

Modeling Dynamic Heterogeneous Graph and Node Importance for Future Citation Prediction

Hao Geng

SKLSDE, School of Computer Science
at Beihang University
Beijing, China
genghao@buaa.edu.cn

Deqing Wang*

SKLSDE, School of Computer Science
at Beihang University
Beijing, China
dqwang@buaa.edu.cn

Fuzhen Zhuang

SKLSDE, School of Computer Science
& Institute of Artificial Intelligence at
Beihang University
Beijing, China
zhuangfuzhen@buaa.edu.cn

Xuehua Ming

SKLSDE, School of Computer Science
at Beihang University
Beijing, China
xhming@buaa.edu.cn

Chenguang Du

SKLSDE, School of Computer Science
at Beihang University
Beijing, China
duchenguang@buaa.edu.cn

Ting Jiang

SKLSDE, School of Computer Science
at Beihang University
Beijing, China
royokong@buaa.edu.cn

Haolong Guo

SKLSDE, School of Computer Science
at Beihang University
Beijing, China
ghl_123@buaa.edu.cn

Rui Liu

SKLSDE, School of Computer Science
at Beihang University
Beijing, China
lr@buaa.edu.cn

ABSTRACT

Accurate citation count prediction of newly published papers could help editors and readers rapidly figure out the influential papers in the future. Though many approaches are proposed to predict a paper's future citation, most ignore the dynamic heterogeneous graph structure or node importance in academic networks. To cope with this problem, we propose a Dynamic heterogeneous Graph and Node Importance network (DGNI) learning framework, which fully leverages the dynamic heterogeneous graph and node importance information to predict future citation trends of newly published papers. First, a dynamic heterogeneous network embedding module is provided to capture the dynamic evolutionary trends of the whole academic network. Then, a node importance embedding module is proposed to capture the global consistency relationship to figure out each paper's node importance. Finally, the dynamic evolutionary trend embeddings and node importance embeddings calculated above are combined to jointly predict the future citation counts of each paper, by a log-normal distribution model according to multi-faced paper node representations. Extensive experiments on two large-scale datasets demonstrate that our model significantly improves all indicators compared to the SOTA models.

*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 ACM must be honored. Abstracting with credit is permitted. To copy otherwise, or republish, to post on servers or to redistribute to lists, requires prior specific permission and/or a fee. Request permissions from permissions@acm.org.
CIKM '22, October 17–21, 2022, Atlanta, GA, USA.

© 2022 Association for Computing Machinery.
ACM ISBN 978-1-4503-9236-5/22/10...\$15.00
<https://doi.org/10.1145/3511808.3557398>

CCS CONCEPTS

• Information systems → Data mining.

KEYWORDS

citation count prediction, dynamic heterogeneous graph, node importance estimation

ACM Reference Format:

Hao Geng, Deqing Wang, Fuzhen Zhuang, Xuehua Ming, Chenguang Du, Ting Jiang, Haolong Guo, and Rui Liu. 2022. Modeling Dynamic Heterogeneous Graph and Node Importance for Future Citation Prediction. In *Proceedings of the 31st ACM International Conference on Information and Knowledge Management (CIKM '22)*, October 17–21, 2022, Atlanta, GA, USA. ACM, New York, NY, USA, 10 pages. <https://doi.org/10.1145/3511808.3557398>

1 INTRODUCTION

Predicting the impact of research papers is of great significance for science researchers to find out the most promising research topic to study and identify significant works from a sea of scientific literature [7]. As there is no precise definition of the impact of a scientific research, citation counts of scientific papers are usually taken as the estimation [8, 26, 39].

However, the task of predicting citation counts is challenging and nontrivial, due to the following reasons. Firstly, as existing publications own citation counts in each year, when a new publication emerging, there does not exist any historical citations, leading to the lack of label information to train the model. Secondly, in the heterogeneous academic network, each paper is associated with heterogeneous information such as authors, venues and fields, and how to make full use of these heterogeneous information turns to be a great challenge. Thirdly, all nodes in the academic network keep continuously evolving with changeable states, making it crucial to properly capture the dynamics. To follow previous footsteps,

we summarize the existing methods for citation count prediction task in the following two categories.

The first category predicts early published papers' citation counts with previous citations available [1, 20, 25, 32, 36, 41]. One part of these methods design parametric patterns to model the paper citation trend, including the log-normal intensity function in Wang et al. [32], the reinforced Poisson process in Shen et al. [25] and the recency-weighted effect in Liu et al. [20]. Other methods utilize neural networks to capture the temporal patterns in historical citations, such as the Recurrent Neural Network (RNN) in Yuan et al. [41] and the seq2seq framework in Abrishami [1]. However, these methods are unable to predict citation counts for newly published papers, since there exists no historical citations.

The second category predicts citation counts for newly published papers without historical citations. As a paper carries multi-dimensional information including its related authors, keywords, reference papers and venues, some researchers extract hand-crafted features to represent a paper [4, 7, 38–40]. For instance, Dong et al. [7] represents a paper by 6 types of factors, including author, reference, topic, venue, social and temporal features. Yan et al. [39] extracts rank-based features such as author rank and venue rank, as a part of paper features. However, these feature engineering methods require expert knowledge and manual labour for feature designing, which is time-consuming and cannot utilize the power of heterogeneous and evolving characteristics of nodes. HINTS [15] is the first work as we known to design an end-to-end framework for predicting new paper citation counts without feature engineering. Specifically, HINTS takes the whole academic network as a sequence of heterogeneous graphs, and combines the temporally aligned Graph Neural Network (GNN), the Recurrent Neural Network (RNN) and a time series generator to learn representations of the sequence. Although creative and rational, it fails to make full use of heterogeneous academic network features and the significance or popularity of a node in the graph, and thus causing the prediction accuracy lower.

In this paper, we propose a new framework based on **Dynamic Heterogeneous Graph and Node Importance Network**, named **DGNI**, which models the dynamic evolutionary trends of the academic network and each paper's node importance on the global scale to predict future citation counts of newly-published papers. DGNI is divided into three parts: dynamic heterogeneous network embedding module, node importance embedding module, and time series generation module.

Firstly, in dynamic heterogeneous network embedding module, we use Heterogeneous Graph Neural Network (HGNN) on snapshots of the academic network from different timestamps, together with RNN-based model to jointly model time series features. So the module can capture the dynamic evolutionary trends of the whole academic network before the publication of the papers.

Since the heterogeneous graph neural network can only capture the local consistency relationship of academic network, neglecting the significance of the global consistency relationship. In node importance embedding module, we propose a node importance embedding module to calculate each paper's node importance on the global scale to capture the global consistency relationship. This is consistent with our intuition that a paper with higher importance

and influence in the academic community tends to receive more citations in the subsequent years.

Finally, in time series generation module, the dynamic evolutionary trends embeddings and node importance embeddings calculated above are transformed into the parameters of the time series generation module, using a simple multilayer perceptron (MLP). Following [15, 32], we use a log-normal distribution model to generate the prediction citation counts sequence for each newly published paper.

In summary, our main contributions can be summarized as follows:

- We solve the challenging cold start problem in time series citation prediction. To be specific, it refers to the prediction of citation count for newly published articles without historical citation count values.
- We propose a novel framework named DGNI for citation time series prediction, which leverages both the local consistency relationship and global consistency relationship of the heterogeneous academic network. Hence our model can make full use of the dynamic academic network and node importance to predict future citation count of newly published papers.
- We conduct extensive experiments on two large-scale real-world academic network datasets, and the experimental results illustrate that our model outperforms the SOTA models by 11.39% improvement in terms of MAE and 12.90% improvement in terms of RMSE.

The rest of this paper is organized as follows: Section 2 introduces necessary definitions and makes a formal definition of the problem we tackle. Section 3 introduces the motivation and framework of our proposed model DGNI, and further elaborate each component of our model. Section 4 evaluates the performance of DGNI by experiments and analyses. Section 5 reviews the related works of citation count prediction, node importance estimation and heterogeneous graph representation learning. Section 6 makes a conclusion to the entire paper.

2 PRELIMINARIES

In this section, we introduce necessary definitions used in the paper and make a formal definition of the problem we study.

2.1 Definitions

2.1.1 Heterogeneous Academic Network. A heterogeneous academic network is a special kind of heterogeneous information network (HIN), which consists of multiple types of nodes and edges to represent academic networks. It can be defined as a graph $\mathcal{G} = (\mathcal{V}, \mathcal{E})$ with node mapping function $\varphi : \mathcal{V} \rightarrow \mathcal{T}$ and edge mapping function $\phi : \mathcal{E} \rightarrow \mathcal{R}$ where $|\mathcal{T}| + |\mathcal{R}| > 2$.

The types of nodes in heterogeneous academic network include paper, venue, field and author. It is widely noticed that the papers are the central nodes and other nodes are neighbors. The types of edge include publish (paper-venue), write (author-paper), contain (paper-field) and cite (paper-paper).

2.1.2 Metapath and metapath-based subgraph. Metapath is defined as a path with the following form: $A_1 \xrightarrow{R_1} A_2 \xrightarrow{R_2} \dots \xrightarrow{R_{l-1}} A_l$ (abbreviated as $A_1 A_2 \dots A_l$), where $A_i \in \mathcal{V}, R_i \in \mathcal{E}$. The metapath

describes a composite relation between node types A_1 and A_l , which expresses specific semantics.

Given a metapath ϕ_p of a heterogeneous graph \mathcal{G} , the metapath-based subgraph of graph \mathcal{G} is defined as a graph composed of all neighbor pairs based on metapath ϕ_p .

2.1.3 Dynamic Heterogeneous Network. A dynamic heterogeneous network is a sequence of heterogeneous academic network from 1 to T year: $\langle \mathcal{G}_t \rangle_{t=1}^T = \{\mathcal{G}^1, \mathcal{G}^2, \dots, \mathcal{G}^T\}$, where $\mathcal{G}^t = (\mathcal{V}^t, \mathcal{E}^t)$ ($1 \leq t \leq T$) is the heterogeneous academic network in t th year.

2.1.4 Node importance. A node importance $s \in \mathbb{R}^+$ is a non-negative real number representing the significance or the popularity of an entity in a knowledge graph. For instance, the gross of the movie or the voting number for the movie on the website can be regarded as the node importance in movie knowledge graphs. The specific importance value of a node is collected from the real scenarios and obtained after the log transformation.

2.2 Problem Formalization

Given a dynamic heterogeneous network $\langle \mathcal{G}_t \rangle_{t=1}^T$ and a target paper p , paper citation time series prediction aims to learn a function $f: (\langle \mathcal{G}_t \rangle_{t=1}^T, p) \rightarrow \{c_p^{T+1}, c_p^{T+2}, \dots, c_p^{T+L}\}$ that predicts the citation time series of the target paper p in the following L years after the publication year T .

3 METHODOLOGY

In this section, we introduce our proposed model DGNI for citation count prediction. First, we describe our motivation to design the architecture of DGNI, then we present the framework of DGNI, as shown in Fig. 1. Next, we elaborate the details of three components in DGNI: dynamic heterogeneous network embedding module, node importance embedding module, and time series generation module.

3.1 Framework of DGNI

The key idea of our model is to learn a continuously evolving vector representation of each node from the snapshots of the academic network in different periods, so the node representation can reflect the node’s evolutionary trend. Powered by such representations, the dynamics of the academic network can be well captured, making it easier to predict future citations for new papers. As shown in Fig. 1, our proposed method DGNI is composed of three modules: dynamic heterogeneous network embedding module, node importance embedding module, and time series generation module.

According to [15], the impact of paper can be predicted by modeling changes of snapshots of dynamic heterogeneous networks in different periods. In our dynamic heterogeneous network embedding module, instead of simply using RGCN [23] to capture the heterogeneity of the academic network, which keeps distinct non-sharing weights for node types and edge types alone and is insufficient to capture heterogeneous properties, we use the SOTA heterogeneous graph neural network HGT [13] to encode the dynamics and heterogeneity in each year’s heterogeneous network snapshot. The HGT model automatically learns different weights for different types of nodes and relations, and aggregates features accordingly. The node representations learned by HGT serve as dynamic network features before the paper is published.

Since the heterogeneous graph representation learning algorithm is based on the neighbor aggregation mechanism, it takes only the information of a very limited neighborhood for each node. So it can only capture the local evolutionary trend patterns of academic network, unable to capture global trend patterns. As there are huge number of nodes in the academic network interacting with each other, and each node contribute differently, it makes sense to encode the global evolutionary trends to model each node’s importance. A paper node with higher importance tends to receive more citations in the following years and vice versa. In node importance embedding module, we take Personalized PageRank (PPR) [21] into the graph neural network to reflect the larger neighbor information of different types of relations, capturing the global evolutionary trends. To make full use of the heterogeneity in the graph, we devise a semantic-level attention mechanism to learn the importance of different meta-paths and fuse them automatically. After that, we obtain each new paper’s node importance patterns.

In time series generation module, we use attention mechanism to fuse the dynamic node features and node importance patterns learned by above modules to generate final embeddings of each paper. The paper embeddings serve as the parameters of a parametric citation count generator. Based on the work [32], we use a log-normal distribution to encode prior knowledge of citation processes and generate citation count time series in the years immediately following publication. We unfold the details of these three modules in the following subsections.

3.2 Dynamic Heterogeneous Network Embedding Module

A dynamic heterogeneous academic network refers to a sequence of heterogeneous networks from several years before a paper’s publication year T , and it reflects the evolutionary trends of the entire academic network in different periods. For each year’s heterogeneous network, to capture the rich semantic heterogeneity information, we use the Transformer[31]-based heterogeneous graph neural network model HGT [13] to learn the node representations of each node, which treats one type of node as query to calculate the importance of other types of nodes around it. As each year’s network snapshot reflects different parts of the whole dynamic academic network, and should be comparable, we use the same HGT model to encode each static heterogeneous network in each year.

Secondly, in academic dynamic network, unlike other dynamic networks such as social dynamic network, nodes won’t change rapidly. In other words, the characteristics of same nodes in adjacent years tend to be similar. Inspired by this, we introduce a Mean Squared Error (MSE) loss to force embeddings of same nodes in adjacent years to be close to each other, named temporal-aligned loss:

$$\mathcal{L}_{time} = \frac{1}{\mathcal{T} - 1} \sum_{t=1}^{\mathcal{T}-1} \frac{1}{|V_t \cap V_{t+1}|} \sum_{i \in V_t \cap V_{t+1}} \|h_i^t - h_i^{t+1}\|_2^2 \quad (1)$$

where h_i^t denotes node i ’s embedding at year t (year starts from 1), V_t denotes the node set in the heterogeneous network of year t , and \mathcal{T} denotes the number of years (observed heterogeneous networks) before the paper’s publication.

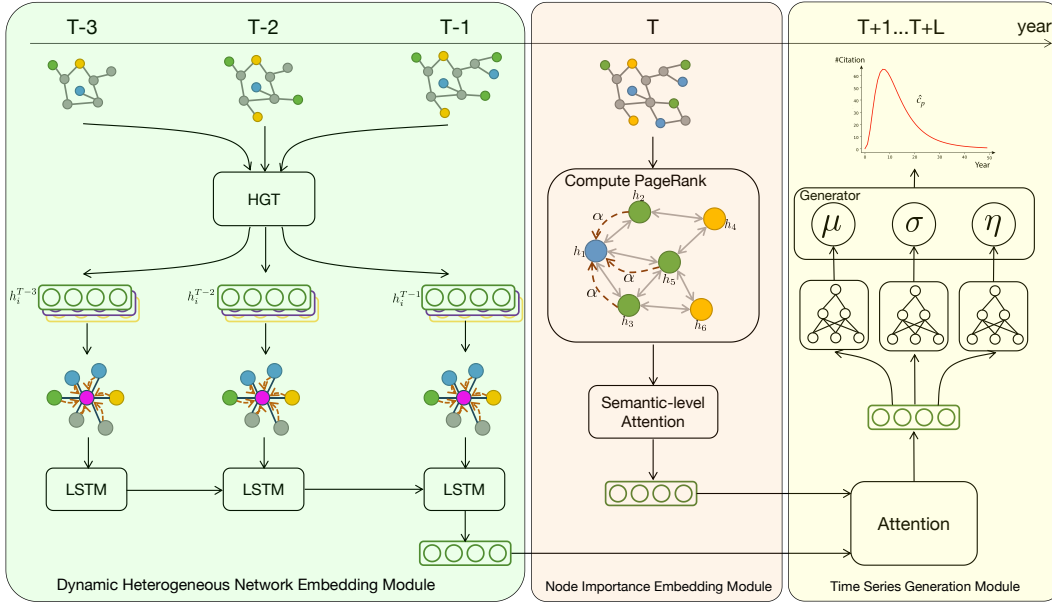


Figure 1: The overall architecture of our proposed model DGNI. To predict the citation count time series of a new paper p published in year T , DGNI first learns the heterogeneous network features from years before T to generate the fake embeddings of p before its publication. Then we compute the node importance of p by our proposed node importance module. After that, the above embeddings are fed into the final time series generation module to generate the cumulative citation of p in $1-L$ years after its publication.

By encoding static heterogeneous networks using HGT model, we obtain the embeddings of each node in each year. However, since a newly published paper has no historical citations, it does not exist in the dynamic academic network. As a solution, we use the metadata nodes (e.g. authors, venues, keywords, etc.) existing in previous years which are linked to the new paper to generate the "fake" embeddings of the paper node in past timestamps. The reason is that a paper's metadata nodes are often with a long history and have a high probability to exist in previous years' academic networks.

To fulfil that, we first explore the snapshot network related to the published year of the paper to find out the linked metadata neighbor nodes, denoted as N_p . Then we look back on each previous year's snapshot network and average the features of these neighbor nodes to get the paper's fake embeddings in each past year. As different types of neighbor nodes may not contribute equally to the impact of the paper, inspired by [15], we apply type-aware trainable weights to preserve the unequal contribution of different kinds of metadata neighbor nodes, as follows:

$$v_p^t = \sum_{r \in \mathbf{R}} \sum_{i \in N_{p,t}^r} W_r \cdot \frac{\mathbf{h}_{i,t}}{|N_{p,t}^r|}, \quad (2)$$

where \mathbf{R} denotes the relation set in the network, $\mathbf{h}_{i,t}$ denotes the feature of node i in year t , $N_{p,t}^r$ denotes the set of neighbor nodes adjacent to the paper p based on relation r , and W_r denotes the learnable weight of relation r shared in all years of network snapshots. After that, the generated fake embeddings of paper p in year t , denoted as v_p^t , can be obtained.

We apply Eq. 2 in every timestamp to obtain a sequence of generated fake embeddings of new paper p as $V_p = \{v_p^t, v_p^{t+1}, \dots, v_p^{T-1}\}$, where t is the first year when paper p 's metadata nodes can be observed. Then, to model the new paper p 's temporal trajectory, we temporally encode the embedding sequence V_p into a single vector h_p^g through the recurrent neural network LSTM [12] as follows:

$$h_p^g = \text{LSTM}(v_p^1, v_p^2, \dots, v_p^T). \quad (3)$$

After that, we obtain h_p^g , the dynamic heterogeneous network feature vector of the paper p before its publication, which reflects new paper's dynamic evolutionary trends. And it will be used in the time series generation module.

3.3 Node Importance Embedding Module

Since the heterogeneous graph learning algorithm is based on a neighbor aggregation mechanism, it takes only the information of a very limited neighborhood of each node to avoid overfitting and over-smoothing. As a result, the heterogeneous graph neural network can only capture the local consistency relationship of the academic network. However, the global consistency relationship is vitally important. For instance, in an academic network, each scholar can be a member of several communities and can be influenced by his neighborhoods with different distances from local consistency relationship to global consistency relationship, so only considering the local relationship tends to be one-sided.

To capture the global consistency relationship, we propose a node importance embedding module to calculate each paper's node importance on the global scale. Intuitively, node importance has a close connection with citation counts. A paper with higher node

importance has higher academic impacts and tends to receive more citations in the following years.

In this work, we take Personalized PageRank (PPR) [21] into the graph neural network to reflect the global consistency relationship of the whole academic network. We define the PPR matrix as:

$$\Pi^{PPR} = \alpha(I_n - (1 - \alpha)D^{-1}A)^{-1}, \quad (4)$$

where α is a teleport probability. The PPR representation of node i refers to the i^{th} row in Π^{PPR} as $\pi(i) := \Pi_{i,:}^{PPR}$. As the academic network used in our experiment is too large to compute on the whole graph, following the work of [6], we use random walk sampling [10] as an efficient and scalable algorithm for computing an approximation of PPR. To guarantee the absolute error lower than ϵ with probability of $1 - \frac{1}{n}$, we need $O(\frac{\log n}{\epsilon^2})$ random walks. Additionally, we also take the instructions in [6] to truncate Π^{PPR} to contain only the top k largest entries for each row $\pi(i)$, denoted as $\Pi_i^{\epsilon,k}$.

However, the above method can only handle the homogeneous graph structure. When it comes to heterogeneous graph, since it contains different types of nodes and links, each node is connected via various types of relations, e.g., meta-paths. Since different meta-paths reflect different aspects of the whole graph, and they take unequal contribution to the final result, we compute PageRank respectively in each metapath-based subgraph. Then, inspired by Graph Attention Network (GATv2) [2], we propose a novel attention mechanism to learn the importance of different meta-paths and fuse multiple semantics revealed by them.

Firstly, we extract metapath-based subgraphs by each meta-path, and compute PageRank matrix Π^{ϕ_p} for each subgraph using Eq. 4. Then by modifying the attention mechanism proposed in GATv2 [2] to

$$e^{\phi_p}(h_i, h_j) = \alpha^{\phi_p} \cdot \text{LeakyReLU}([W_{\phi_p} h_i || W_{\phi_p} h_j || \Pi_i^{\phi_p} || \Pi_j^{\phi_p}]), \quad (5)$$

where α^{ϕ_p} denotes the attention vector of meta-path ϕ_p , h_i denotes the raw feature vector of node i , W_{ϕ_p} denotes the transformation weight matrix of meta-path ϕ_p , aiming at projecting the raw node feature into the meta-path vector space, and $\Pi_i^{\phi_p}$ denotes the PageRank pattern of node i in meta-path ϕ_p . By Eq. 5, we can naturally incorporate the global PageRank patterns into the GAT layer in each subgraph.

Then, we normalize the attention scores from all neighbors $j \in \mathcal{N}_i^{\phi_p}$ within meta-path ϕ_p , and aggregate these features by learned weights:

$$\alpha^{\phi_p}(h_i, h_j) = \frac{\exp(\text{LeakyReLU}(e^{\phi_p}(h_i, h_j)))}{\sum_{k \in \mathcal{N}_i^{\phi_p}} \exp(\text{LeakyReLU}(e^{\phi_p}(h_i, h_k)))}, \quad (6)$$

$$z_i^{\phi_p} = \text{LeakyReLU}\left(\sum_{j \in \mathcal{N}_i^{\phi_p}} \alpha^{\phi_p}(h_i, h_j) \cdot W_{\phi_p} \cdot h_j\right). \quad (7)$$

After that, we get node embedding $z_i^{\phi_p}$ for each meta-path ϕ_p , incorporated with both PageRank patterns and metapath-level representations. Next, in order to fuse these metapath-level representations to get the final node embeddings, we use an attention mechanism similarly, but at the level of meta-path.

Specifically, we do a normalization on each meta-path by averaging the node embeddings $z_i^{\phi_p}$ learned before. Then we use an attention vector q to transform these embeddings into the importance of specific meta-path, denoted as w^{ϕ_p} :

$$w^{\phi_p} = q^T \cdot \frac{1}{|\mathcal{V}|} \sum_{i \in \mathcal{V}} \text{ReLU}(W \cdot z_i^{\phi_p}). \quad (8)$$

To obtain the weight of meta-path ϕ_i , we normalize the above importance of all meta-paths by softmax function:

$$\alpha^{\phi_p} = \frac{\exp(w^{\phi_p})}{\sum_{p=1}^P \exp(w^{\phi_p})}, \quad (9)$$

where P denotes the number of all meta-paths. Since the learned α^{ϕ_p} can be interpreted as the contribution of each meta-path ϕ_p , we take it as coefficient to fuse semantic-specific embeddings belonging to different meta-paths, to obtain the final embedding as follows:

$$h_i^c = \sum_{p=1}^P \alpha_{\phi_p} \cdot h_i^{\phi_p}. \quad (10)$$

The final node embedding h_i^c reflects the node importance of paper i , capturing the global consistency relationship. It will be used for the time series generation module, together with node features reflecting new paper's dynamic evolutionary trends described in Sec. 3.2.

3.4 Time Series Generation Module

Through above two modules, for a target paper p , we can generate its dynamic evolutionary trend patterns and node importance patterns in the whole academic network. To combine node dynamic trend patterns and node importance patterns, we use an attention mechanism:

$$h_p = \alpha_g h_p^g + \alpha_c h_p^c \quad (11)$$

$$\alpha_g = \frac{\exp(\lambda^T h_p^g)}{\exp(\lambda^T h_p^g) + \exp(\lambda^T h_p^c)}, \quad (12)$$

$$\alpha_c = \frac{\exp(\lambda^T h_p^c)}{\exp(\lambda^T h_p^g) + \exp(\lambda^T h_p^c)}, \quad (13)$$

where h_p^g refers to the dynamic trend vector of paper p calculated in Sec. 3.2, h_p^c refers to the node importance vector of paper p calculated in Sec. 3.3, and λ denotes the trainable attention weight.

As described in [15, 32], a new paper's influence will reach the peak within a few years after publication, and will gradually decrease on account of novelty fading and continuous appearance of new ideas and new research topics attracting researchers' attention, known as aging effect. Therefore, we model the citation trajectory of a paper as a log-normal probability along time t :

$$P_p(t) = \frac{1}{\sqrt{2\pi}\sigma_p t} \exp\left[-\frac{(\ln t - \mu_p)^2}{2\sigma_p^2}\right], \quad (14)$$

where μ_p denotes the mean of the normal distribution, which describes the time required for an article to reach the peak of citation trajectory. σ_p denotes the variance of the normal distribution, which describes the decay rate of paper p 's citation decrement.

As discussed in [15], the "fitness" makes significant contributions to a paper's citations, so another parameter η_p is used to model it. Integrated across η_p , the cumulative number of citations of a paper can be generated by the cumulative distribution function:

$$C_p^t = \alpha \left[\exp(\eta_p * \Phi(\frac{\ln t - \mu_p}{\sigma_p})) - 1 \right], \quad (15)$$

where η_p is the parameter which weights the citation count to model the difference between papers. α is a scalar that adjusts the weight of the result, which is a hyper-parameter that will be fixed during the model training process. $\Phi(x)$ is defined as:

$$\Phi(x) = (2\pi)^{-1/2} \int_{-\infty}^x e^{-y^2/2} dy. \quad (16)$$

To get three parameters μ_p , σ_p and η_p to generate each new paper's citation time series, we use three Multilayer Perceptron models (MLP) to transform the final node embedding h_p of the target paper p to generate these parameters:

$$\mu_p = \text{MLP}_1(h_p), \quad (17)$$

$$\sigma_p = \text{MLP}_2(h_p), \quad (18)$$

$$\eta_p = \text{MLP}_3(h_p). \quad (19)$$

After that, we can obtain the cumulative citation of target paper p in $1-L$ years after publication, denoted as a sequence $C_p = \{C_p^T, C_p^{T+1}, \dots, C_p^{T+L}\}$.

3.5 Loss Function

The cumulative citation sequence of target paper p in $1-L$ years after publication is calculated by Eq. 15, denoted as $\{C_p^T, C_p^{T+1}, \dots, C_p^{T+L}\}$. The loss function is composed of two parts: prediction loss and temporal-aligned loss.

The temporal-aligned loss is discussed in Sec. 3.2. As nodes in academic dynamic network won't change rapidly and the characteristics of same nodes in adjacent years tend to be similar, the temporal-aligned loss aims to force embeddings of same nodes in adjacent years to be close to each other, as described in Eq. 1.

For prediction loss, we adopt the Mean Square Error (MSE) to compare the predicted time series with the ground-truth as prediction loss:

$$\mathcal{L}_{pred} = \frac{1}{P} \sum_{p=1}^P \frac{1}{L} \sum_{t=T+1}^{T+L} (C_p^t - \hat{c}_p^t)^2, \quad (20)$$

where C_p^t denotes the ground-truth citation count of paper p in the t^{th} year, T denotes the number of prediction years, and P denotes the total number of papers for prediction. Since the citation counts of different papers vary widely, we conduct log transformation on ground-truth citation counts to smooth rapid changes, i.e. $\hat{c}_p^t = \log(C_p^t + 1)$, and make predictions on the logged version of true values.

The overall model will be optimized by prediction loss (Eq. 20) and temporal-aligned loss (Eq. 1) at the same time. The total loss is defined as follow:

$$\mathcal{L} = \mathcal{L}_{pred} + \beta \mathcal{L}_{time} \quad (21)$$

where β is the hyper-parameter to adjust the proportion of temporal-aligned loss in total loss.

Table 1: The Statistics of Datasets.

Dataset	#node				#edge
	#paper	#author	#keyword	#venue	
APS	311,533	161,051	41,126	9	3,250,651
AMiner	1,026,795	831,151	34,833	3,673	10,366,576

4 EXPERIMENTS

In this section, we evaluate the performance of our proposed model DGNI by experiments on two real-world large-scale datasets. We describe our experimental settings and then show numerical comparison results with other citation count prediction baselines. To help readers understand how DGNI works, we breakdown the model in ablation studies and conduct visualization analyses to prove DGNI's efficacy.

4.1 Experimental Setup

4.1.1 Datasets. We conduct experiments on two real-world datasets *APS* and *AMiner*. The statistics of nodes and edges about the two datasets are shown in Table 1. Below we take a brief introduction:

APS¹ (*American Physical Society*) is a dataset that covers publications in the journal of the American Physical Society, including three node types: paper, author, venue. To generate keyword nodes, we extract keywords from the title of papers following the pre-processing procedure proposed by [24]. In the experiment, We use papers from 2003 to 2008 to build the training set to train the model, papers in 2009 to build the validation set, and papers in 2010 to build the testing set.

AMiner² is a dataset that covers publications in computer science venues [28], including four node types: paper, author, venue and keyword in its V11 version. The training, validation and testing datasets are of the same configuration as APS.

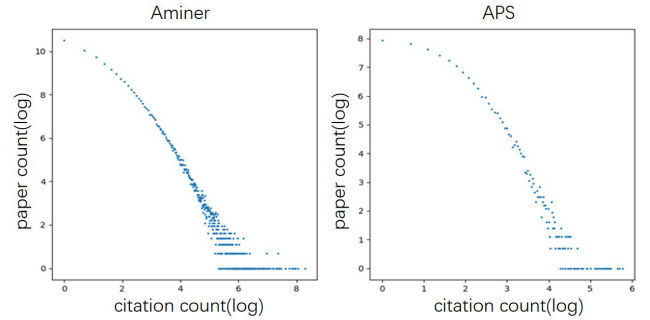


Figure 2: Distribution of cumulative citation counts within five years after publication.

The distribution of cumulative citation counts of papers in two datasets is shown in Fig. 2. We can note that the cumulative citations of papers subject to the long-tailed distribution: most papers are rarely cited after publication, and only a fraction of them can receive considerable citations.

¹<https://journals.aps.org/datasets>

²<https://aminer.org/citation>

4.1.2 Baselines. Since the "cold start" citation count time series prediction task is a novel problem, there is only one work (i.e. HINTS [15]) to compare with. Besides, we consider 4 other citation count time series prediction methods for comparison, which are briefly described below.

- **Gradient Boosting Machine (GBM):** A gradient boosting model used to model scientific features and predict citation time series. Following [15], we extract scientific features that are available in our problem setting or data, to predict citation time series with XGBoost [5].
- **DeepCas [18]:** This model conducts random walk across an information cascade graph to predict popularity. In our experiments, the ego network of a new paper in the publication year is used as the initial cascade graph.
- **HINTS [15]:** A state-of-the-art model for "cold start" citation count time series prediction. This model uses R-GCN [23] to encode dynamic heterogeneous graph and a log-normal distribution to generate the citation time series.

4.1.3 Evaluation Metrics. Following [3, 15, 18], we use the Mean Absolute Error (MAE) and Root Mean Squared Error (RMSE) to evaluate the accuracy of predictions, which are common choices for regression tasks. MAE and RMSE are defined as follows:

$$\text{MAE}(c^t, c^{\hat{t}}) = \frac{1}{P} \sum_{p=1}^P |c_p^{\hat{t}} - c_p^t|, \quad (22)$$

$$\text{RMSE}(c^t, c^{\hat{t}}) = \sqrt{\frac{1}{P} \sum_{p=1}^P (c_p^{\hat{t}} - c_p^t)^2}, \quad (23)$$

where c_p^t and $c_p^{\hat{t}}$ denote the ground-truth and the prediction of citation counts of paper p in the t^{th} year after publication, and P denotes the total number of papers for prediction.

4.1.4 Implementation Details. We implement DGNI using PyTorch 1.10.0. For the dynamic heterogeneous network embedding module, the number of historical dynamic heterogeneous networks we use to model is 3, the dimension of output features of HGT is 32, the layer count and attention heads are set to 2 and 4. The node features in dynamic heterogeneous network are randomly initialized using Xavier initialization. The hidden dimension and the layer count of GRU are 32 and 3 respectively.

For node importance embedding module, the number of attention heads and attention layers are 4 and 2 respectively, the hidden dimension is set to 32. The meta-paths we use are based on paper nodes: PAP (Paper-Author-Paper), PVP (Paper-Venue-Paper), PKP (Paper-Keyword-Paper). For time series generation module, the hidden dimensions of three fully-connected layers are all set to 20. The weight of temporal-aligned loss β is set to 0.5.

We set learning rate to 0.001 and model parameter optimizer as Adam. We set batch size to 3000 for both APS and AMiner datasets. All the models predict the citation series of the paper in the first 5 years after publication, and the averages are used as the prediction result. We run each experiment 10 times following the same configuration with different random seeds and take the average of all results as final result.

4.2 Numerical Comparison Results

4.2.1 Comparison with Baselines. The prediction results are shown in Table 2. We can summarize that DGNI achieves significant improvement on all the baselines on two datasets in terms of both MAE and RMSE. Since HINTS is the baseline with best performance as we know, for our proposed model DGNI, in terms of MAE, DGNI outperforms HINTS by 10.93% on APS and 11.39% on AMiner. And for RMSE, DGNI outperforms HINTS by 12.90% on APS and 11.31% on AMiner. The results prove the effectiveness and efficacy of our DGNI.

Compared with HINTS, our use of HGT as the dynamic heterogeneous network encoder has stronger feature expression ability than the use of simple R-GCN. Besides, our proposed node importance embedding module can capture the node importance of different papers and pay more attention on papers with higher reputation, thus boosting the prediction of citation count. The more detailed analyses of these two components are elaborated in Ablation Study in Sec. 4.2.2.

Additionally, it can be found from annual prediction results that DGNI can achieve the best results also in the annual results and can achieve better performance in early citation prediction than long-term citation prediction.

4.2.2 Ablation Study on DGNI Components. To present more detailed analyses of these components in our proposed DGNI model and find out why DGNI works, we compare DGNI with two variants on APS and AMiner datasets to evaluate the effectiveness of the modules of our model. The two variants are described as follows:

- **DGNI-graph:** A variant of our framework, which removes the node importance embedding module and only uses dynamic heterogeneous network for prediction.
- **DGNI-inp:** A variant of our framework, which removes the dynamic heterogeneous network embedding module and only uses node importance embedding module for prediction.

The experimental results are shown in Table 3. From the results, we can come to the following conclusions: (1) The whole DGNI can achieve almost best results than all variants on APS and AMiner datasets. It verifies the effectiveness of all our proposed modules in DGNI. (2) The node importance embedding module has a great impact on the results, because after removing the module, the model accuracy on both datasets gets lower. The reason might be that the node importance can reflect the global consistency relationship of the network and can guide the allocation of focus on different papers. Furthermore, in terms of AMiner dataset, the variant of DGNI-inp has the best performance, the reason might be that in AMiner dataset the global consistency relationship is more important. (3) The use of dynamic heterogeneous network plays a key role in the citation prediction. As after removing the dynamic heterogeneous network embedding module, the model performance on both datasets has a huge decrease. The reason might be that heterogeneous graph neural model can capture the rich structure information in the network. And structure information is the core of heterogeneous academic network and vital for citation prediction.

Table 2: Comparison Results of Different Methods over Two Datasets.

Dataset	Model	MAE						RMSE					
		year1	year2	year3	year4	year5	overall	year1	year2	year3	year4	year5	overall
APS	GBM	0.898	0.885	0.900	0.921	1.041	0.934	1.098	1.088	1.105	1.124	1.274	1.139
	DeepCas	0.931	0.923	0.885	0.855	0.832	0.904	1.125	1.139	1.126	1.103	1.062	1.104
	HINTS	0.769	0.809	0.825	0.831	0.828	0.805	0.936	0.994	1.019	1.032	1.035	1.023
	DGNI	0.608	0.689	0.734	0.766	0.788	0.717	0.748	0.850	0.909	0.949	0.978	0.891
AMiner	GBM	0.584	0.920	0.989	1.310	1.260	1.018	0.691	1.031	1.224	1.535	1.621	1.224
	DeepCas	0.948	1.052	1.008	0.898	0.968	0.981	1.054	1.260	1.302	1.245	1.265	1.258
	HINTS	0.610	0.710	0.751	0.775	0.788	0.764	0.769	0.905	0.966	1.001	1.024	0.991
	DGNI	0.491	0.629	0.704	0.759	0.803	0.677	0.606	0.782	0.898	0.986	1.003	0.879

Table 3: The Ablation Results of DGNI on Two Datasets.

Dataset	Metric	DGNI-graph	DGNI-inp	DGNI
APS	MAE	0.728	0.805	0.717
	RMSE	0.909	0.987	0.891
AMiner	MAE	0.699	0.684	0.677
	RMSE	0.914	0.866	0.879

4.3 Visualization Analysis

4.3.1 Feature Dimension Reduction Analysis. In order to verify that DGNI can learn expressive embeddings for each paper, we use T-SNE [30] to project final embeddings h_p on AMiner dataset into a two-dimensional space, as shown in Fig. 3.

The figure represent the log-scale cumulative citation count in 5 years after publication. The blue points indicate low-cited papers, while red points indicate high-cited papers. It can be seen that the citation counts from the red points (left-bottom) to the blue points (right-top) decrease gradually, so DGNI can model papers with different citation counts effectively. But there are still many lowly cited papers and high-cited papers distribute together. The reason is that the specific number of citations of the paper are more related to the quality of the paper itself. Therefore, DGNI is more discriminative on the macro scale.

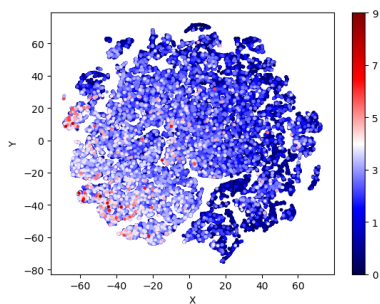


Figure 3: The T-SNE projection result of final embeddings on AMiner dataset. Each point represent the log-scale cumulative citation count.

4.3.2 Normal Distribution Parameter Analysis. Then, we visualize the relationship between the parameters of the normal distribution and the cumulative number of citations, the result is shown as Fig. 4. The cumulative citation counts gradually increase from the left to the right. At the same time, the more the citation counts, the larger the parameter of the normal distribution η_p . In addition, it can be seen that the high-cited papers usually have smaller parameters σ , which indicates that high-cited papers have a larger weight and higher growth rate. The conclusion is in line with the physical meaning of the parameters we proposed in Sec. 3.4.

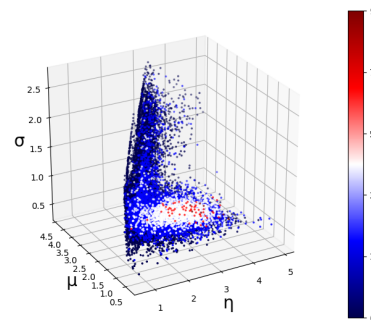


Figure 4: The visualization of parameters. Each point represent the log-scale cumulative citation count.

4.3.3 Prediction Error Analysis. In order to analyze the prediction error between papers of different citation counts, we experiment on AMiner dataset and compare the gap between actual value and predicted value. Specifically, we evenly divide the dataset into three parts by papers' citation numbers: low citation number interval (0%-33%), medium citation number interval (33%-66%) and high citation number interval (66%-100%). For visualization analysis, we plot the average of prediction result and actual result in the following 5 years after publication, on the line chart.

As shown in Fig. 5, the DGNI model has a better prediction accuracy in the medium citation count range, but a higher prediction error in the low citation count range and high citation count range. The reason might be the long-tailed distribution of paper citations. Exactly, most papers included in the real-world dataset can only receive 0 or 1 citations in the following 5 years after publication,

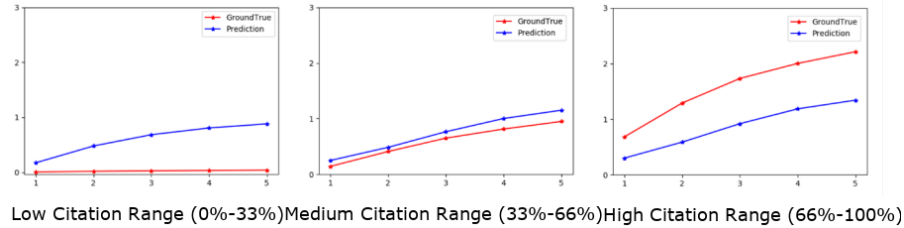


Figure 5: The predicted citation counts compared with ground-truth of papers with different citation range.

while only a small part of papers can receive higher citation numbers. However, the low citation result is not necessarily depended on authors, venues or fields, since a prominent scholar may also publish papers without citations, and papers in top-ranked journals or conferences often receive no citations either. As a result, on the one hand, the independence between papers and their related meta-data brings challenges to the prediction of papers with low citation counts. On the other hand, since papers in high citation count range only constitute a small part of the whole dataset, leading to the poor prediction performance in such a scenario.

5 RELATED WORK

This section reviews three lines of related work: citation time series prediction, heterogeneous graph representation learning, and node importance estimation.

5.1 Citation Time Series Prediction

Citation count prediction includes two categories: using early citation after publication to predict and using information before publication to predict. Parametric approaches uses early citations to model citation trends as a parametric pattern [17, 20, 25, 36]. Some researchers use machine learning methods to model the early citations and citation graph after publication, e.g. Abrishami and Aliakbary [1] used the schema of encoder-decoder to convert the early citations into future citation trends; Li et al. [18] modeled the early citations of the paper as an information cascade network. There are similar methods like [27, 42]. However, this kind of method relies on the early citation within 1-3 years after publication, so it can not deal with the cold start problem.

Some recent works focus on predicting future citations of new papers [7, 39]. Li et al. [19] used peer review text to predict future citations. Xu et al. [37] used heterogeneous academic network to predict the citations of the paper in ten years. The work [15] is the first work to generate citation time series for newly-published paper without any leading citation values, and they use dynamic GNN to model the dynamic academic networks before the publication of papers. Despite most of the citation prediction of newly-published papers are proposed, they are not using information effectively. In this paper, we propose a new framework, which uses academic networks and node importance to predict citation of newly-published papers.

5.2 Node Importance Estimation

Early node importance estimation methods used the degree of nodes in the networks to measure the importance of nodes [11, 29]. PageRank [21] is a classic algorithm based on random walk model to propagate the importance of each node to another node with a

certain probability. By traversing the entire graph, the importance of each node can be quickly calculated. With the rapid development of deep learning in recent years, some methods based on graph neural networks are proposed. Park et al. [22] estimated node importance by using graph attention mechanism, and fusing information of neighbor nodes. Huang et al. [14] utilized a relational graph transformer to learn semantic features, and node2vec [9] to learn structure features. Then combine these features to get the final node importance values. However, these methods can only predict the node importance at a single time rather than time series. In addition, they cannot model the evolution of dynamic network.

5.3 Heterogeneous Graph Representation Learning

Recent years have witnessed the emerging success of graph neural networks (GNNs) for modeling graph structured data [35]. While most GNNs only work for homogeneous graphs, to represent heterogeneous structures and capture the dynamics of network time series, which are more associated with real-world scenarios, researchers have further developed heterogeneous GNNs [33] and dynamic GNNs [16]. For example, Schlichtkrull et al. [23] propose RGCN using multiple weight matrices to project the node embeddings into different relation spaces to capture the heterogeneity of the graph. Wang et al. [34] propose HAN using a hierarchical attention mechanism to capture both node and semantic importance. Hu et al. [13] propose HGT treating one type of node as query to calculate the importance of other types of nodes around it, by multi-head attention mechanism.

6 CONCLUSION

In this paper, we develop a framework named DGNI which adaptively fuses the dynamic evolutionary trends and the node importance in the academic network to predict citation time series with parameter-based generator, tackling the problem of cold start citation count prediction. The contrast experiments and ablation experiments have been conducted on two real-world datasets to demonstrate the superiority of our framework and effectiveness of all the components respectively. For future work, we will consider the interaction between the importance of nodes and more effective generator.

ACKNOWLEDGMENTS

This research was supported by by National Key R&D Program of China under Grant 2019YFA0707204 and the National Natural Science Foundation of China under Grant Nos. 62176014.

REFERENCES

- [1] Ali Abrishami and Sadeq Aliakbary. 2019. Predicting citation counts based on deep neural network learning techniques. *Journal of Informetrics* 13, 2 (2019), 485–499.
- [2] Shaked Brody, Uri Alon, and Eran Yahav. 2022. How Attentive are Graph Attention Networks?. In *International Conference on Learning Representations*. <https://openreview.net/forum?id=F72ximsx7C1>
- [3] Qi Cao, Huawei Shen, Keting Cen, Wentao Robin Ouyang, and Xueqi Cheng. 2017. DeepHawkes: Bridging the Gap between Prediction and Understanding of Information Cascades. *Proceedings of the 2017 ACM on Conference on Information and Knowledge Management* (2017).
- [4] Carlos Castillo, Debora Donato, and A. Gionis. 2007. Estimating Number of Citations Using Author Reputation. In *SPIRE*.
- [5] Tianqi Chen and Carlos Guestrin. 2016. Xgboost: A scalable tree boosting system. In *KDD*. 785–794.
- [6] Julie Choi. 2022. Personalized Pagerank Graph Attention Networks. In *ICASSP 2022-2022 IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)*. IEEE, 3578–3582.
- [7] Yuxiao Dong, Reid A Johnson, and Nitesh V Chawla. 2016. Can scientific impact be predicted? *IEEE Transactions on Big Data* 2, 1 (2016), 18–30.
- [8] James A. Evans and Jacob Reimer. 2009. Open Access and Global Participation in Science. *Science* 323 (2009), 1025 – 1025.
- [9] Aditya Grover and Jure Leskovec. 2016. node2vec: Scalable feature learning for networks. In *Proceedings of the 22nd ACM SIGKDD international conference on Knowledge discovery and data mining*. 855–864.
- [10] William L. Hamilton, Zitao Ying, and Jure Leskovec. 2017. Inductive Representation Learning on Large Graphs. In *NIPS*.
- [11] Taher H Haveliwala. 2002. Topic-sensitive PageRank. In *WWW*.
- [12] Sepp Hochreiter and Jürgen Schmidhuber. 1997. Long Short-Term Memory. *Neural Computation* 9 (1997), 1735–1780.
- [13] Ziniu Hu, Yuxiao Dong, Kuansan Wang, and Yizhou Sun. 2020. Heterogeneous graph transformer. In *WWW*. 2704–2710.
- [14] Han Huang, Leilei Sun, Bowen Du, Chuanren Liu, Weifeng Lv, and Hui Xiong. 2021. Representation Learning on Knowledge Graphs for Node Importance Estimation. In *KDD*. 646–655.
- [15] Song Jiang, Bernard Koch, and Yizhou Sun. 2021. HINTS: Citation Time Series Prediction for New Publications via Dynamic Heterogeneous Information Network Embedding. In *WWW*. 3158–3167.
- [16] Seyed Mehran Kazemi, Rishab Goel, Kshitij Jain, Ivan Kobyzev, Akshay Sethi, Peter Forsyth, Pascal Poupard, and Karsten M. Borgwardt. 2020. Representation Learning for Dynamic Graphs: A Survey. *J. Mach. Learn. Res.* 21 (2020), 70:1–70:73.
- [17] Qing Ke, Emilio Ferrara, Filippo Radicchi, and Alessandro Flammini. 2015. Defining and identifying sleeping beauties in science. In *WWW*. 7426–7431.
- [18] Cheng Li, Jiaqi Ma, Xiaoxiao Guo, and Qiaozhu Mei. 2017. Deepcas: An end-to-end predictor of information cascades. In *WWW*. 577–586.
- [19] Siqing Li, Wayne Xin Zhao, Eddy Jing Yin, and Ji-Rong Wen. 2019. A neural citation count prediction model based on peer review text. In *EMNLP-IJCNLP*. 4914–4924.
- [20] Xin Liu, Junchi Yan, Shuai Xiao, Xiangfeng Wang, Hongyuan Zha, and Stephen Chu. 2017. On predictive patent valuation: Forecasting patent citations and their types. In *AAAI*.
- [21] Lawrence Page, Sergey Brin, Rajeev Motwani, and Terry Winograd. 1999. *The PageRank citation ranking: Bringing order to the web*. Technical Report. Stanford InfoLab.
- [22] Namyoung Park, Andrey Kan, Xin Luna Dong, Tong Zhao, and Christos Faloutsos. 2019. Estimating node importance in knowledge graphs using graph neural networks. In *KDD*. 596–606.
- [23] M. Schlichtkrull, Thomas Kipf, Peter Bloem, Rianne van den Berg, Ivan Titov, and Max Welling. 2018. Modeling Relational Data with Graph Convolutional Networks. *ArXiv abs/1703.06103* (2018).
- [24] Jingbo Shang, Jialu Liu, Meng Jiang, Xiang Ren, Clare R. Voss, and Jiawei Han. 2018. Automated Phrase Mining from Massive Text Corpora. *IEEE Transactions on Knowledge and Data Engineering* 30 (2018), 1825–1837.
- [25] Huawei Shen, Dashun Wang, Chaoming Song, and Albert-László Barabási. 2014. Modeling and predicting popularity dynamics via reinforced poisson processes. In *AAAI*.
- [26] Roberta Sinatra, Dashun Wang, Pierre Deville, Chaoming Song, and A L Barabasi. 2016. Quantifying the evolution of individual scientific impact. *Science* 354 (2016).
- [27] Mayank Singh, Ajay Jaiswal, Priya Shree, Arindam Pal, Animesh Mukherjee, and Pawan Goyal. 2017. Understanding the impact of early citers on long-term scientific impact. In *JCDL*. 1–10.
- [28] Jie Tang, Jing Zhang, Limin Yao, Juanzi Li, Li Zhang, and Zhong Su. 2008. Arnetminer: extraction and mining of academic social networks. In *KDD*. 990–998.
- [29] Hanghang Tong, Christos Faloutsos, and Jia-Yu Pan. 2008. Random walk with restart: fast solutions and applications. *Knowledge and Information Systems* 14, 3 (2008), 327–346.
- [30] Laurens Van der Maaten and Geoffrey Hinton. 2008. Visualizing data using t-SNE. *Journal of machine learning research* 9, 11 (2008).
- [31] Ashish Vaswani, Noam M. Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N. Gomez, Lukasz Kaiser, and Illia Polosukhin. 2017. Attention is All you Need. *ArXiv abs/1706.03762* (2017).
- [32] Dashun Wang, Chaoming Song, and Albert-László Barabási. 2013. Quantifying long-term scientific impact. *Science* 342, 6154 (2013), 127–132.
- [33] Xiao Wang, Deyu Bo, Chuan Shi, Shaohua Fan, Yanfang Ye, and Philip S. Yu. 2020. A Survey on Heterogeneous Graph Embedding: Methods, Techniques, Applications and Sources. *ArXiv abs/2011.14867* (2020).
- [34] Xiao Wang, Houye Ji, Chuan Shi, Bai Wang, Peng Cui, Pinggang Yu, and Yanfang Ye. 2019. Heterogeneous Graph Attention Network. *The World Wide Web Conference* (2019).
- [35] Zonghan Wu, Shirui Pan, Fengwen Chen, Guodong Long, Chengqi Zhang, and Philip S. Yu. 2019. A Comprehensive Survey on Graph Neural Networks. *IEEE Transactions on Neural Networks and Learning Systems* 32 (2019), 4–24.
- [36] Shuai Xiao, Junchi Yan, Changsheng Li, Bo Jin, Xiangfeng Wang, Xiaokang Yang, Stephen M Chu, and Hongyuan Zha. 2016. On Modeling and Predicting Individual Paper Citation Count over Time. In *IJCAI*. 2676–2682.
- [37] Jianguo Xu, Mengjun Li, Jiang Jiang, Bingfeng Ge, and Mengsi Cai. 2019. Early prediction of scientific impact based on multi-bibliographic features and convolutional neural network. *IEEE Access* 7 (2019), 92248–92258.
- [38] Rui Yan, Cong Huang, Jie Tang, Yan Zhang, and Xiaoming Li. 2012. To better stand on the shoulder of giants. In *JCDL '12*.
- [39] Rui Yan, Jie Tang, Xiaobing Liu, Dongdong Shan, and Xiaoming Li. 2011. Citation count prediction: learning to estimate future citations for literature. In *CIKM*. 1247–1252.
- [40] Tian Yu, Guang Yu, Peng-Yu Li, and Liang Wang. 2014. Citation impact prediction for scientific papers using stepwise regression analysis. *Scientometrics* 101 (2014), 1233–1252.
- [41] Sha Yuan, Jie Tang, Yu Zhang, Yifan Wang, and Tong Xiao. 2018. Modeling and predicting citation count via recurrent neural network with long short-term memory. *arXiv preprint arXiv:1811.02129* (2018).
- [42] Qihang Zhao. 2020. Utilizing Citation Network Structure to Predict Citation Counts: A Deep Learning Approach. *arXiv preprint arXiv:2009.02647* (2020).