

Explicit Semantic Cross Feature Learning via Pre-trained Graph Neural Networks for CTR Prediction

Feng Li*, Bencheng Yan*, Qingqing Long, Pengjie Wang[‡], Wei Lin, Jian Xu and Bo Zheng[†]

Alibaba Group

{adam.lf, bencheng.ybc, lantu.lqq, pengjie.wpj, xiyu.xj, bozheng}@alibaba-inc.com, lwsaviola@163.com

ABSTRACT

Cross features play an important role in click-through rate (CTR) prediction. Most of the existing methods adopt a DNN-based model to capture the cross features in an implicit manner. These implicit methods may lead to a sub-optimized performance due to the limitation in explicit semantic modeling. Although traditional statistical explicit semantic cross features can address the problem in these implicit methods, such features still suffer from some challenges, including *lack of generalization* and *expensive memory cost*. Few works focus on tackling these challenges. In this paper, we take the first step in learning the explicit semantic cross features and propose Pre-trained Cross Feature learning Graph Neural Networks (PCF-GNN), a GNN based pre-trained model aiming at generating cross features in an explicit fashion. Extensive experiments are conducted on both public and industrial datasets, where PCF-GNN shows competence in both performance and memory-efficiency in various tasks.

CCS CONCEPTS

• Information systems → Recommender systems.

KEYWORDS

CTR prediction; Pre-trained GNNs; Cross Features; Explicit Fashion

ACM Reference Format:

Feng Li*, Bencheng Yan*, Qingqing Long, Pengjie Wang[‡], Wei Lin, Jian Xu and Bo Zheng[†]. 2021. Explicit Semantic Cross Feature Learning via Pre-trained Graph Neural Networks for CTR Prediction. In *Proceedings of the 44th International ACM SIGIR Conference on Research and Development in Information Retrieval (SIGIR '21)*, July 11–15, 2021, Virtual Event, Canada. ACM, New York, NY, USA, 5 pages. <https://doi.org/10.1145/3404835.3463015>

1 INTRODUCTION

Modeling the complex relationships among features is the key to click-through rate (CTR) prediction. A major way to characterize such relationships is to introduce the cross features. For example, a

user who is a basketball player likes to click the recommended item "Nike-Air Jordan". While a user who is a programmer may be more likely to click digital products. These suggest that the interactions (i.e., cross features) $\langle user_occupation, item \rangle$ can be taken as a strong signal for CTR prediction. In practice, such cross features have achieved remarkable results in modeling feature interactions and improving the CTR performance [15, 17].

There are two ways to model the cross features, i.e., implicit semantic modeling and explicit semantic modeling. Most of the existing works focus on implicit semantic modeling. The main idea of these works is to adopt a deep neural network (DNN) and hope the capacity of DNN can characterize the cross features in an implicit manner (e.g., Wide & Deep [2], DeepFM [4], DCN[16], Cross-GCN [3], and Fi-GNN [9]). Although this kind of methods has achieved great success in cross features modeling, it is restricted to the implicit semantic modeling, and we can never guarantee the learned cross features are what we want, leading to a sub-optimized performance (Note in these methods, although the involved two features may be explicitly fed into a same DNN-based layer, the learned cross features can only contain implicit semantic information.). Therefore explicit semantic information is an important signal to make up the limitation of the above methods.

However, few works focus on learning the explicit semantic information of cross features effectively. On the contrary, almost all the existing methods adopt simple statistical explicit semantic cross features (SESCF) manually to describe the explicit semantic information [13]. Specifically, SESCOF can be modeled by counting how many clicks or calculating the click rate among features from history [13]. For example, we can count the times that a user who is a basketball player clicks "Nike-Air Jordan" in the history, and take this counting value as the explicit semantic cross feature $\langle user_occupation=basketball\ player, item=Nike-Air\ Jordan \rangle$ for the samples having the same user occupation and items. Obviously, such a counting value describes the feature relationship explicitly. A higher value refers to a higher correlation among these features. However, simply adopting such SESCOF poses great challenges.

(1) Lack of generalization. SESCOF mainly relies on the statistical counting from history and has no ability to infer the new pairs of cross features which are never shown before. For example, the latest "Nike-Air Jordan" may be heavily browsed by the user who is a basketball player recently, and rarely shown to a user who is a programmer. In this case, the SESCOF about $\langle user_occupation, item \rangle$ cannot give any suggestions for the recommendation of the latest "Nike-Air Jordan" to a programmer user. Hence, there is a large requirement for predicting new pairs and a generalized model on explicit semantic cross features still remains a significant challenge.

* These authors contributed equally to this work and are co-first authors.

[‡] This author gave a lot of guidance in this work.

[†] Corresponding author.

Permission to make digital or hard copies of all or part of this work for personal or classroom use is granted without fee provided that copies are not made or distributed for profit or commercial advantage and that copies bear this notice and the full citation on the first page. Copyrights for components of this work owned by others than the author(s) must be honored. Abstracting with credit is permitted. To copy otherwise, to republish, to post on servers or to redistribute to lists, requires prior specific permission and/or a fee. Request permissions from permissions@acm.org.

SIGIR '21, July 11–15, 2021, Virtual Event, Canada

© 2021 Copyright held by the owner/author(s). Publication rights licensed to ACM.

ACM ISBN 978-1-4503-8037-9/21/07...\$15.00

<https://doi.org/10.1145/3404835.3463015>

(2) Expensive memory cost. In real applications, we usually need to keep a *<cross feature, counting value>* (e.g., «basketball player, "Nike-Air Jordan">, 100) mapping table for online severing. The challenge is that maintaining such a mapping table may need excessively vast storage resources especially for a web-scale application [2, 17]. For instance, if there are N_1 occupations and N_2 items, the number of rows in the mapping table is related to $N_1 N_2$. In other words, the corresponding memory cost is $O(N_1 N_2)$ (i.e., Cartesian product). Therefore, modeling the interaction among features which have a large-scale vocabulary (e.g., User ID) may need huge memory cost which definitely hurts the severing efficiency.

None of the existing works try to address the above challenges in the view of the explicit semantic modeling. To tackle these challenges, we propose a Pre-trained Cross Feature learning Graph Neural Networks (PCF-GNN), a GNN [11] based pre-trained model. Specifically, we first transform the interactions among features in history to an informative graph. Because abundant interactions can be naturally modeled as edges in a graph, and the graph can further provide us the convenience to understand the complex interactions. In this graph, each node refers to a single feature (e.g., user_occupation), each edge refers to a feature interaction, and the edge attributes are assigned with the counting value (i.e., SESCOF). Then a self-supervised edge-attributes prediction task is designed to pre-train a GNN (a most powerful architecture in node relationship modeling [6, 12]) to allow the model to have the domain knowledge about the explicit semantic information of cross features and have the ability to infer the attributes of new edges (addressing challenge one). Furthermore, such a design makes us free from storing the huge mapping table since we can simply infer rather than store the explicit semantic cross features (i.e., edge attributes) (addressing challenge two). Detailed analysis can be found in Section 2.5.

We summarize our contributions as follows: (1) We propose PCF-GNN, a pre-trained GNN model to learn explicit semantic cross features. To the best of our knowledge, we take the first step in learning the explicit semantic cross features, and the pre-trained GNN is also first introduced in this area. (2) A pre-training pipeline is designed and a weighted loss based on interactions between features is proposed, to capture explicit semantic cross features. (3) We carry out extensive experiments on both public and industrial datasets, where PCF-GNN shows competence in both performance and memory-efficiency.

2 MODEL: PCF-GNN

In this section, we introduce our model. Generally speaking, we adopt a two-stage pipeline, i.e., pre-training stage (Section 2.1–2.3) and downstream task application stage (Section 2.4).

2.1 Graph Construction

We first construct the graph based on the interactions among features in history. As shown in Fig. 1, given samples of users' click history, we take the interactions *<User, Item>* and *<Item, Shop>* as an example (in practice, abundant kinds of interactions are considered), and construct an informative graph where nodes refer to features, edges refer to feature interactions and the edge attributes refer to statistical explicit semantic cross features. Formally, in this paper, the attribute $a_{u,i}$ of edge $e_{u,i}$ can be calculated as,

$$a_{u,i} = \text{Count}(u, i) | (\text{click} = 1) / \text{Count}(u, i) \quad (1)$$

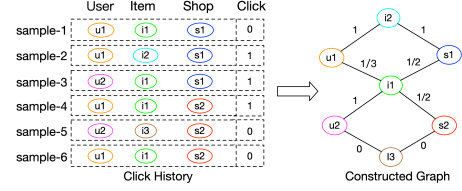


Figure 1: An example graph constructed from history.

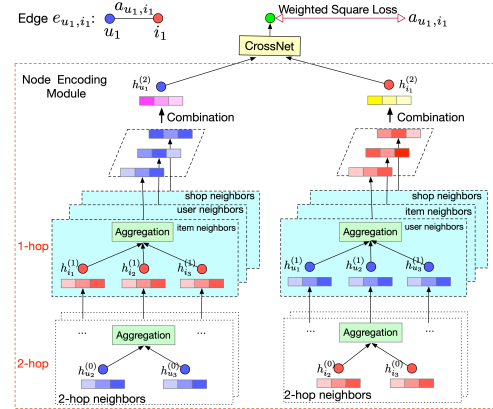


Figure 2: An overview of the proposed PCF-GNN.

where $\text{Count}(u, i)$ denotes the number of co-occurrence that features u and i are presented in the same history samples. Note $a_{u,i}$ characterizes the interactions between features u and i explicitly.

2.2 Pre-training Task Design

The key idea is that we need to give the domain knowledge about the explicit semantic cross features to the pre-training task. Specifically, since the edge attributes in the constructed graph refer to the explicit semantic cross features, we can design a self-supervised task, i.e., predicting the edge attributes. Formally, the pre-training task can be represented as $\|p_{u,v} - a_{u,v}\|^2$ where $p_{u,v}$ is the output of the PCF-GNN. In this way, with the guidance of the edge attribute, the learned $p_{u,v}$ can explicitly characterize the cross features.

2.3 Architecture Design

As shown in Fig 2, the architecture of PCF-GNN can be divided into two parts, i.e., node encoding module and CrossNet.

2.3.1 Node Encoding Module. This module is a GNN which can capture the neighbor structure and learn the node embeddings. Specifically, considering the constructed graph is a heterogeneous graph, we propose a multi-relation based aggregation function,

$$m_{i,r}^{(k)} = \text{AGGREGATE} \left(\left\{ h_j^{(k-1)}, j \in N_r(i) \right\} \right) \quad (2)$$

$$h_i^{(k)} = \text{COMBINATION} \left(h_i^{(k-1)}, \left\{ m_{i,r}^{(k)}, r \in [1, R] \right\} \right), \quad (3)$$

where $h_i^{(k)}$ denotes the output vector for node i at the k -th layer, $N_r(i)$ refers to the neighbors have the relationship r with the node i . $m_{i,r}^{(k)}$ refers to the embedding of $N_r(i)$. The $h_i^{(0)}$ are the attributes of nodes. For simplicity, in practice, the node attributes are defined as a learnable vector. The $h_i^{(K)}$ refers to the output of node i at the last layer of the node encoding module. Actually, we take $h_i^{(K)}$ as the pre-trained embedding of node i . In this paper, We take identical

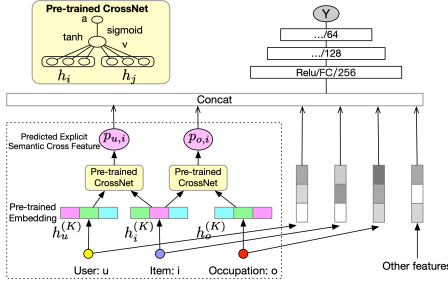


Figure 3: PCF-GNN application on the CTR prediction task. AGGREGATE and COMBINATION as GraphSAGE [5] while being flexible to more sophisticated methods.

2.3.2 CrossNet. The CrossNet transforms the neighbor aware node-level embedding $h_i^{(K)}$ into a edge (i.e., cross feature) space and gives a prediction for the edge attribute. Formally, for each edge $e_{u,v}$, the edge attribute can be predicted by,

$$p_{u,v} = \text{CrossNet}(h_u^{(K)}, h_v^{(K)}) \quad (4)$$

where $p_{u,v}$ is the predicted edge attribute, and **CrossNet can be a simple dot production or a multilayer perceptron (MLP). In this paper, we implement CrossNet as a single perceptron layer.**

2.3.3 Weighted Square Loss. When obtaining the edge attribute prediction $p_{u,v}$, we can simply apply a square loss to train the PCF-GNN, i.e., $\text{loss} = \sum_{e_{u,v}} ||p_{u,v} - a_{u,v}||^2$. The problem is that, in this loss, each edge takes as an equal contribution to the loss optimization. However, different edges take different importance in CTR tasks. For example, the pair *<basketball player, "Nike-Air Jordan">* may appear more times than the pair *<basketball player, mouse>* in history. If we treat these two pairs equally in loss calculation, the learned PCF-GNN may be difficult to distinguish which kind of cross features is important for the interaction modeling and may mislead the CTR model to recommend a mouse to a basketball player. To address this problem, we propose a simple but effective weighted square loss. The key idea is allocating different weight to different edges by their co-occurrence,

$$\text{loss} = \sum_{e_{u,v}} \log(\text{Count}(u, v) + t) ||p_{u,v} - a_{u,v}||^2 \quad (5)$$

where \log is a smooth operation which is used to avoid an extremely large $\text{Count}(u, v)$, and t is a small positive constant ($t = 1$ as default) to avoid a zero weight produced by \log .

2.4 PCF-GNN Application

When PCF-GNN is pre-trained, the learned knowledge about the explicit semantic information can be well characterized by $p_{u,v}$. Then we can use this knowledge to help the downstream task, especially for the CTR prediction task. Here we take the CTR prediction task as an example to show how to combine PCF-GNN with downstream models. As shown in Fig 3, the right part is a standard CTR model which follows an Embedding&MLP architecture [2, 4, 13]. Then we can take the predicted $p_{u,v}$ as an additional feature, and concatenate it with the embeddings of the CTR model (as shown in the dotted part in Fig 3). In this stage, by default, we drop the node encoding module for the efficient purpose and directly keep the output $h_i^{(K)}$ of this module as the pre-trained embedding and

the parameters in PCF-GNN are fixed. (Note when considering the *fine-tuning* strategy, the complete PCF-GNN including the node encoding module needs to be introduced in the CTR model). In this way, the explicit semantic cross features learned by PCF-GNN can help to improve the performance of the CTR model.

2.5 Discussion

In this section, we provide a brief discussion of our model on addressing the two challenges mentioned in Section 1.

Generalization. In PCF-GNN, the explicit semantic cross features (i.e., edge attribute) are inferred by pre-trained node embeddings. It can infer the explicit semantic cross features of new pairs (i.e., new edges). For example, even if there have no interactions between a programmer and the latest "Nike-Air Jordan" in history, PCF-GNN can still infer the new cross pair *<programmer, the latest "Nike-Air Jordan">* by the pre-trained node embeddings of the programmer and the latest "Nike-Air Jordan", which finally can help the recommendation of the latest "Nike-Air Jordan" to a programmer user.

Low memory cost. As introduced in Section 1, theoretically, the memory cost of statistical explicit semantic cross features is $O(N_1 N_2)$. While the main cost of inferred explicit cross features comes from the pre-trained node embedding table, i.e., $O((N_1 + N_2)d)$ where d is the embedding dimension. Considering N_1 and N_2 are usually a large number in the web-scale application, we can reduce the memory cost from $O(N_1 N_2)$ to $O((N_1 + N_2)d)$.

3 EXPERIMENTS

3.1 Experiment Setup

Datasets. We conduct experiments on a large-scale industrial **Alibaba** dataset (containing 60 billion samples) collected by the online platform and a public **MovieLens** dataset (containing 1000209 samples). Note each dataset is split into pre-trained data used for constructing graphs, downstream training data, and downstream test data respectively.

Baselines. We compare PCF-GNN with the following baselines. (1) Implicit semantic modeling methods, including Wide&Deep [2], DeepFM [4], AutoInt [14], FiBiNET [7]. (2) GNN-based methods that capture implicit semantic cross features, including Fi-GNN [10]. (3) Graph pre-trained models, which are originally designed to capture neighbor structure rather than explicit semantic cross features, i.e., GraphSAGE [5] and PGCN [6]. Note methods in (1)(2) can be directly taken as the downstream models for the CTR task.

Hyperparameter Settings. We take 8-dimensional embeddings for all pre-trained methods. Parameters of other baselines follow recommended settings in relevant codes. For PCF-GNN, we take 2-layer networks with a hidden layer sized 64.

3.2 Performance on CTR Task

3.2.1 Performance on the public dataset. To investigate whether PCF-GNN captures explicit semantic cross features and further improves the performance, we conduct the experiments on the downstream CTR task where the additional features generated by SESCf or pre-trained models are combined into the downstream models (similar to Fig 3). Here we construct the graph based on the interactions between User_id and Genres on MovieLens. There is a total of 60,574 edges and 5,992 nodes. Then we use AUC [1] as the

Downstream Task	Downstream Model	No ESCF	SESCF	GraphSAGE	PGCN	PCF-GNN
CTR	Wide&Deep	0.7144	0.7253	0.7254	0.7259	0.7266
	DeepFM	0.7238	0.7268	0.7285	0.7288	0.7293
	FiBiNET	0.7171	0.7210	0.7251	0.7254	0.7261
	AutoInt	0.7238	0.7280	0.7276	0.7280	0.7282
	Fi-GNN	0.7256	0.7301	0.7277	0.7283	0.7315
Task1	LightGBM	0.7110	0.7107	0.7126	0.7127	0.7239
Task2	LightGBM	0.5791	0.5862	0.5839	0.5931	0.5990

Table 1: Results of different tasks on MoiveLens dataset.

Model	DNN	DNN + SESCf	DNN + PCF-GNN
AUC	0.7442	0.7466	0.7477

Table 2: CTR performance on industrial Alibaba dataset.

metric (Note it is noticeable that a slightly higher AUC at **0.001-level** is regarded significant for the CTR task [14, 17]). Results are reported in Table 1 (Note ESCF refers to explicit semantic cross features). We summarize the following findings: (1) The proposed PCF-GNN outperforms all baselines, which shows the effectiveness of our model. (2) Compared with No ESCF, even if most of the downstream models capture cross features implicitly, the models with ESCF (including SESCf and PCF-GNN) are able to improve CTR results, which shows the usefulness of explicit semantic modeling. (3) Compared with SESCf, PCF-GNN can further improve the CTR performance. It shows PCF-GNN has the superiority of explicit semantic modeling. (4) Compared with other pre-trained models (i.e., GraphSAGE and PGCN), although due to capture the structure of the constructed interaction-based graph they can improve CTR performance, PCF-GNN achieves the best performance. It illustrates the design of our pre-training task and our model can capture explicit semantic cross features better.

3.2.2 Performance on the in-house industrial dataset. We further conduct experiments on an industrial dataset (Alibaba). Here, we construct a heterogeneous graph based on 11 kinds of cross features. There is a total of 20 billion edges and 0.7 billion nodes. The results are reported in Table 2. Note DNN refers to Wide & Deep [2] which is taken as the CTR model. We can find that PCF-GNN still achieves the best performance. It indicates that PCF-GNN can well model the explicit semantic cross features.

3.3 Performance on Other Downstream Tasks

To show the universality of PCF-GNN, we evaluate PCF-GNN on various related downstream tasks. Detailedly, we select some features (i.e., age and occupation) on MovieLens as labels (Note other features can also be considered as labels) and conduct experiments on predicting the corresponding labels, i.e., predicting user’s age (defined as Task1) and user’s occupation (defined as Task2). We use LightGBM [8] as the downstream classifier, and accuracy as the metric. Results are shown in Table 1. We observe that PCF-GNN significantly improves the results of all downstream tasks, compared with other competitors. It shows the effectiveness of the learned explicit cross feature in various related downstream tasks.

3.4 Evaluation of Model Generalization

Here we conduct experiments to evaluate the generalization of models (addressing challenge one). Specifically, we first calculate the Hit Rate (HR) of explicit semantic cross features on test data of MovieLens where there may exist some new cross feature pairs.

Model	HR	New	Δ_{New}	Org	Δ_{Org}
DNN	–	0.6972	–	0.7144	–
DNN + SESCf	78.27%	0.6916	-0.0056	0.7253	0.0109
DNN + PCF-GNN	91.81%	0.7010	0.0038	0.7266	0.0122

Table 3: Results about generalization.

Method	SESCf	PCF-GNN	Change
Memory Cost (GB)	84	37	-56%

Table 4: Memory cost of different methods.

	Base	Base+GNN	Base+GNN+WL	Base+GNN+WL+FT
AUC	0.7256	0.7261	0.7266	0.7268

Table 5: Results of the ablation study.

Then we evaluate the AUC of different methods on the test samples which only contain new cross feature pair (defined as New). We also provide the results on the whole test samples (defined as Org) for comparison. The results are reported in Table 3. Note DNN refers to Wide & Deep [2] and Δ_{XX} means the gap with DNN on the test dataset XX. We can find that (1) Since SESCf cannot infer the explicit semantic cross feature of new pairs, SESCf can cover only 78.27% of test samples. While PCF-GNN can cover 91.81% of test samples. (2) SESCf shows a negative impact on the New dataset. While PCF-GNN can always achieve improvement both on the New and Org datasets. These demonstrate that PCF-GNN has a good generalization ability and can infer the explicit semantic cross feature of new pairs to improve the model performance.

3.5 Evaluation of Memory Cost

In this section, we conduct experiments on the industrial Alibaba dataset to analyze the memory cost of PCF-GNN (addressing challenge two). The memory cost of modeling the explicit semantic cross features of different methods is reported in Table 4. With the help of PCF-GNN, we can save 56% memory cost compared with SESCf which efficiently improves memory storage. Note when each feature has interactions with all other features, theoretically, the memory cost can be reduced from $O(N_1N_2)$ to $O((N_1 + N_2)d)$.

3.6 Ablation Study

We conduct experiments to analyze the influence of Weighted Loss (Section 2.3.3), Fine-Tune (Section 2.4), and the GNN architecture in the node encoding module (Section 2.3.1), to provide better understandings of PCF-GNN. Specifically, we conduct experiments on MovieLens with Wide & Deep as the CTR model to verify the analysis. Results are shown in Table 5 (Note WL refers to weighted loss and FT refers to fine-tune. Base refers to PCF-GNN without GNN, WL, and FT). We conclude that the weighted loss, Fine-tune, and GNN can improve the model performance.

4 CONCLUSION

In this paper, we propose PCF-GNN, a pre-trained recommendation model based on GNNs. We first transform the interactions among features in history into an informative graph. Then we introduce pre-trained GNN models to learn efficient and effective explicit relationships among features. We propose a novel training loss to capture explicit semantic cross features. Experiments on both online and public datasets show our model is competent in both performance and memory-efficiency.

REFERENCES

- [1] Andrew P Bradley. 1997. The use of the area under the ROC curve in the evaluation of machine learning algorithms. *Pattern recognition* 30, 7 (1997), 1145–1159.
- [2] Heng-Tze Cheng, Levent Koc, Jeremiah Harmsen, Tal Shaked, Tushar Chandra, Hrishu Aradhye, Glen Anderson, Greg Corrado, Wei Chai, Mustafa Ispir, et al. 2016. Wide & deep learning for recommender systems. In *Proceedings of the 1st workshop on deep learning for recommender systems*. 7–10.
- [3] Fuli Feng, Xiangnan He, Hanwang Zhang, and Tat-Seng Chua. 2020. Cross-GCN: Enhancing Graph Convolutional Network with k -Order Feature Interactions. *arXiv preprint arXiv:2003.02587* (2020).
- [4] Huifeng Guo, Ruiming Tang, Yunming Ye, Zhenguo Li, and Xiuqiang He. 2017. DeepFM: a factorization-machine based neural network for CTR prediction. *arXiv preprint arXiv:1703.04247* (2017).
- [5] Will Hamilton, Zhitaoying, and Jure Leskovec. 2017. Inductive representation learning on large graphs. In *Advances in neural information processing systems*. 1024–1034.
- [6] Weihua Hu, Bowen Liu, Joseph Gomes, Marinka Zitnik, Percy Liang, Vijay Pande, and Jure Leskovec. 2019. Strategies for Pre-training Graph Neural Networks. In *International Conference on Learning Representations*.
- [7] Tongwen Huang, Zhiqi Zhang, and Junlin Zhang. 2019. FiBiNET: combining feature importance and bilinear feature interaction for click-through rate prediction. In *Proceedings of the 13th ACM Conference on Recommender Systems*. 169–177.
- [8] Guolin Ke, Qi Meng, Thomas Finley, Taifeng Wang, Wei Chen, Weidong Ma, Qiwei Ye, and Tie-Yan Liu. 2017. Lightgbm: A highly efficient gradient boosting decision tree. *Advances in neural information processing systems* 30 (2017), 3146–3154.
- [9] Zekun Li, Zeyu Cui, Shu Wu, Xiaoyu Zhang, and Liang Wang. 2019. Fi-gnn: Modeling feature interactions via graph neural networks for ctr prediction. In *Proceedings of the 28th ACM International Conference on Information and Knowledge Management*. 539–548.
- [10] Bin Liu, Ruiming Tang, Yingzhi Chen, Jinkai Yu, Huifeng Guo, and Yuzhou Zhang. 2019. Feature generation by convolutional neural network for click-through rate prediction. In *The World Wide Web Conference*. 1119–1129.
- [11] Qingqing Long, Yilun Jin, Guojie Song, Yi Li, and Wei Lin. 2020. Graph Structural-topic Neural Network. In *Proceedings of the 26th ACM SIGKDD International Conference on Knowledge Discovery & Data Mining*. 1065–1073.
- [12] Qingqing Long, Yiming Wang, Lun Du, Guojie Song, Yilun Jin, and Wei Lin. 2019. Hierarchical Community Structure Preserving Network Embedding: A Subspace Approach. In *Proceedings of the 28th ACM International Conference on Information and Knowledge Management*. 409–418.
- [13] Ying Shan, T Ryan Hoens, Jian Jiao, Haijing Wang, Dong Yu, and JC Mao. 2016. Deep crossing: Web-scale modeling without manually crafted combinatorial features. In *Proceedings of the 22nd ACM SIGKDD international conference on knowledge discovery and data mining*. 255–262.
- [14] Weiping Song, Chence Shi, Zhiping Xiao, Zhijian Duan, Yewen Xu, Ming Zhang, and Jian Tang. 2019. AutoInt: Automatic feature interaction learning via self-attentive neural networks. In *Proceedings of the 28th ACM International Conference on Information and Knowledge Management*. 1161–1170.
- [15] Jizhe Wang, Pipei Huang, Huan Zhao, Zhibo Zhang, Binqiang Zhao, and Dik Lun Lee. 2018. Billion-scale commodity embedding for e-commerce recommendation in alibaba. In *Proceedings of the 24th ACM SIGKDD International Conference on Knowledge Discovery & Data Mining*. 839–848.
- [16] Ruoxi Wang, Bin Fu, Gang Fu, and Mingliang Wang. 2017. Deep & cross network for ad click predictions. In *Proceedings of the ADKDD'17*. 1–7.
- [17] Guorui Zhou, Xiaoqiang Zhu, Chenru Song, Ying Fan, Han Zhu, Xiao Ma, Yanghui Yan, Junqi Jin, Han Li, and Kun Gai. 2018. Deep interest network for click-through rate prediction. In *Proceedings of the 24th ACM SIGKDD International Conference on Knowledge Discovery & Data Mining*. ACM, 1059–1068.