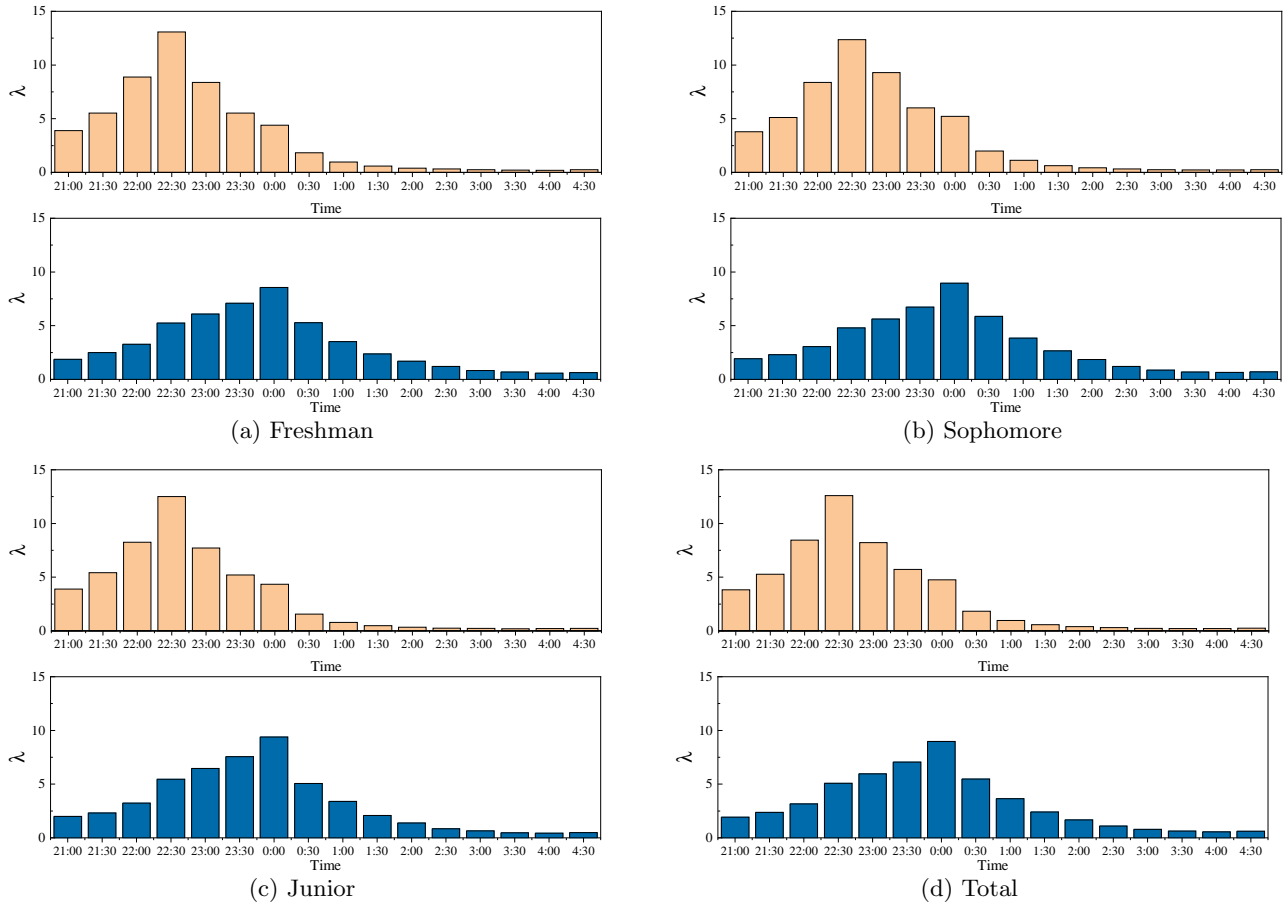


Figure 3: Clustering results of Poisson mixture model with two components.

We set the threshold of component weight at 0.5 to classify student  $i = 1, 2, \dots, N$  into one of the two components, i.e., if  $\omega_{i1} \geq 0.5$ , then student  $i$  is classified into the group of staying up (where  $m = 1$  corresponds to the staying up group). We first experiment individually for freshmen, sophomores, and juniors, and then experiment with all students together. Figure 3 shows the  $\lambda$  of each component and Table 2 show the corresponding number of each components. Two distinct sleep patterns can be seen for each group, as the  $\lambda$  of Poisson distribution represents the number of times things happen per unit of time. Each pattern is an approximate bell curve. They can be distinguished according to two important characteristics shown as follows:

- First, one of the mixture components has a peak at 10:30 pm (shown as the yellow bar chart) and the other one has a peak at 0:00 am (shown as the blue chart).
- Second, the mixture component with 0:00 am peak has a thicker tail compared with the other one with 10:30 pm peak.

Clearly, these two patterns reflect two different types of students: students who stay up later and students who don't stay up late. According to the regulations of the university where we carry out our experiments, all lights in dormitory areas will be turned off at 10:30 to urge students to sleep early, which is consistent with our experimental results, verifying the effectiveness of our proposed method. Moreover, it can be seen in Table 2 that from the freshmen to the juniors, the proportion of students who do not stay up late gradually increased, eventually surpassing the proportion of students staying up late. The possible reason is that they are willing to live a more regular life because of increased academic pressure and forthcoming job-hunting.

### 4.3 Student Profiling

In this section, besides sleep status ( $S$ ) defined above through the mixture probabilistic model, we profile students from other five aspects in eight dimensions, including interest ( $R$ ,  $T$  and  $A$ ), orderliness ( $Br$  and  $Ba$ ), finance ( $F$ ), Academic performance ( $Ac$ ) and Gender ( $G$ ). As shown in Table 1, we sort students according to specific indicators and classify students with 50% as a threshold, resulting in that all variables are binary. The details are introduced in this section. The interest of students is profiled from three aspects: reading habits  $R$ , Interest access time  $T$ , and internet preferences  $A$ . First, for reading habit  $R$ , we divide students into two categories: those who like to read and those who don't, based on the number of books borrowed during the experiment. Second, for Interest access time  $T$ , we divide students into those who like to surf on the Internet and those who don't, based on their average Interest access time. Finally, for app preference status  $A$ , apps in our dataset can be classified into two types: games app and video app according to their functions, and compare the time that users spend on both apps to identify their  $A$ .

The orderliness of students is profiled from two aspects: breakfast orderliness  $Br$  and bath orderliness  $Ba$ . For orderliness profiling, the indicator of  $Br$  is the number of taking breakfasts during the experiment and the indicator of  $Ba$  is the variance of bath date interval during the experiment. Students are classified into two categories according to the corresponding indicators, with a threshold of 50%.

In addition to interest and orderliness, we also profile students from academic performance, financial situation, and gender perspectives. For academic performance, students are divided into two categories based on their Grade-Point Average (GPA). For the financial situation, we divided students into two categories according to their daily spending during the experiment. For gender, students are divided into males and females.

## 4.4 Bayesian Network Analysis

### 4.4.1 Bayesian Network

As mentioned above, our dataset includes three types of students: freshmen, sophomores, and juniors. To capture the rules between different groups, we first build BNs for these three groups separately. In this research, we use a hill-climbing strategy to build BNs through optimizing the BDeu score.

1. First, for each group, we learn 200 BNs through a hill-climbing strategy with restarts strategy, and the results are shown in Figure 4. We learn 200 BNs and 300 BNs, respectively, and find that the score of high-quality BNs in these two cases are equal.
2. Second, we select one-third of the BNs with the high score for the null model to build the consensus network [1]. According to the frequency of edge occurring in these BNs, we design null models with bootstrap sampling to filter the invalid edge (shown in Figure 5) [30]. We determine the threshold by keeping the occurring frequency above the mean value plus two times of standard deviation of the random case and in each group, the value is 25, 27, and 25, respectively.
3. Third, we reconstruct the consensus BNs by those edges that occur more frequently than thresholds.

All BNs are shown in Figure 6. Moreover, we synthesize a total network based on three consensus-BNs according to the occurring frequency of each edge shown in Figure 6(d). For the same edge but different directions, we choose the edge with higher occurrence frequency in high score BNs. For example, in three consensus-BNs, the edge from  $T$  to  $A$  and the edge from  $A$  to  $T$  have the same occurrence frequency. Because the edge from  $A$  to  $T$  has a higher occurrence frequency in BNs with a high score, we keep this edge in total BNs.

### 4.4.2 Structure Analysis

Table 3: Structure analysis for  $S$  in the three BNs.

	$G$	$R$	$A$	$T$	$Br$	$Ba$	$F$	$Ac$
Freshman		✓	✓	✓	✓		✓	✓
Sophomore	✓		✓	✓	✓			✓
Junior			✓	✓				
Times in $\pi(S)$	1		2	2				
Times in $MB(S)$	1	1	3	3	2		1	2
Times in $ch(S)$					2			2

- ✓: The variable is in  $MB(S)$ .
- ✓: The variable is in  $\pi(S)$ .
- ✓: The variable is in  $ch(S)$ .

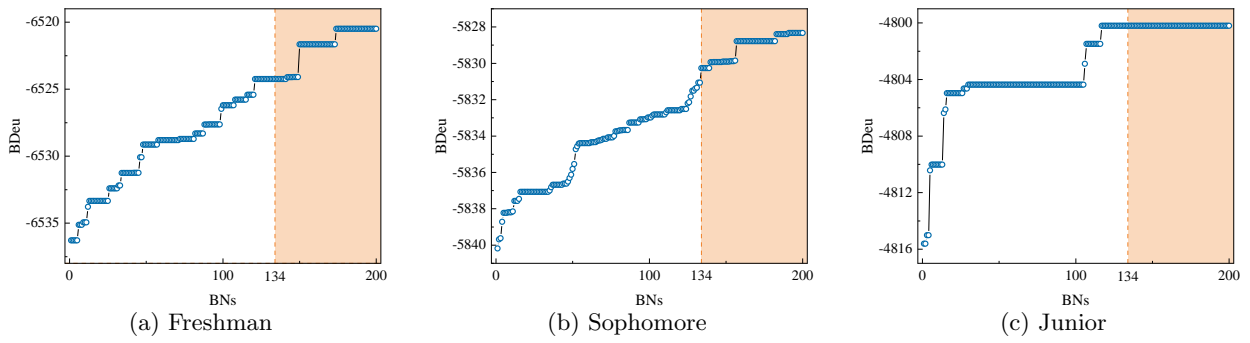


Figure 4: Scores of 200 BNs from restart strategy. We select one-third of the BNs with the high score (shown in the yellow area) for the null model and the consensus network.

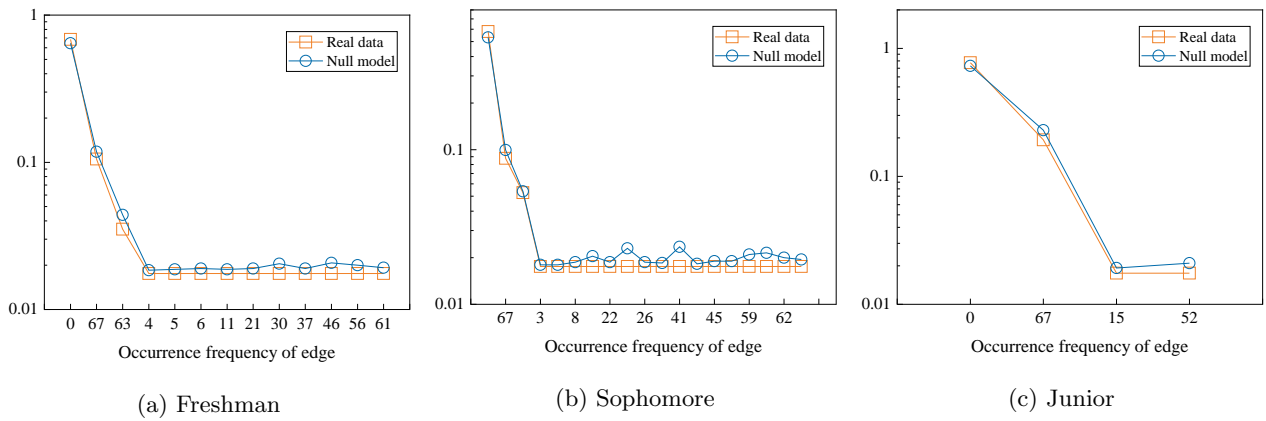


Figure 5: Null model for each group.

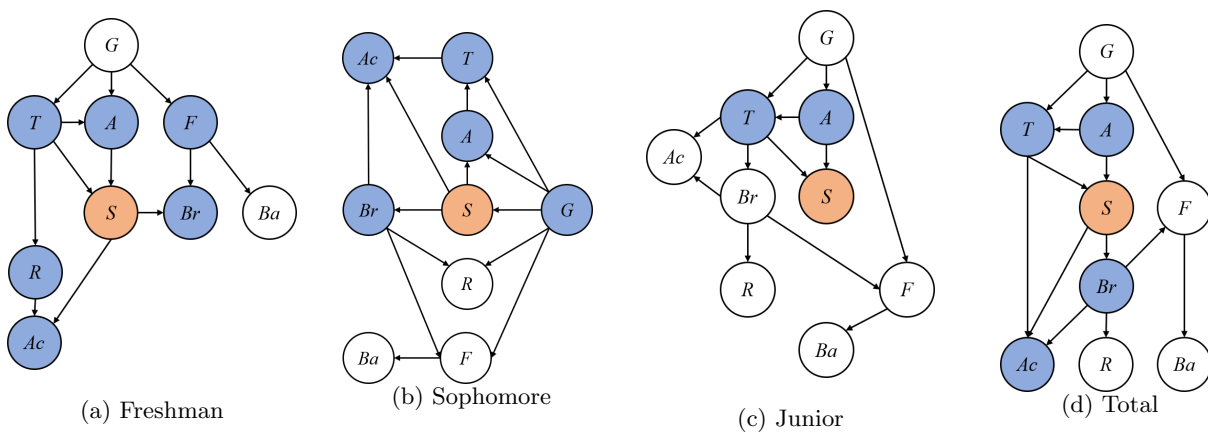


Figure 6: Consensus BN for each group. The orange node is sleep status and its Markov Blanket is represented by blue node.

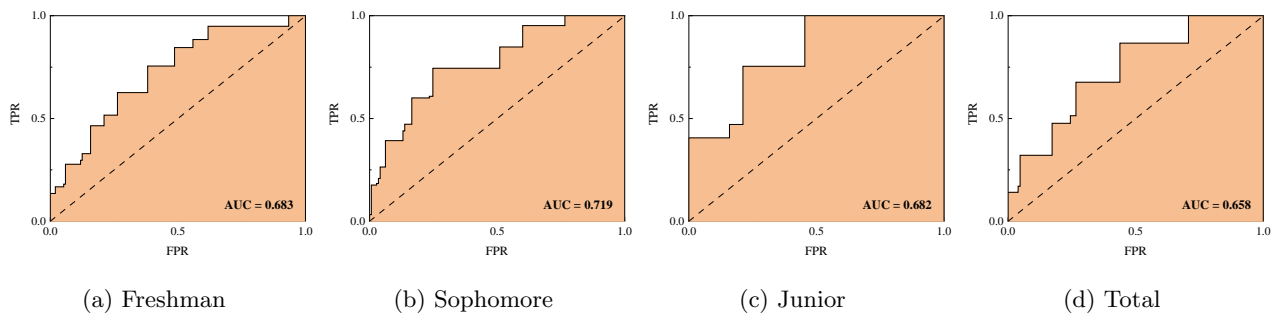


Figure 7: Prediction performance of each BNs. We use ROC to measure the prediction performance. AUC is indicated by yellow and diagonal dotted lines indicate 0.5 in random cases.

The three networks shown in Figure 6 represent three types of students: freshmen, sophomores, and juniors. It is obvious that the behaviors of different groups are different. According to the Markov property of the Bayesian network, the parents node  $\pi(X_i)$  are the most related to the output of the targeted node  $X_i$ . In this case, we summarize those BNs in Figure 6 and show the results in Table 3.

First, both  $A$  and  $T$  are frequent in  $\pi(S)$ , representing that a strong connection exists between sleep habit and surfing habit. Specifically, from the perspective of App preference, the probabilities of staying up late for people who love watching videos are 0.75, 0.63, and 0.70, in freshmen, sophomores, and juniors, respectively. For people who love playing games, the probabilities are 0.5, 0.39, and 0.47. Video lovers prefer to stay up late than game fans. Then, we summarize the most popular games among college students, and the top three are Glory of Kings, PUBG (Playerunknown’s Battlegrounds), and LOL (League of Legends), which account for more than 80 percent (This result is consistent with the findings of previous studies that people who primarily engage in group play have superior adherence to people who primarily play alone [26]). A common feature of these games is that they are all web-based battle games. Compared to previous single-player games, these games bring all players together to collaborate or compete with each other through an online platform. Players play with real people instead of computers, which is closer to social behavior. In this case, if the user’s enemies or teammates are sleeping, then this game naturally loses fun, which explains why game fans don’t like staying up late. This phenomenon can also be analyzed from another perspective. According to the theory mentioned in [18], ‘spectator experience’ and a sense of social presence are very important for players of online games, which can also explain the finding of our experiment. In contrast, watching a video is a behavior that a student can do well on their own, even in the middle of the night. Note that the above results show that people who love to watch videos are more likely to stay up late rather than that videos are to keep people stay up late compared with games, due to the differences in the number of each type of students.

Does the seemingly rational behavior of game lovers contradict the popular view of game addiction? We analyze the relation between app preference  $A$  and total surfing length  $T$  due to that  $A$  appears frequently in the parent node of  $T$ . As mentioned before, according to the total length of surfing time, we divide students into long-time internet users and short-time internet users. The probabilities of game fans

surfing the Internet for a long time are 0.57, 0.60, and 0.57 in three groups, which is higher than 0.25, 0.22, and 0.22 for video lovers. This is in line with the current popular view of game addiction. At the same time, the above mentioned can also explain why the short-time internet user prefers to stay up late. (Note that this finding does not conflict with the conclusions in [33] that a longer duration of digital media use was associated with reduced total sleep time and later bedtime because they focus on the digital media two hours before bedtime instead of all day).

Second, for the child nodes of  $S$ , the variables of breakfast status  $Br$  and academic performance  $Ac$  appear frequently according to Table 3. In the matter of having breakfast, the probability of students staying up late is 0.2 lower than normal students, which is consistent with common sense that students who like to stay up late usually get up late to miss breakfast. For the academic performance, students who go to bed early are more likely to achieve a good grade, which is consistent with previous findings, both in biology [27] and cognitive science [4].

#### 4.4.3 Inference

The Markov property of the Bayesian network implies that conditioned on the Markov blanket of a node (shown in Figure 6), the probability distribution of the node is independent of the rest of the network shown as the following equation:

$$X \perp \{U - MB(X) - \{X\}\} | MB(X) \quad (8)$$

where  $U$  is the set of all random variables in BNs and  $MB(X)$  represents the Markov blanket of variable  $X$ , which is a set of nodes that consists of the parents of the node, the children of the node, and any other parents of the children of that node.

In our research, we design an experiment of predictive inference to estimate the predictability of sleep status. In other words, we want to test whether the sleep state  $S$  can be predicted rather than pursue prediction performance. So, we use MLE to fit our datasets and use Area Under The Curve (AUC) to evaluate the prediction performance. For AUC, if the AUC is bigger than 0.5, representing that the prediction result is better than the random model. In other words, the  $S$  is predictive. From Figure 7, the AUC value for each group is 0.683, 0.719 and 0.682, which reflects the predictability of sleep pattern.

## 5. DISCUSSION AND CONCLUSION

In this paper, we design a Possion-based probabilistic mixture model to identify students who are used to stay up late based on the Internet access patterns. We apply this model to a real-world dataset and classify students into two groups: students who are used to stay up late and students who sleep on time. We profile students from five aspects in eight dimensions, including interest (reading status, app preference status, and surfing length status), orderliness (breakfast orderliness status and bath orderliness status), finance (financial status), academic performance (academic performance status) and gender. Then we build Bayesian networks to explore the relationship between these characteristics and sleeping habits and find that surfing habits have a big impact on sleep habits. Finally, we test the predictability of sleeping habits based on campus behavior features.

The assumption of our experiment that students will access the Internet through mobile phones or computers before bedtime is reasonable because the Internet pervades every aspect of our lives, including entertainment, study, social contact, and so on. However, the underlying assumptions made in this study raise a couple of limitations. First, we identify students staying up late based on a hypothesis that students use mobile phones or computers to access the Internet before going to bed. So, it is hard to detect the people who don't have this habit. Second, more data of life details need to be collected for drawing a valid and solid conclusion. There are multiple avenues for future work. First, we intend to expand our dataset and investigate this issue from more aspects. Second, we only check if the sleep status can be predicted rather than pursuing a precise prediction. Next, we plan to design a prediction model with a good performance and integrate this model into the modern educational management system and apply real-time data to detect the sleep status of students.

## 6. REFERENCES

- [1] R. Agrahari, A. Foroushani, T. R. Docking, L. Chang, G. Duns, M. Hudoba, A. Karsan, and H. Zare. Applications of bayesian network models in predicting types of hematological malignancies. *Scientific reports*, 8(1):6951, 2018.
- [2] W. Al-Salman, Y. Li, and P. Wen. Detecting sleep spindles in eegs using wavelet fourier analysis and statistical features. *Biomedical Signal Processing and Control*, 48:80–92, 2019.
- [3] Z. Alimoradi, C.-Y. Lin, A. Broström, P. H. Bülöw, Z. Bajalan, M. D. Griffiths, M. M. Ohayon, and A. H. Pakpour. Internet addiction and sleep problems: A systematic review and meta-analysis. *Sleep Medicine Reviews*, 47:51–61, 2019.
- [4] T. Althoff, E. Horvitz, R. W. White, and J. Zeitzer. Harnessing the web for population-scale physiological sensing: A case study of sleep and performance. In *WWW '17 Proceedings of the 26th International Conference on World Wide Web*, pages 113–122, 2017.
- [5] K. Bartel, R. Scheeren, and M. Gradisar. Altering adolescents' pre-bedtime phone use to achieve better sleep health. *Health Communication*, 34(4):456–462, 2019.
- [6] M. B. Becerra, B. S. Bol, R. Granados, and C. Hassija. Sleepless in school: The role of social determinants of sleep health among college students. *Journal of American College Health*, (1):1–7, 2018.
- [7] H. D. Bedru, S. Yu, X. Xiao, D. Zhang, L. Wan, H. Guo, and F. Xia. Big networks: A survey. *Computer Science Review*, 37:100247, 2020.
- [8] B. Bjorvatn, J. Mrdalj, I. W. Saxvig, T. Aasnæs, S. Pallesen, and S. Waage. Age and sex differences in bedroom habits and bedroom preferences. *Sleep medicine*, 32:157–161, 2017.
- [9] F. C. Brown, W. C. Buboltz Jr, and B. Soper. Relationship of sleep hygiene awareness, sleep hygiene practices, and sleep quality in university students. *Behavioral medicine*, 28(1):33–38, 2002.
- [10] W. C. Buboltz Jr, F. Brown, and B. Soper. Sleep habits and patterns of college students: a preliminary study. *Journal of American college health*, 50(3):131–135, 2001.
- [11] C. Burr, J. Morley, M. Taddeo, and L. Floridi. Digital psychiatry: Risks and opportunities for public health and wellbeing. *IEEE Transactions on Technology and Society*, 1(1):21–33, 2020.
- [12] C.-F. Chien, S.-L. Chen, and Y.-S. Lin. Using bayesian network for fault location on distribution feeder. *IEEE Transactions on Power Delivery*, 17(3):785–793, 2002.
- [13] D. Combs, J. L. Goodwin, S. F. Quan, W. J. Morgan, S. Shetty, and S. Parthasarathy. Insomnia, health-related quality of life and health outcomes in children: a seven year longitudinal cohort. *Scientific reports*, 6:27921, 2016.
- [14] J. N. Cousins, K. Sasmita, and M. W. L. Chee. Memory encoding is impaired after multiple nights of partial sleep restriction. *Journal of Sleep Research*, 27(1):138–145, 2018.
- [15] J. N. Cousins, E. van Rijn, J. L. Ong, K. F. Wong, and M. W. L. Chee. Does splitting sleep improve long-term memory in chronically sleep deprived adolescents? *nj Science of Learning*, 4(1):8, 2019.
- [16] F. de Arriba Pérez, J. M. S. Gago, and M. C. Rodríguez. Calculation of sleep indicators in students using smartphones and wearables. In *New Advances in Information Systems and Technologies*, pages 169–178. Springer, 2016.
- [17] A. P. Dempster, N. M. Laird, and D. B. Rubin. Maximum likelihood from incomplete data via the em algorithm. *Journal of the Royal Statistical Society: Series B (Methodological)*, 39(1):1–22, 1977.
- [18] N. Ducheneaut, N. Yee, E. Nickell, and R. J. Moore. Alone together? exploring the social dynamics of massively multiplayer online games. In *Proceedings of the SIGCHI conference on Human Factors in computing systems*, pages 407–416, 2006.

- [19] I. N. Fossum, L. T. Nordnes, S. S. Storemark, B. Bjorvatn, and S. Pallesen. The association between use of electronic media in bed before going to sleep and insomnia symptoms, daytime sleepiness, morningness, and chronotype. *Behavioral sleep medicine*, 12(5):343–357, 2014.
- [20] N. Friedman, D. Geiger, and M. Goldszmidt. Bayesian network classifiers. *Machine learning*, 29(2-3):131–163, 1997.
- [21] J. Grønli, I. K. Byrkjedal, B. Bjorvatn, Ø. Nødtvedt, B. Hamre, and S. Pallesen. Reading from an ipad or from a book in bed: the impact on human sleep. a randomized controlled crossover trial. *Sleep medicine*, 21:86–92, 2016.
- [22] D. Grossman and P. Domingos. Learning bayesian network classifiers by maximizing conditional likelihood. *Twenty-first international conference on Machine learning - ICML '04*, page 46, 01 2004.
- [23] T. Guo, F. Xia, S. Zhen, X. Bai, D. Zhang, Z. Liu, and J. Tang. Graduate employment prediction with bias. In *Thirty-Second AAAI Conference on Artificial Intelligence*. AAAI Press, 2020.
- [24] D. Heckerman, D. Geiger, and D. M. Chickering. Learning bayesian networks: The combination of knowledge and statistical data. *Machine learning*, 20(3):197–243, 1995.
- [25] F. V. Jensen et al. *An introduction to Bayesian networks*, volume 210. UCL press London, 1996.
- [26] M. D. Kaos, R. E. Rhodes, P. Hämäläinen, and T. N. Graham. Social play in an exergame: how the need to belong predicts adherence. In *Proceedings of the 2019 CHI Conference on Human Factors in Computing Systems*, pages 1–13, 2019.
- [27] P. Kelley, S. W. Lockley, R. G. Foster, and J. Kelley. Synchronizing education to adolescent biology: 'let teens sleep, start school later'. *Learning, Media and Technology*, 40(2):210–226, 2015.
- [28] Y.-M. Liu, L. Wang, H.-C. Chu, and C. Yang. Development of a mobile sleep-management system for improving students' lifestyles based on a self-regulated learning strategy. In *2017 6th IIAI International Congress on Advanced Applied Informatics (IIAI-AAI)*, pages 623–628. IEEE, 2017.
- [29] J. C. Lo, J. L. Ong, R. L. Leong, J. J. Gooley, and M. W. Chee. Cognitive performance, sleepiness, and mood in partially sleep deprived adolescents: The need for sleep study. *Sleep*, 39(3):687–698, 2016.
- [30] M. J. McGeachie, H.-H. Chang, and S. T. Weiss. Cgbayesnets: Conditional gaussian bayesian network learning and inference with mixed discrete and continuous data. *PLOS Computational Biology*, 10(6):e1003676, 2014.
- [31] G. McLachlan and T. Krishnan. *The EM algorithm and extensions*, volume 382. John Wiley & Sons, 2007.
- [32] K. Okano, J. R. Kaczmarzyk, N. Dave, J. D. Gabrieli, and J. C. Grossman. Sleep quality, duration, and consistency are associated with better academic performance in college students. *NPJ science of learning*, 4(1):1–5, 2019.
- [33] K. M. Orzech, M. A. Grandner, B. M. Roane, and M. A. Carskadon. Digital media use in the 2 h before bedtime is associated with sleep variables in university students. *Computers in human behavior*, 55:43–50, 2016.
- [34] J. Park, R. Yu, F. Rodriguez, R. Baker, P. Smyth, and M. Warschauer. Understanding student procrastination via mixture models. *International Educational Data Mining Society*, 2018.
- [35] M. E. Patrick, J. Griffin, E. D. Huntley, and J. L. Maggs. Energy drinks and binge drinking predict college students sleep quantity, quality, and tiredness. *Behavioral sleep medicine*, 16(1):92–105, 2018.
- [36] D. Peters, K. Vold, D. Robinson, and R. A. Calvo. Responsible ai-two frameworks for ethical design practice. *IEEE Transactions on Technology and Society*, 1(1):34–47, 2020.
- [37] A. C. Schneider, D. Zhang, and Q. Xiao. Adolescent sleep characteristics and body-mass index in the family life, activity, sun, health, and eating (flashe) study. *Scientific Reports*, 10(1):1–10, 2020.
- [38] E. B. Simon and M. P. Walker. Sleep loss causes social withdrawal and loneliness. *Nature communications*, 9(1):1–9, 2018.
- [39] P. Spirtes, C. Glymour, and R. Scheines. *Causation, prediction, and search*. 1993.
- [40] R. Tavernier and T. Willoughby. Sleep problems: predictor or outcome of media use among emerging adults at university? *Journal of Sleep Research*, 23(4):389–396, 2014.
- [41] R. Tavernier and T. Willoughby. A longitudinal examination of the bidirectional association between sleep problems and social ties at university: The mediating role of emotion regulation. *Journal of youth and adolescence*, 44(2):317–330, 2015.
- [42] L.-L. Tsai and S.-P. Li. Sleep patterns in college students: Gender and grade differences. *Journal of psychosomatic research*, 56(2):231–237, 2004.
- [43] J. Vroon, C. Zaga, D. Davison, J. Kolkmeier, and J. Linssen. Snoozle—a robotic pillow that helps you go to sleep: Hri 2017 student design competition. In *Proceedings of the Companion of the 2017 ACM/IEEE International Conference on Human-Robot Interaction*, pages 399–400. ACM, 2017.
- [44] H. Zhao, W. Gui, H. Huang, Y. Liu, H. Ding, W. Fan, S. Huang, W. Yang, X. Wang, and G. Chen. Association of long-term sleep habits and hypertension: a cross-sectional study in chinese adults. *Journal of Human Hypertension*, pages 1–10, 2019.
- [45] X. Zhao, J. Li, W. Liu, J. Zhang, and Y. Li. Design of the sleeping aid system based on face recognition. *Ad Hoc Networks*, 99:102070, 2020.

# Meta-Learning with Graph Neural Networks: Methods and Applications

Debmalya Mandal<sup>1</sup>, Sourav Medya<sup>2</sup>, Brian Uzzi<sup>2</sup>, Charu Aggarwal<sup>3</sup>

<sup>1</sup>Data Science Institute, Columbia University, New York

<sup>2</sup>Kellogg School of Management, Northwestern University

<sup>3</sup>IBM T. J. Watson Research Center, Yorktown Heights, New York

dm3557@columbia.edu, {sourav.medya,uzzi}@kellogg.northwestern.edu,  
charu@us.ibm.com

## ABSTRACT

Graph Neural Networks (GNNs), a generalization of deep neural networks on graph data have been widely used in various domains, ranging from drug discovery to recommender systems. However, GNNs on such applications are limited when there are few available samples. Meta-learning has been an important framework to address the lack of samples in machine learning, and in recent years, researchers have started to apply meta-learning to GNNs. In this work, we provide a comprehensive survey of different meta-learning approaches involving GNNs on various graph problems showing the power of using these two approaches together. We categorize the literature based on proposed architectures, shared representations, and applications. Finally, we discuss several exciting future research directions and open problems.

## 1. INTRODUCTION

The methods of artificial intelligence (AI) and machine learning have found tremendous success in various applications, ranging from natural language processing [17] to cancer screening [66]. Such success of AI systems can be attributed to various architectural innovations, and the ability of deep neural networks (DNN) to extract meaningful representations from Euclidean data (e.g. image, video etc.). However, in many applications, the data is graph-structured. For example, in drug discovery, the goal is to predict whether a given molecule is a potential candidate for a new drug, where the input molecules are represented by graphs. In a recommender system, the interaction between the users and the items are represented by a graph, and such non-Euclidean data is crucial in designing a better system.

The proliferation of graph structured data in various applications has led to Graph Neural Networks (GNNs) which are generalizations of DNN for graph-structured inputs. The main goal of GNNs is to learn effective representations of the graphs. Such representations map the vertices, edges, and/or graphs to a low-dimensional space, so that the structural relationships in the graph are reflected by the geometric relationships in the representations [29]. In recent years, GNNs have been applied in diverse domains, often with surprising positive results like discovery of a new antibiotic [57], accurate traffic forecasting [14], etc.

Despite of recent success of GNNs in various domains, GNN frameworks have their own shortcomings. One of the major challenges in applying GNNs, particularly for large graph-structured datasets,

is the limited number of samples. Furthermore, real-world systems like recommender systems often need to handle diverse types of problems, and must adapt to a new problem with very few observations. In recent years, meta-learning has turned out to be an important framework to address these shortcomings of deep learning systems. The main idea behind meta-learning is to design learning algorithms that can leverage prior learning experience to adapt to a new problem quickly, and learn a useful algorithm with few samples. Such approaches have been quite successful in diverse applications like natural language processing [41], robotics [48], and healthcare [74].

Recently, several meta learning methods to train GNNs have been proposed for various applications. There are two main challenges in applying meta-learning to graph-structured data. First, an important challenge is to determine the type of representation that is shared across different tasks. As GNNs are used for a wide range of tasks from node classification to graph classification, the learned shared representation needs to consider the type of tasks to be solved and this makes the choice and design of architecture quite important for meta-learning. Second, in a multi-task setting, we usually have few samples from each task. Thus, the support and query examples have often limited overlap in terms of similarity. For example, in node classification tasks, the nodes rarely are similar in the support and query set of a given task. On the other hand, in link prediction, the support and query edges are often located far away from each other in the graph. Therefore, a major challenge in applying meta-learning to GNNs is to model the dependencies among nodes (or edges) that are far apart (both distance-wise and similarity-wise) from each other in the graph. In this survey, we review the growing literature on meta learning with GNNs. There are several thorough individual surveys on GNNs [77, 67] and meta-learning [30], but we believe this survey is the first effort to categorize and comprehensively review the existing papers on meta learning with GNNs.

### 1.1 Our Contributions

Besides providing background on meta-learning and architectures based on GNNs individually, our major contributions can be summarized as follows.

- **Comprehensive review:** We provide a comprehensive review of meta learning techniques with GNNs on several graph problems. We categorize the literature based on methods, representations and applications and show various scenarios where limitations of GNNs are addressed via meta learning.
- **Future directions:** We discuss how meta learning and GNNs

can address some of the challenges in several areas: (i) combinatorial graph problems, (ii) graph mining problems, and (iii) other emerging applications such as traffic flow prediction, molecular property prediction, and network alignment.

The rest of this paper is organized as follows. Section 2 provides background on a few key graph neural network architectures. Section 3 outlines the background on meta-learning and major theoretical advances. A comprehensive categorization of the papers that use the framework of meta-learning equipped with GNNs on important graph related problems is described in Sections 4 and 5. First, Section 4 covers applications of meta-learning framework for solving some classical graph problems. The problem discussed here doesn't explicitly propose a multi-task setting, rather the meta-learning framework is applied to a fixed graph. In Section 5 we cover the literature on graph meta learning when there are multiple tasks and the graph might change with the tasks. Although various GNNs have been proposed for graph meta-learning, they can be categorized broadly based on the type of shared representation, which can be either at a local level (node/edge based) or at the global level (graph based). Table 1 provides an overview of various papers categorized by the type of shared representation and the application domains. Table 2 presents the papers described in Section 5 based on the corresponding meta-learning approaches. Section 6 covers a broad range of applications of meta-learning on GNNs and Section 7 suggests some exciting future directions.

## 2. GRAPH NEURAL NETWORKS

Generalizing deep learning on graphs has resulted in an exciting area of Graph Neural networks (GNNs). GNNs embed or represent nodes as points in a vector space with the help of structural and attribute information from the neighbourhood of a node and the node itself. They encode this information via non-linear transformations and aggregation functions into a final representation. The proposed architectures can be broadly categorized into two types: (i) *convolution on neighborhood*, and (ii) *location-aware*.

(i) **Convolution on neighborhood:** The primary examples of architectures that are based on *convolution on neighborhood* include GCN [36], GRAPH SAGE [28], and GAT [61]. These architectures mostly create representations of nodes through a *convolution* operation  $\psi$  over its neighborhood, i.e., the embedding,  $z_{v,G} = \psi(N_G^k(v))$  where the ( $k$ -hop) neighborhood (set of nodes) of the node  $v$  in the graph  $G$  is  $N_G^k(v)$ . Thus, two nodes with similar neighborhoods are likely to have similar embeddings.

(ii) **Location-aware:** The examples of GNNs that are location aware framework include PGNN [71] and GRAPH REACH [49]. In this approach, if two nodes are located close (usually by number of hops) to each other in the graph then they are expected to have similar embeddings. If the graph has a high clustering coefficient, then one-hop neighbors of a node share many other neighbors among them as well. Therefore, if two nodes are close to each other, they have a high likelihood of having similar neighborhoods. Many real graphs have *small-world* and *scale-free* properties and have high clustering coefficients. Next, we briefly describe the key architectures of GNNs.

**GCN [36]:** A primary contribution in applying neural architectures on graphs has been made by [36] with the introduction of Graph Convolutional Networks (GCNs). GCNs are analogous version of convolutional neural networks (CNNs) on graphs. Inspired by the idea of representing a pixel with information from its nearby pixels (filter in CNNs), graph convolutions also apply the key idea of aggregating feature information from a node's local neighborhood. More formally, GCNs are neural network architectures that

produces a  $d$ -dimensional embeddings for each node by taking as input adjacency matrix  $A$  and node features  $X$ ;  $\text{GCN}(A, X) : \mathbb{R}^{n \times n} \times \mathbb{R}^{n \times p} \rightarrow \mathbb{R}^{n \times d}$ . The idea is to aggregate feature information from a node's neighborhood (can be generalized to multiple hops) and its own features to produce the final embedding. A 2-layer (neighbourhood is 2-hops) GCN can be defined as follows:

$$\text{GCN}(A, X) = \sigma(\hat{A}\sigma(\hat{A}XW^{(1)}W^{(2)}))$$

where  $\hat{A} = \tilde{D}^{-\frac{1}{2}}\tilde{A}\tilde{D}^{-\frac{1}{2}}$  is the normalized adjacency matrix with  $\tilde{D}$  as weighted degree matrix and  $\tilde{A} = I_n + A$  with  $I_n$  being an  $n \times n$  identity matrix and  $\sigma$  is an activation function. Moreover,  $W^{(i)}$  is a weight matrix for the  $i$ -th layer to be learned during training, with  $W^{(1)} \in \mathbb{R}^{p \times d'}$ ,  $W^{(2)} \in \mathbb{R}^{d' \times d}$ , and  $d$  ( $d'$ ) being the number of neural network nodes in the output (hidden) layer.

**GRAPH SAGE [28]:** Hamilton et al. [28] propose an inductive framework with an aggregation function that is able to share weight parameters ( $\mathbb{W}^k$ ) across nodes, can be generalized to unseen nodes and scale to large datasets. To learn representation  $h_v^k$  of a node  $v$ , it iterates over all nodes which are in their  $K$ -hop neighborhood. While iterating over node  $v$ , it *aggregates* (with  $\text{AGGREGATE}_k$ ) the current representations of  $v$ 's neighbors ( $\mathbf{h}_N^k(v)$ ) and *concatenate* with the current representation of  $v$  ( $h_v^{k-1}$ ), which is then fed through a fully connected layer with an activation function. Intuitively, with more iterations, nodes incrementally receive information from neighbors of higher depth (i.e., distance). More specifically for  $k$ -th iteration,

$$\begin{aligned} \mathbf{h}_N^k(v) &= \text{AGGREGATE}_k \left( \left\{ h_u^{k-1}, \forall u \in N(v) \right\} \right) \\ \mathbf{h}_v^k &= \sigma \left( \mathbb{W}^k \cdot \text{CONCAT} \left( \mathbf{h}_N^k(v), h_v^{k-1} \right) \right) \end{aligned}$$

**GAT [61]:** Graph Attention Networks (GATs) [61] learn edge weights using attention mechanisms. GAT does not assume that the contributions of neighbouring nodes are all equal unlike in GRAPH SAGE [28]. GAT learns the relative importance/weights between two connected nodes. The graph convolutional operation ( $k$ -th iteration) is defined as follows:

$$\mathbf{h}_v^k = \sigma \left( \sum_{u \in N(v) \cup v} \alpha_{v,u}^k \mathbb{W}^k h_u^{k-1} \right)$$

where  $\alpha_{v,u}$  measures the strength between the node  $v$  and its neighbour  $u \in N(v)$ . GAT has been shown to outperform both GCN and GRAPH SAGE in node classification task both in transductive as well as inductive settings in benchmark datasets.

**PGNN [71]:** Unlike in GRAPH SAGE where the representation of a node depends on its  $k$ -hop neighborhood, PGNN follows a different paradigm and aims to incorporate positional information of a node with respect to the nodes in the entire network. The key idea is that the position of a node can be captured via a low-distortion embedding by quantifying the distance between that node and a set of anchor nodes. The framework first samples multiple sets of anchor nodes. It also learns a non-linear aggregation scheme to combine the features of the nodes in each anchor set. The aggregation is normalized by the distance between the node and the anchor-set.

**Other variations:** There are several other variations and improvements of GNNs that are based on different mechanisms: GAT is further extended by Gated Attention Network (GAAN) [72] through a self-attention mechanism which computes an additional attention score for each attention head. Graph Autoencoders [9, 37] encode nodes/graphs into a latent vector space and further reconstruct the graph related data depending on the application from this encoding in an unsupervised fashion; Recurrent GNNs [53, 39] apply

the same set of parameters recurrently over nodes to extract high-level node representations. For a comprehensive survey on GNNs, please refer to [67].

## 2.1 Applications

GNNs outperform traditional approaches for semi-supervised learning tasks (e.g. node classification) on graphs. The high level applications of GNNs can be categorized in three major tasks: node classification, link prediction, and graph classification. For node classification and link prediction tasks, traditionally four benchmark datasets are used: Cora, Citeseer, Pubmed, and protein-protein interaction (PPI) dataset. Shchur et al. [56] and Errica et al. [22] provide a detailed comparison of performances of the key architectures on node and graph classification tasks. GNNs are also used in the link prediction task that has applications in many domains such as friend or movie recommendation, knowledge graph completion, and metabolic network reconstruction [73].

## 3. BACKGROUND ON META-LEARNING

Meta-learning has turned out to be an important framework to address the problem of limited data in various machine learning applications. The main idea behind meta-learning is to design learning algorithms that can leverage prior learning experience to adapt to a new problem quickly, and learn a useful algorithm with few samples [55]. Such approaches have been quite successful in diverse applications like natural language processing [41], robotics [48], and healthcare [74].

### 3.1 Framework

In standard supervised learning, we are given a training dataset  $\mathcal{D} = \{\mathbf{x}_i, y_i\}_{i=1}^n$ , a loss function  $\ell$ , and we aim to find a predictive model of the form  $\hat{y} = f_\theta(\mathbf{x})$ .

$$\theta^* = \operatorname{argmin}_\theta \mathcal{L}(\mathcal{D}, \theta) = \operatorname{argmin}_\theta \sum_{i=1}^n \ell(f_\theta(\mathbf{x}_i), y_i)$$

In meta-learning, we are given samples from a number of different tasks and the goal is to learn an algorithm that generalizes across tasks. In particular, the tasks are drawn from a distribution  $p(\mathcal{T})$ , and the meta-objective is to find a common parameter that works across the distribution of tasks.

$$\omega^* = \operatorname{argmin}_\omega \sum_{\substack{\mathcal{T}_i \sim p(\mathcal{T}) \\ \mathcal{D}_i \sim \mathcal{T}_i}} \mathcal{L}_i(\mathcal{D}_i, \omega) \quad (1)$$

In the meta-test phase, we are given a target task (say task 0) and we use the meta-knowledge  $\omega^*$  to obtain the best parameter for the target with few samples.

$$\theta_0^* = \operatorname{argmin}_\theta \mathcal{L}_0(\mathcal{D}_0, \theta | \omega^*)$$

### 3.2 Training

Many popular meta-learning algorithms are based on gradient descent on the meta-parameter  $\omega$  [23, 52]. In order to understand how to perform gradient descent with respect to  $\omega$ , it is insightful to frame Equation (1) as a bi-level optimization problem.

$$\omega^* = \operatorname{argmin}_\omega \sum_{\substack{\mathcal{T}_i \sim p(\mathcal{T}) \\ \mathcal{D}_i \sim \mathcal{T}_i}} \mathcal{L}(\mathcal{D}_i, \theta_i^*(\omega), \omega)$$

$$\text{s.t. } \theta_i^*(\omega) = \operatorname{argmin}_\theta \mathcal{L}_i(\theta, \omega, \mathcal{D}_i) \quad \forall i$$

If we have a model for the inner-optimization method, then a gradient of the objective with respect to  $\omega$  can be computed by using

the chain rule e.g.

$$\nabla_\omega \mathcal{L}(\mathcal{D}_i, \theta_i^*(\omega), \omega) = \nabla_{\theta_i^*(\omega)} \mathcal{L}(\mathcal{D}_i, \theta_i^*(\omega), \omega) \frac{d\theta_i^*(\omega)}{d\omega}$$

However, often the inner objective function is non-convex, and hard to solve. So model agnostic meta learning (MAML), introduced by Finn et al. [23] suggests to first take a gradient step for each task  $i$  as follows:

$$\theta'_i = \theta_i(\omega) - \alpha \nabla_\theta \mathcal{L}_i(\theta_i(\omega), \omega, \mathcal{D}_i)$$

Then MAML replaces  $\theta_i^*(\omega)$  in the outer objective, i.e.,

$$\omega = \omega - \beta \nabla_\omega \sum_i \mathcal{L}(\mathcal{D}_i, \theta'_i, \omega)^1$$

We now instantiate the MAML algorithm for the task of classifying nodes of a graph. Recall the GCN framework from Section 2. Here the  $t$ -th task is classification of nodes of a graph  $G_t$  with adjacency matrix  $A_t$  and node-feature matrix  $X_t$ . Then a standard two-layer GCN for node classification problem is given as follows:

$$f(X_t, A_t, W_t) = \operatorname{softmax} \left( \hat{A}_t \operatorname{ReLU} \left( \hat{A}_t X_t W_t^{(1)} \right) W_t^{(2)} \right) \quad (2)$$

Given labels of the nodes  $Y_t$ , such a network is often trained with the cross-entropy loss:

$$\mathcal{L}_t(X_t, A_t, W_t) = - \sum_{\ell} \sum_f Y_{\ell f} \ln f(X_t, A_t, W_t)_{\ell f}$$

Usually, the parameters  $W_t$  are trained by stochastic gradient descent. Here, we wish to identify a meta parameter vector  $W_*$ , which is close to the parameters of different tasks (i.e.  $\|W_t - W_*\|_F \leq \delta$  for some  $\delta > 0$ ). The benefit of learning such meta-parameters  $W_*$  is that, on a new task  $s$ , we can initialize task-parameter  $W_s$  as  $W_*$  and the new task would require very few samples to train. Algorithm 1 describes the MAML algorithm instantiated for the case of node classification with GCN based representation.

---

#### ALGORITHM 1: Model Agnostic Meta-Learning for GCN

---

**Input:** Step sizes  $\alpha$  and  $\beta$ .

Initialize  $W_*$ .

**do**

    Sample a batch of  $T$  tasks  $\{G_i\} \sim p(\cdot)$ .

    Sample a batch of  $T$  datasets  $\{\mathcal{D}_i = (A_i, X_i, Y_i)\}$  where  $\mathcal{D}_i \sim G_i$ .

**for each task  $t$  in  $T$  do**

        Update  $W_t = W_* - \alpha \nabla_W \mathcal{L}_t(X_t, A_t, W) \Big|_{W=W_*}$ .

    Update  $W_* = W_* - \beta \nabla_W \sum_t \mathcal{L}(X_t, A_t, W) \Big|_{W=W_t}$ .

**while Not Convergence**

**return** Meta-parameter  $W_*$ .

---

### 3.3 Representation Learning

Another perspective of meta-learning, which will be particularly important for the context of graph neural networks, is learning a shared representation across different tasks. Here we assume that, given an input  $x$ , the training data from the  $t$ -th task is generated as  $y_t = f_t \circ h(x) + \eta_t$ , where  $\eta_t$  is some iid noise. Effectively, the function  $h$  maps input  $x$  to a shared representation and then a task-specific function  $f_t$  is applied to generate the task-specific representation.

<sup>1</sup>We write  $\theta_i(\omega)$  to denote the meta-parameter  $\omega$  adapted to task  $i$ .

During the meta-training phase, we attempt to learn the shared function  $h$ . Suppose we are given  $T$  datasets  $\mathcal{D}_t = \{(x_{ti}, y_{ti})_{i=1}^{n_t}\}$  for  $t = 1, \dots, T$ . Then we solve the following optimization problem to recover  $h$ .

$$\operatorname{argmin}_{h, \{f_t\}_{t=1}^T} \sum_{t=1}^T \sum_{i=1}^{n_t} \mathcal{L}_t(y_{ti}, f_t(h(x_{ti}))) + \mathcal{R}(h) + \sum_t \mathcal{R}(f_t) \quad (3)$$

Here  $\mathcal{R}(\cdot)$  is some regularization function, and let  $\hat{h}, \{\hat{f}_t\}_{t=1}^T$  be its solution. In general, the optimization problem defined in Equation (3) is hard to solve unless we make specific assumption about the types of functions. For example, even if we assume  $f_t$  is same across the tasks and in fact an identity function, the problem defined in Equation (3) can involve learning a general neural network based shared representation  $h$ . For the special case of linear models, this problem can be solved efficiently (e.g. by using matrix regression [59]). In this survey, we focus on gradient based methods for learning the shared representation  $h$  in equation (3), which has been quite successful in practice. In the meta-test phase, we are given samples from a new task  $s$  e.g.  $\{(x_{si}, y_{si})_{i=1}^{n_s}\}$ . We substitute  $\hat{h}$ , the estimate of the common representation function  $h$ , and learn the new task-specific function  $f_s$ .

$$\hat{f}_s \leftarrow \operatorname{argmin}_{f_s} \sum_{i=1}^{n_s} \mathcal{L}_s(y_{si}, f_s(\hat{h}(x_{si}))) + \mathcal{R}(f_s)$$

We now instantiate this framework for the task of classifying nodes of a graph. As before, we use two-layer GCN where the model is defined in Equation (2). However, we now assume that the first layer is shared across different tasks and only the second layer is trained for a new task. In particular, we assume  $W_t = [W^*; W_t^{(2)}]$ . Although, the optimization problem in Equation (3) is NP-hard to solve with this particular type of representation, we can write down an algorithm to solve for the meta-parameter  $W^*$  using gradient descent. Algorithm 2 describes this algorithm and returns the meta-parameter  $W^*$ .

---

**ALGORITHM 2:** Shared Representation Learning for GCN

---

**Input:** Step sizes  $\alpha$  and  $\beta$ , datasets  $\mathcal{D}_t = \{(x_{ti}, y_{ti})_{i=1}^{n_t}\}$  for  $t = 1, \dots, T$ .

Initialize  $W^*$   
Initialize  $W_t^{(2)}$  for  $t = 1, \dots, T$ .

Set  $W_t = [W^*, W_t^{(2)}]$ .

**do**  
  **for each task  $t$  in  $[T]$  do**  
    Update  
     $W_t^{(2)} = W_t^{(2)} - \alpha \nabla_W \mathcal{L}_t(\mathcal{D}_t, [W^*; W]) \Big|_{W=W_t^{(2)}}$   
    Update  $W^* = W^* - \beta \nabla_W \sum_t \mathcal{L}(\mathcal{D}_t, [W; W_t^{(2)}]) \Big|_{W=W^*}$

**while Not Convergence**  
**return Meta-parameter  $W^*$ .**

---

### 3.4 Theory

Despite immense success, we are yet to fully understand the theoretical foundations of meta-learning algorithms. Baxter [5] first prove generalization bound for multitask learning problem, by considering a model where tasks with shared representation are sampled from a generative model. Pontil et al. [51], and Maurer et al. [46] develop general uniform-convergence based framework to analyze multitask representation learning. However, they assume oracle access to a global empirical risk minimizer. Recently, there

have been promising attempts to understand meta learning from representation learning. The main idea is that the tasks share a common shared representation and a task-specific representation [60, 59, 21]. If the shared representation is learned from the training tasks, then the task-specific representation for the new task can be learned with only a few samples. Finally, there have been interesting recent work trying to understand gradient-based meta-learning. [24, 4, 35, 16] analyze gradient based meta-learning in the framework of online convex optimization (OCO). They assume that the parameters of the tasks are close to a shared parameter to bound regret in the OCO framework.

## 4. META-LEARNING ON FIXED GRAPHS

In this section, we review applications of meta-learning for solving some classical problems on graphs. Here we consider the setting when the underlying graph is fixed and the node/edge features do not change with different tasks. In fact, we are not in a multitask framework where there are a number of tasks and few samples are available from each task. Rather, the framework of meta-learning is applied to various graph problems by creating multiple tasks either considering the nodes or the edges.

### 4.1 Node Embedding

The goal of node embedding is to learn representations for the nodes in the graph so that any downstream application can directly work with these representations, without considering the original graph. This problem is often challenging in practice because the degree distributions of most graphs follow a power law distribution and there are many nodes with very few connections. Liu et al. [43] address this issue by applying meta-learning to the problem of node embedding of graphs. They set up a regression problem with a common prior to learn the node embeddings. Since the base representations of high-degree nodes are accurate, they are used as meta training set to learn the common prior. The low degree nodes have only a few neighbors (samples), the regression problem for learning their representations is formulated as a meta-testing problem, and the common prior is adapted with a small number of samples for learning the embeddings of such nodes.

### 4.2 Node Classification

The node classification task aims to infer the missing labels of nodes of a given partially labeled graph. This problem often appears in diverse contexts such as document categorization and protein classification [58, 6], and has received significant attention in recent years. However, often many classes are novel i.e., they have a small number of labeled nodes. This makes meta-learning or few-shot learning particularly suitable for this problem.

Zhou et al. [76] have applied a meta-learning framework for the node classification problem on graphs by learning a transferable representation using data from classes that have many labeled examples. Then, during the meta-test phase, this shared representation is used to make predictions for novel classes with few labeled samples. Ding et al. [19] improve upon the previous method by considering a prototype representation of each class and meta-learning the prototype representation as an average of weighted representations of each class. Lan et al. [38] address the same problem via meta-learning but in a different setting where the nodes do not have attributes. Their method only uses the graph structure to obtain latent representation of nodes for the task. Subsequently, Liu et al. [42] point out that it is important to also learn the dependencies among the nodes in a task, and propose to use nodes with high centrality scores (or hub nodes) to update the representations learned by a GNN. This is done by selecting a small set of hub

Representation	Graph applications		
	Node classification	Link Prediction	Graph Classification
Node/Edge Level	Meta-GNN [76], GPN [19], RALE [42], AMM-GNN [62] SAME [8], SELAR [32] GFL [70], Meta-GDN [20]	MetaR [12], GEN [1] SAME [8], SELAR [32]	SAME [8]
Graph Level	MI-GNN [64]	Meta-graph [7]	AS-MAML [44], Spectral [11]

Table 1: Organization of the papers on Meta-learning and GNNs based on applications and underlying graph-related representations. The abbreviations of the frameworks (methods) are as follows. GPN: Graph Prototypical Networks, MetaR: Meta Relational learning, GEN: Graph Extrapolation Networks, RALE: Relative and Absolute Location Embedding, AMM-GNN: Attribute Matching Meta-learning Graph Neural Networks, SAME: Single-task Adaptation for Multi-task Embeddings, SELAR: SELF-supervised Auxiliary Learning, GFL: Graph Few-shot Learning, GDN: Graph Deviation Networks, MI-GNN: Meta-Inductive framework for Graph Neural Network, AS-MAML: Adaptive Step Model Agnostic Meta Learning.

nodes and for each node  $v$ , considering all the paths to the node  $v$  from the set of hub nodes. It helps to encode the absolute location in the graph. Parallel to these developments, Yao et al. [70] consider a metric-learning based approach where the label of a node is predicted to be the nearest class-prototype in a transferable metric space. They first learn a class-specific representation using a GNN, and then learn a task-specific representation using hierarchical graph representations.

Finally, the few-shot node classification task has also been used in the presence of noisy or inaccurate labels in the support sets of different tasks. Ding et al. [18] propose a method (Graph Hallucination Network) that creates a set by taking a specified number of samples from a class. Then the method learns to produce a confidence score on the accuracy of the label of each node in the set. By using these weights/scores, the final cleaner (i.e., less noisy) node representations are generated. The rest of the algorithm follows the standard MAML framework.

### 4.3 Link Prediction

The objective of the link prediction problem is to identify pairs of nodes that will either form a link or not. Meta-learning has been shown to be useful for learning new relationship via edges/links especially in multi-relational graphs.

In multi-relational graphs, an edge is represented by a triple of two end points and a relation. Such graphs appear in many important domains such as drug-drug interaction prediction. The goal of link prediction in multi-relation graphs is to predict new triples given one end point of a relation  $r$  with observing a few triples about  $r$ . This problem is challenging as only few associative triples are usually available. Chen et al. [12] use meta-learning to solve the link prediction problem in two steps. First, they design a Relation-Meta Learner which learns shared structure across a number of relations. Such a meta-learner generates relation meta from heads’ and tails’ embeddings in the support set. Second, they use an embedding learner that calculates the truth values of triples in support set via end points’ embeddings and relation meta.

Multi-relational graphs are even more difficult to manage with their dynamic nature (addition of new nodes) over time and the learning is even more difficult when these newly evolved nodes have only few links among them. Baek et al. [1] introduce a few-shot out-of-graph link prediction technique, where they predict the links between the seen and unseen nodes as well as between the unseen nodes. The main idea is to randomly split the entities in a given graph into the meta-training set for simulated unseen entities, and the meta-test set for real unseen entities.

Finally, Hwang et al. [32] show the effectiveness of graph neu-

ral networks on downstream tasks such as node classification and link prediction via a self-supervised auxiliary learning framework combined with meta-learning. The auxiliary task such as meta-path prediction does not need labels and thus the method becomes self-supervised. In the meta learning framework, various auxiliary tasks are used to improve generalization performance of the underlying primary task (e.g., link prediction). The proposed method effectively combines the auxiliary tasks and automatically balances them to improve performance on the primary task. The method is also flexible to work with any graph neural network architecture without additional data.

## 5. META-LEARNING ON GRAPH NEURAL NETWORKS

We now discuss the growing and exciting literature on graph meta learning where there are multiple tasks and the underlying graph can change across the tasks. The changes in graphs occur when either the node/edge features change, or the underlying network structure changes with the tasks. In the context of meta-learning, several architectures have been proposed in recent years. However, a common thread underlying all of them is a shared representation of the graph, either at a local node/edge level, or at a global graph level. Based on the type of shared representation, we categorize the existing works into two groups. Most of the existing literature adopt the MAML algorithm [23] to train the proposed GNNs. The outer loop of MAML updates the shared parameter, whereas the inner loop updates the task-specific parameter for the current task. Table 2 lists the shared and the task-specific parameters for all the papers in this section.

### 5.1 Node/Edge Level Shared Representation

First, we consider the setting where the shared representation is local i.e. node or edge based. Huang et al. [31] consider the node classification problem where the input graphs as well as the labels can be different across tasks. They learn a representation for each node  $u$  in two steps. First, the method extracts a subgraph  $S_u$  corresponding to the set of nodes  $\{v : d(u, v) \leq h\}$  where  $d(u, v)$  is the distance of the shortest path between nodes  $u$  and  $v$ . Then it feeds the subgraph  $S_u$  through a GCN to learn a representation for node  $u$ . The theoretical motivation behind considering the graph  $S_u$  is that the influence of a node  $v$  on  $u$  decreases exponentially as the shortest-path distance between them increases. Once the nodes are encoded, one can learn any function  $f_\theta$  that maps the encodings to class labels. Huang et al. [31] use MAML to learn this function with very few samples on a new task, enjoying the benefits of node-

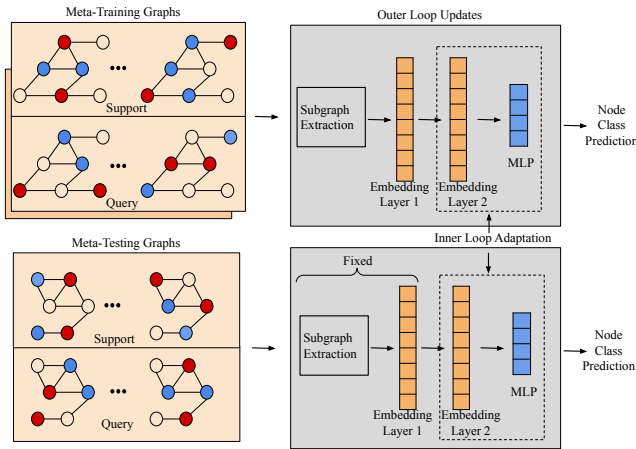


Figure 1: A prototype of the meta learning framework with GNNs for solving node classification problem. This is based on the architectures proposed by [31] and [62]. Following [31], the neighborhoods of each node are used for node embedding. Embedding layer 1 is trained in the outer loop of MAML, whereas the other layers are adapted for particular tasks.

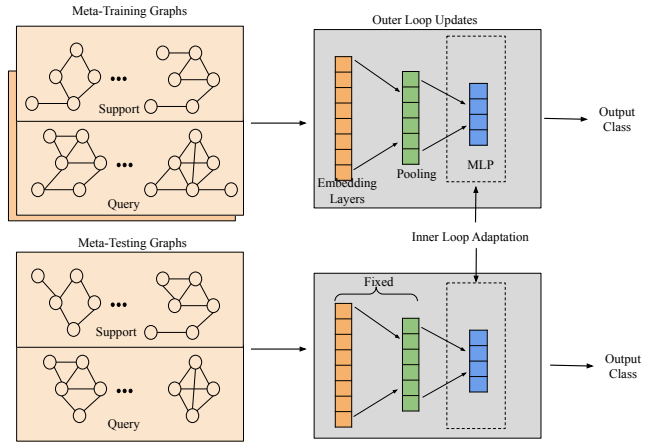


Figure 2: A prototype of the meta learning framework with GNNs for solving graph classification problem. This is based on the architectures proposed by [44], and [8]. The embedding and pooling layers learn global representation of the input graph, and are trained in the outer loop of MAML. The final multi-layer perceptron (MLP) is used for the classification task and is adapted to the particular task at meta-test.

Papers	Meta-learning parameters	
	Inner Loop (Task-Specific)	Outer Loop (Shared)
[31]	Node embeddings	Classification
[62]	Node embeddings	Feature matrix
[11]	Graph feature, Super-class	graph label/ actual class
[44]	Graph feature, Graph embedding	graph embedding/ Classification
[8]	Node Embedding	Output Layer
[7]	VGAE Initialization	Graph Signature (GCN + MLP)
[43]	High-degree node embedding	node specific embedding

Table 2: Organization of the papers in Section 6 based on the corresponding meta-learning approaches.

level shared representations in node classification.

Wang et al. [62] also consider the few-shot node classification problem for a setting where the network structure is fixed, but the features of the nodes change with tasks. In particular, given a base graph with node feature matrix  $X \in \mathbb{R}^{n \times d}$ , the proposed model learns a new feature matrix  $X_t = X \odot \alpha_t(\phi) + \beta_t(\phi)$  for the  $t$ -th task, and then use a GNN  $f_\theta(X_t)$  to learn the node representations for the  $t$ -th task. During training, the outer loop updates the  $\phi$  parameters, whereas the inner loop of MAML only updates the  $\theta$ -parameter. This enables quick adaptation to the new task.

Wen et al. [64] study the problem of node classification in an inductive setting, where the graph instances in testing and training do not overlap. Their method involves computing a task prior given a graph (i.e., its representation) using multi-layer perceptron (MLP). These representations are useful for the graph-level adaptation. They used the traditional MAML paradigm in their approach for the task-level adaptation.

## 5.2 Graph Level Shared Representation

In this subsection, we discuss the setting when the shared representation is global i.e. graph-level. A canonical application of this representation is the *graph classification* problem, where the goal is to classify a given graph to one of many possible classes. This problem appears in many applications, ranging from bioinformatics to social network analysis [69]. However, in many settings, the number of samples/graphs available for a particular task is few and the graph classification task often requires a large number of samples for high quality prediction. These challenges can be addressed via meta-learning. The existing papers on using meta-learning for graph classification usually learn an underlying shared representation and adapt the representation for a new task.

Chauhan et al. [11] propose the few-shot graph classification task based on graph spectral measures. In particular, they train a feature-extractor  $F_\theta(\cdot)$  to extract features from the graphs in meta-training. For classification, they first use a unit  $C^{\text{sup}}$  to first predict the super-class probability of a graph which is a clustering of abundant base class labels. Then they use  $C^{\text{att}}$ , an attention network to predict the actual class label. During the meta-test phase, the weights of the networks  $F_\theta(\cdot)$  and  $C^{\text{sup}}$  are fixed, and the network  $C^{\text{att}}$  is retrained on the new test classes. As the feature extractor  $F_\theta$  is the common shared structure, and is not retrained on the test tasks, this approach requires few samples from new classes.

Although Chauhan et al. [11] propose a novel meta-learning architecture for graph classification, there are several limitations. First, the architecture assumes significant overlap between the super-class structure of the test and the training set. Second, the fixed feature extractor cannot be updated for the new tasks. Ma et al. [44] design a better meta-learning technique by allowing the feature extractor to adapt efficiently for new tasks. They apply two networks – embedding layers ( $\theta_e$ ), followed by classification layers ( $\theta_c$ ) to classify a given graph. However, for a new task, both  $\theta_e$  and  $\theta_c$  are updated. In particular, the authors use MAML [23] to update the parameters and use a reinforcement learning based controller to determine how the inner loop is run i.e., what is the optimal adaptation step for a

new task. The parameters of the controller is updated using the graph's embedding quality and the meta-learner's training state. Jiang et al. [33] solve the problem of few-shot graph classification via a paradigm in meta learning called metric learning approach [63] that is different from MAML. In the training phase, the idea is to get a mean representations of the instances in each class in the support set. The prediction for query is based on the nearest neighbour. Here the graph representations were obtained by the Graph Isomorphism Network (GIN) model. To capture the global structure of the graph, they used different weights for different GIN layers in the final aggregation scheme. To encode the crucial local structures that might have importance in deciding the graph label, the paper embeds subgraphs and includes their representations with different attention weights.

Finally, Buffelli et al. [8] attempt to develop a framework that can adapt to three different tasks – graph classification, node classification, and link prediction. Like [11, 44] they use two different layers; one generates node embeddings and converts the graph to a representation, and another is a multi-head output layer for the three types of tasks. The node embedding layer is trained during the initialization phase of MAML and the multi-head output layer is updated in the inner loop of MAML based on the type of task.

Bose et al. [7] consider the few-shot link prediction problem, where the goal is to predict labels of links/edges that contain only a small fraction of their true labels. They assume that the graphs are generated from a common distribution  $p(\cdot)$  and learn a meta link prediction model that can be quickly adapted to a new graph  $G \sim p(\cdot)$ . In particular, the authors use Variational Graph Autoencoder (VGAE) [37] to model the base link prediction model. There are two sets of parameters – global initialization parameters for the VGAE, and local graph signature  $s_G = \psi(G)$  which is obtained by passing the graph  $G$  through GCN and then using a  $k$ -layer MLP. The training is done using MAML and only the graph signature is updated for the test graph.

## 6. OTHER APPLICATIONS

We have discussed applications of meta-learning equipped with GNNs on node classification, link prediction, and graph classification. In fact, this framework is quite general and can be applied to many other relevant important problems.

**Anomaly Detection:** The problem of anomaly detection often suffers from scarcity of labels, as obtaining labels for anomalies is usually labor intensive. Ding et al. [20] study anomaly detection when there are scarcity of labels, and different tasks involve different graphs. The proposed method used traditional architectures of GNNs to embed nodes and predict the anomaly score by adding another layer after the embedding is obtained. Finally it exploits the traditional MAML framework to deploy the meta-learner. The inner loop optimizes the parameters for a specific task, i.e., graph. The outer-loop optimizes the generic parameter for all graphs.

**Network Alignment (NA):** NA aims to map or link entities from different networks and relevant in many application domains such as cross-domain recommendation and advertising. Zhou et al. [75] address this alignment problem via meta-learning. If two different networks share some common nodes or anchors, then these networks are partially aligned networks. A virtual link between two anchors is called anchor link. In NA, given a set of networks and some known anchor nodes (or links), the goal is to identify all the other (unknown) potential anchor nodes (or links). The main idea in [75] is to frame this problem as one shot classification problem and use the meta-metric learning from known anchor nodes to obtain latent priors for linking unknown anchor nodes.

**Traffic Prediction:** Recently, the traffic prediction problem [50] has been addressed via meta-learning. In traffic prediction, the main challenges are modeling complex spatio-temporal correlations of traffic and capturing the diversity of such correlations varying locations. Pan et al. [50] address these challenges with a meta-learning based model. Their method predicts traffic in all locations at the same time. The proposed framework consists of a sequence-to-sequence architecture that uses an an encoder to learn traffic history and a decoder to make predictions. For the encoder and decoder components a combination of graph attention networks and recurrent neural networks is used to model diverse spatial and temporal correlations respectively.

## 7. FUTURE DIRECTIONS

The application of meta-learning using GNNs for graph specific applications is a growing and exciting area of research. In this section, we suggest several future directions for research.

### 7.1 Combinatorial Optimization Problems on Graphs

Combinatorial optimization problems appearing in graphs have applications in many domains such as viral marketing in social networks [34], health-care [65], and infrastructure development [47], and several architectures based on GNNs have been proposed for solving them [15, 40, 26, 45]. These optimization problems are often NP-hard, and polynomial-time algorithms, with or without approximation guarantees, are often desirable and used in practice. However, some techniques [40, 45] based on GNNs need to generate candidate solution nodes/edges before generating the actual solution set. Note that, labels in the form of importance of each node in a solution set of these problems are often difficult to get. Meta-learning can be used when there are scarcity of labels. Furthermore, these combinatorial problems often share similar structures. For instance, the influence maximization problem [34] have similarity with the Max Cover problem. However, even performing a greedy iterative algorithm to generate solutions/labels for influence maximization problem is computationally expensive. The idea of using meta-learning in solving a harder combinatorial problem (unseen task) with a fewer node labels will be to learn on the easier problems (seen tasks) where labels can be generated at a lower cost. Solving combinatorial optimization problems on graphs via neural approaches has recently gained a lot of attention and we refer the readers to [10] for further reading.

### 7.2 Graph Mining Problems

There has been recent attempt to solve classical graph mining problems with GNNs. For instance, a popular problem is to learn similarity between two graphs, i.e., to find graph edit distance (similarity) between two graphs [2]. When the notion of similarity changes and there are not enough data to learn via a standard supervised learning method, can meta-learning be helpful? Another popular graph mining problem is detecting the Maximum Common Subgraph (MCS) between two input graphs with applications in biomedical analysis and malware detection. In drug design, common substructures in compounds can reduce the number of human-conducted experiments. However, MCS computation is NP-hard, and state-of-the-art exact MCS solvers are not scalable to large graphs. Designing learning based models [3] for the MCS problem while utilizing as few labeled MCS instances as possible remains to be a challenging task and meta-learning could be helpful in mitigating this challenge.

### 7.3 Theory

We point out several important theoretical questions in the context of meta learning with GNNs. The most natural question is understanding the benefits of transfer learning in GNNs. Garg et al. [25] and Scarselli et al. [54] have recently established generalization bounds for GNNs. On the other hand, in the context of meta-learning, Tripuraneni et al. [60] consider functions of the form  $f_j \cdot h$ , where  $f_j \in \mathcal{F}$  is the task-specific function and  $h$  is the shared function. Then the number of samples required in the meta-test phase grows as  $C(\mathcal{F})$ , which can be significantly lower than learning  $f_j \cdot h$  from scratch. It would be interesting to see if one can prove similar speedup results for GNNs by generalizing the results of [25] and [54]. Another interesting question is determining the right level of shared representation and figuring out the expressiveness of such structures. The seminal work of Xu et al. [68] proves that variants of GNNs such as GCN and GraphSAGE are no more discriminative than the Weisfeiler-Lehman (WL) test. Since GNNs for meta-learning further limit the type of architecture used, an interesting question is whether it comes with any additional cost on expressiveness. Finally, the methods discussed in Section 5 differ in one crucial way – whether they fine-tune and update the shared meta-parameter on a new task or whether they keep the shared meta-parameter fixed. Recently, Chua et al. [13] show that fine-tuning the meta-parameter could be beneficial in some situations, particularly when the number of samples on the new task is large. In the context of meta learning on GNNs, it would be interesting to understand when such fine-tuning helps to improve the performance on a new task.

### 7.4 Applications

We have already discussed a few applications of meta-learning with frameworks of GNNs in Section 6. This generic framework is quite relevant for many important problems in the field.

**Network alignment:** A potential problem where meta-learning could be helpful is network alignment (NA) [75]. In NA, the main goal is to map or link entities from different networks and the existing approaches is quite difficult to scale. An interesting direction of research would consider meta-learning to overcome this scalability challenge.

**Molecular property prediction:** GNNs have been also used in predicting molecular properties. However, one of the main challenges is that molecules are heterogeneous structure where each atom has connection with different neighboring atoms via different types of bonds. Secondly, often a limited amount of data on labeled molecular property are available; and thus, to predict new molecular properties, meta-learning techniques [27] can be relevant and effective.

**Dynamic graphs:** In many applications, graphs arise with their dynamic nature, i.e., nodes and edges along with their attributes can change (addition or deletion) over time. Most of the papers discussed above use frameworks that are built on meta-learning and GNNs for static graphs. An interesting direction would be to extend this framework for dynamic graphs. Dynamic nature brings new challenges such as difficulty in obtaining labels for newly added nodes or edges. For instance, in knowledge graphs, newly added edges introduces new relationships. The other challenge is efficiency as managing and making predictions on evolving networks are difficult tasks as its own. Meta-learning would be useful to address these challenges.

## 8. CONCLUSION

In this survey, we have performed a comprehensive review of the

works that are combination of graph neural networks (GNNs) and meta-learning. Besides outlining backgrounds on GNNs and meta-learning, we have organized the past research in an organized manner in multiple categories. We have also provided a thorough review, summary of methods, and applications in these categories. Furthermore, we have described several future research directions where meta learning with GNNs can be useful. The application of meta-learning to GNNs is a growing and exciting field and we believe many graph problems will benefit immensely from the combination of the two approaches.

## References

- [1] Jinheon Baek, Dong Bok Lee, and Sung Ju Hwang. “Learning to extrapolate knowledge: Transductive few-shot out-of-graph link prediction”. In: *NeurIPS* (2020).
- [2] Yunsheng Bai, Hao Ding, Song Bian, Ting Chen, Yizhou Sun, and Wei Wang. “SimGNN: A Neural Network Approach to Fast Graph Similarity Computation”. In: *WSDM*. 2019.
- [3] Yunsheng Bai, Derek Xu, Alex Wang, Ken Gu, Xueqing Wu, Agustin Marinovic, Christopher Ro, Yizhou Sun, and Wei Wang. “Fast detection of maximum common subgraph via deep q-learning”. In: *arXiv preprint arXiv:2002.03129* (2020).
- [4] Maria-Florina Balcan, Mikhail Khodak, and Ameet Talwalkar. “Provable guarantees for gradient-based meta-learning”. In: *ICML*. 2019, pp. 424–433.
- [5] Jonathan Baxter. “A model of inductive bias learning”. In: *JAIR* 12 (2000), pp. 149–198.
- [6] Karsten M Borgwardt, Cheng Soon Ong, Stefan Schönauer, SVN Vishwanathan, Alex J Smola, and Hans-Peter Kriegel. “Protein function prediction via graph kernels”. In: *Bioinformatics* 21 (2005), pp. i47–i56.
- [7] Avishek Joey Bose, Ankit Jain, Piero Molino, and William L Hamilton. “Meta-graph: Few shot link prediction via meta learning”. In: *arXiv preprint arXiv:1912.09867* (2019).
- [8] Davide Buffelli and Fabio Vandin. “A Meta-Learning Approach for Graph Representation Learning in Multi-Task Settings”. In: *arXiv preprint arXiv:2012.06755* (2020).
- [9] Shaosheng Cao, Wei Lu, and Qiongkai Xu. “Deep neural networks for learning graph representations”. In: *AAAI* 30.1 (2016).
- [10] Quentin Cappart, Didier Chételat, Elias Khalil, Andrea Lodi, Christopher Morris, and Petar Veličković. “Combinatorial optimization and reasoning with graph neural networks”. In: *arXiv preprint arXiv:2102.09544* (2021).
- [11] Jatin Chauhan, Deepak Nathani, and Manohar Kaul. “Few-Shot Learning on Graphs via Super-Classes based on Graph Spectral Measures”. In: *ICLR* (2020).
- [12] Mingyang Chen, Wen Zhang, Wei Zhang, Qiang Chen, and Huajun Chen. “Meta Relational Learning for Few-Shot Link Prediction in Knowledge Graphs”. In: *EMNLP-IJCNLP*. 2019, pp. 4208–4217.
- [13] Kurtland Chua, Qi Lei, and Jason D Lee. “How fine-tuning allows for effective meta-learning”. In: *arXiv:2105.02221* (2021).
- [14] Zhiyong Cui, Kristian Henrickson, Ruimin Ke, and Yin Hai Wang. “Traffic graph convolutional recurrent neural network: A deep learning framework for network-scale traffic learning and forecasting”. In: *IEEE TITS* (2019), pp. 4883–4894.

- [15] Hanjun Dai, Elias B Khalil, Yuyu Zhang, Bistra Dilkina, and Le Song. “Learning combinatorial optimization algorithms over graphs”. In: *NeurIPS* (2017).
- [16] Giulia Denevi, Carlo Ciliberto, Riccardo Grazi, and Massimiliano Pontil. “Learning-to-learn stochastic gradient descent with biased regularization”. In: *ICML*. 2019, pp. 1566–1575.
- [17] Jacob Devlin, Ming-Wei Chang, Kenton Lee, and Kristina Toutanova. “Bert: Pre-training of deep bidirectional transformers for language understanding”. In: *NAACL-HLT* (2019).
- [18] Kaize Ding, Jianling Wang, Jundong Li, James Caverlee, and Huan Liu. “Weakly-supervised Graph Meta-learning for Few-shot Node Classification”. In: *arXiv:2106.06873* (2021).
- [19] Kaize Ding, Jianling Wang, Jundong Li, Kai Shu, Chenghao Liu, and Huan Liu. “Graph prototypical networks for few-shot learning on attributed networks”. In: *CIKM*. 2020, pp. 295–304.
- [20] Kaize Ding, Qinghai Zhou, Hanghang Tong, and Huan Liu. “Few-Shot Network Anomaly Detection via Cross-Network Meta-Learning”. In: *Proceedings of the Web Conference 2021*. 2021, 2448–2456.
- [21] Simon S Du, Wei Hu, Sham M Kakade, Jason D Lee, and Qi Lei. “Few-shot learning via learning the representation, provably”. In: *arXiv preprint arXiv:2002.09434* (2020).
- [22] Federico Errica, Marco Podda, Davide Bacciu, and Alessio Micheli. “A fair comparison of graph neural networks for graph classification”. In: *arXiv preprint arXiv:1912.09893* (2019).
- [23] Chelsea Finn, Pieter Abbeel, and Sergey Levine. “Model-agnostic meta-learning for fast adaptation of deep networks”. In: *ICML*. 2017, pp. 1126–1135.
- [24] Chelsea Finn, Aravind Rajeswaran, Sham Kakade, and Sergey Levine. “Online meta-learning”. In: *ICML*. 2019, pp. 1920–1930.
- [25] Vikas Garg, Stefanie Jegelka, and Tommi Jaakkola. “Generalization and representational limits of graph neural networks”. In: *ICML*. PMLR. 2020, pp. 3419–3430.
- [26] Maxime Gasse, Didier Chételat, Nicola Ferroni, Laurent Charlin, and Andrea Lodi. “Exact combinatorial optimization with graph convolutional neural networks”. In: *NeurIPS* (2019).
- [27] Zhichun Guo, Chuxu Zhang, Wenhao Yu, John Herr, Olaf Wiest, Meng Jiang, and Nitesh V Chawla. “Few-Shot Graph Learning for Molecular Property Prediction”. In: *The Web Conference* (2021).
- [28] Will Hamilton, Zitao Ying, and Jure Leskovec. “Inductive representation learning on large graphs”. In: *NeurIPS*. 2017, pp. 1024–1034.
- [29] William L Hamilton, Rex Ying, and Jure Leskovec. “Representation learning on graphs: Methods and applications”. In: *arXiv preprint arXiv:1709.05584* (2017).
- [30] Timothy Hospedales, Antreas Antoniou, Paul Micaelli, and Amos Storkey. “Meta-learning in neural networks: A survey”. In: *arXiv preprint arXiv:2004.05439* (2020).
- [31] Kexin Huang and Marinka Zitnik. “Graph meta learning via local subgraphs”. In: *NeurIPS* (2020).
- [32] Dasol Hwang, Jinyoung Park, Sunyoung Kwon, Kyung-Min Kim, Jung-Woo Ha, and Hyunwoo J Kim. “Self-supervised Auxiliary Learning for Graph Neural Networks via Meta-Learning”. In: *arXiv preprint arXiv:2103.00771* (2021).
- [33] Shunyu Jiang, Fuli Feng, Weijian Chen, Xiang Li, and Xiangan He. “Structure-Enhanced Meta-Learning For Few-Shot Graph Classification”. In: *arXiv preprint arXiv:2103.03547* (2021).
- [34] David Kempe, Jon Kleinberg, and Éva Tardos. “Maximizing the spread of influence through a social network”. In: *KDD*. 2003.
- [35] M Khodak, M Balcan, and A Talwalkar. “Adaptive Gradient-Based Meta-Learning Methods”. In: *Neural Information Processing Systems*. 2019.
- [36] Thomas N Kipf and Max Welling. “Semi-supervised classification with graph convolutional networks”. In: *ICLR* (2017).
- [37] Thomas N Kipf and Max Welling. “Variational graph auto-encoders”. In: *arXiv preprint arXiv:1611.07308* (2016).
- [38] Lin Lan, Pinghui Wang, Xuefeng Du, Kaikai Song, Jing Tao, and Xiaohong Guan. “Node classification on graphs with few-shot novel labels via meta transformed network embedding”. In: *NeurIPS* (2020).
- [39] Yujia Li, Daniel Tarlow, Marc Brockschmidt, and Richard Zemel. “Gated graph sequence neural networks”. In: *ICLR* (2016).
- [40] Zhuwen Li, Qifeng Chen, and Vladlen Koltun. “Combinatorial optimization with graph convolutional networks and guided tree search”. In: *NeurIPS* (2018).
- [41] Xiaodong Liu, Pengcheng He, Weizhu Chen, and Jianfeng Gao. “Multi-Task Deep Neural Networks for Natural Language Understanding”. In: *ACL*. 2019.
- [42] Zemin Liu, Yuan Fang, Chenghao Liu, and Steven CH Hoi. “Relative and Absolute Location Embedding for Few-Shot Node Classification on Graph”. In: *AAAI* (2021).
- [43] Zemin Liu, Wentao Zhang, Yuan Fang, Xinming Zhang, and Steven CH Hoi. “Towards locality-aware meta-learning of tail node embeddings on networks”. In: *CIKM*. 2020, pp. 975–984.
- [44] Ning Ma, Jiajun Bu, Jieyu Yang, Zhen Zhang, Chengwei Yao, Zhi Yu, Sheng Zhou, and Xifeng Yan. “Adaptive-Step Graph Meta-Learner for Few-Shot Graph Classification”. In: *Proceedings of the 29th ACM International Conference on Information & Knowledge Management*. 2020, 1055–1064.
- [45] Sahil Manchanda, Akash Mittal, Anuj Dhawan, Sourav Medya, Sayan Ranu, and Ambuj Singh. “GCOMB: Learning Budget-constrained Combinatorial Algorithms over Billion-sized Graphs”. In: *NeurIPS* (2020).
- [46] Andreas Maurer, Massimiliano Pontil, and Bernardino Romera-Paredes. “The benefit of multitask representation learning”. In: *Journal of Machine Learning Research* 17.81 (2016), pp. 1–32.
- [47] Sourav Medya, Jithin Vachery, Sayan Ranu, and Ambuj Singh. “Noticeable network delay minimization via node upgrades”. In: *VLDB* (2018).
- [48] Anusha Nagabandi, Kurt Konolige, Sergey Levine, and Vikash Kumar. “Deep dynamics models for learning dexterous manipulation”. In: *CoRL*. 2020.
- [49] Sunil Nishad, Shubhangi Agarwal, Arnab Bhattacharya, and Sayan Ranu. “GraphReach: Locality-Aware Graph Neural Networks using Reachability Estimations”. In: *IJCAI*. 2021.

- [50] Zheyi Pan, Wentao Zhang, Yuxuan Liang, Weinan Zhang, Yong Yu, Junbo Zhang, and Yu Zheng. “Spatio-Temporal Meta Learning for Urban Traffic Prediction”. In: *TKDE* (2020).
- [51] Massimiliano Pontil and Andreas Maurer. “Excess risk bounds for multitask learning with trace norm regularization”. In: *COLT*. 2013, pp. 55–76.
- [52] Sachin Ravi and Hugo Larochelle. “Optimization as a model for few-shot learning”. In: *ICLR*. 2017.
- [53] Franco Scarselli, Marco Gori, Ah Chung Tsoi, Markus Hagenbuchner, and Gabriele Monfardini. “The graph neural network model”. In: *IEEE transactions on neural networks* 20.1 (2008), pp. 61–80.
- [54] Franco Scarselli, Ah Chung Tsoi, and Markus Hagenbuchner. “The Vapnik–Chervonenkis dimension of graph and recursive neural networks”. In: *Neural Networks* 108 (2018), pp. 248–259.
- [55] Jürgen Schmidhuber. “Evolutionary principles in self-referential learning, or on learning how to learn: the meta-meta-... hook”. PhD thesis. Technische Universität München, 1987.
- [56] Oleksandr Shchur, Maximilian Mumme, Aleksandar Bojchevski, and Stephan Günnemann. “Pitfalls of graph neural network evaluation”. In: *arXiv preprint arXiv:1811.05868* (2018).
- [57] Jonathan M Stokes, Kevin Yang, Kyle Swanson, Wengong Jin, Andres Cubillos-Ruiz, et al. “A deep learning approach to antibiotic discovery”. In: *Cell* (2020), pp. 688–702.
- [58] Jie Tang, Jing Zhang, Limin Yao, Juanzi Li, Li Zhang, and Zhong Su. “Arnetminer: extraction and mining of academic social networks”. In: *KDD*. 2008.
- [59] Nilesh Tripuraneni, Chi Jin, and Michael I Jordan. “Provable meta-learning of linear representations”. In: *arXiv preprint arXiv:2002.11684* (2020).
- [60] Nilesh Tripuraneni, Michael Jordan, and Chi Jin. “On the Theory of Transfer Learning: The Importance of Task Diversity”. In: *NeurIPS* 33 (2020).
- [61] Petar Veličković, Guillem Cucurull, Arantxa Casanova, Adriana Romero, Pietro Lio, and Yoshua Bengio. “Graph attention networks”. In: *ICLR* (2018).
- [62] Ning Wang, Minnan Luo, Kaize Ding, Lingling Zhang, Jun-dong Li, and Qinghua Zheng. “Graph Few-shot Learning with Attribute Matching”. In: *CIKM*. 2020, pp. 1545–1554.
- [63] Yan Wang, Wei-Lun Chao, Kilian Q Weinberger, and Laurens van der Maaten. “SimpleShot: Revisiting nearest-neighbor classification for few-shot learning”. In: *arXiv:1911.04623* (2019).
- [64] Zhihao Wen, Yuan Fang, and Zemin Liu. “Meta-Inductive Node Classification across Graphs”. In: *arXiv:2105.06725* (2021).
- [65] Bryan Wilder, Han Ching Ou, Kayla de la Haye, and Milind Tambe. “Optimizing Network Structure for Preventative Health”. In: *AAMAS*. 2018.
- [66] Nan Wu, Jason Phang, Jungkyu Park, Yiqiu Shen, et al. “Deep neural networks improve radiologists’ performance in breast cancer screening”. In: *IEEE transactions on medical imaging* 39.4 (2019), pp. 1184–1194.
- [67] Zonghan Wu, Shirui Pan, Fengwen Chen, Guodong Long, Chengqi Zhang, and S Yu Philip. “A comprehensive survey on graph neural networks”. In: *IEEE transactions on neural networks and learning systems* (2020).
- [68] Keyulu Xu, Weihua Hu, Jure Leskovec, and Stefanie Jegelka. “How powerful are graph neural networks?” In: *ICLR* (2018).
- [69] Pinar Yanardag and SVN Vishwanathan. “Deep graph kernels”. In: *Proceedings of the 21th ACM SIGKDD international conference on knowledge discovery and data mining*. 2015, pp. 1365–1374.
- [70] Huaxiu Yao, Chuxu Zhang, Ying Wei, Meng Jiang, Suhang Wang, Junzhou Huang, Nitesh Chawla, and Zhenhui Li. “Graph few-shot learning via knowledge transfer”. In: *AAAI*. 2020, pp. 6656–6663.
- [71] Jiaxuan You, Rex Ying, and Jure Leskovec. “Position-aware Graph Neural Networks”. In: *ICML*. 2019, pp. 7134–7143.
- [72] Jiani Zhang, Xingjian Shi, Junyuan Xie, Hao Ma, Irwin King, and Dit-Yan Yeung. “Gaan: Gated attention networks for learning on large and spatiotemporal graphs”. In: *UAI* (2018).
- [73] Muhan Zhang and Yixin Chen. “Link prediction based on graph neural networks”. In: *NeurIPS* (2018).
- [74] Xi Sheryl Zhang, Fengyi Tang, Hiroko H Dodge, Jiayu Zhou, and Fei Wang. “Metapred: Meta-learning for clinical risk prediction with limited patient electronic health records”. In: *KDD*. 2019.
- [75] Fan Zhou, Chengtai Cao, Goce Trajcevski, Kunpeng Zhang, Ting Zhong, and Ji Geng. “Fast network alignment via graph meta-learning”. In: *INFOCOM*. 2020, pp. 686–695.
- [76] Fan Zhou, Chengtai Cao, Kunpeng Zhang, Goce Trajcevski, Ting Zhong, and Ji Geng. “Meta-Gnn: On Few-Shot Node Classification in Graph Meta-Learning”. In: *CIKM*. 2019, pp. 2357–2360.
- [77] Jie Zhou, Ganqu Cui, Zhengyan Zhang, Cheng Yang, Zhiyuan Liu, Lifeng Wang, Changcheng Li, and Maosong Sun. “Graph neural networks: A review of methods and applications”. In: *arXiv preprint arXiv:1812.08434* (2018).