

Rich-Item Recommendations for Rich-Users: Exploiting Dynamic and Static Side Information

Amar Budhiraja*
Microsoft Research
amar.budhiraja@microsoft.com

Gaurush Hiranandani*
UIUC
gaurush2@illinois.edu

Darshak Chhatbar
Microsoft Research
t-dachha@microsoft.com

Aditya Sinha
Microsoft Research
t-adisi@microsoft.com

Navya Yarrabelly
Microsoft Research
navya.yarrabelly@gmail.com

Ayush Choure
Microsoft Research
aychoure@microsoft.com

Oluwasanmi Koyejo
UIUC
sanmi@illinois.edu

Prateek Jain
Microsoft Research
prajain@microsoft.com

July 28, 2020

Abstract

In this paper, we study the problem of recommendation system where the users/items to be recommended are "rich" data structures with multiple entity types and with multiple sources of side-information in the form of graphs. We provide a general formulation for the problem that captures the complexities of modern real-world recommendations and generalizes many existing formulations. In our formulation, each user/document that requires a recommendation and each item/tag that is to be recommended, both are modeled by a set of static entities and a dynamic component. The relationships between entities are captured by several weighted bipartite graphs. To effectively exploit these complex interactions and learn the recommendation model, we propose MEDRES— a multiple graph-CNN based novel deep-learning architecture. MEDRES uses AL-GCN, a novel graph convolution network block, that harnesses strong representative features from the underlying graphs. Moreover, in order to capture highly heterogeneous engagement of different users with the system and constraints on the number of items to be recommended, we propose a novel ranking metric pAp@k along with a method to optimize the metric directly. We demonstrate effectiveness of our method on two benchmarks:

*Equal Contribution.

a) citation data, b) Flickr data. In addition, we present *two* real-world case studies of our formulation and the MEDRES architecture. We show how our technique can be used to naturally model the message recommendation problem and the teams recommendation problem in the Microsoft Teams (MSTeams) product and demonstrate that it is 5-6% points more accurate than the production-grade models.

1 Introduction

Recommendation systems are a mainstay of most online systems and have been modeled using various approaches such as (inductive) matrix completion, nearest neighbor based predictions, content filtering [16, 10, 2]. However, typical formulations simplify the entire system significantly, e.g., modeling users and items as a single static entity (collaborative filtering). While such simple abstractions provide tractable formulations that can be addressed using rigorous algorithms, they ignore key subtleties of the product environment. In a product system, information is available in multiple and varied forms, all of which may be important for the recommendation module to perform to its full potential. To capture all the information in the system and yet avoid the need for changing the model formulation with the addition of a new form of information, several practical recommender systems adhere to the *content filtering* based approaches. In these approaches, all joint features for (user, item) tuples, extracted from the product, are fed into a classifier [18]. However, designing features for these approaches is a difficult problem, especially for real-world systems with several sources of information. Consequently, such methods suffer from issues such as poor accuracy and a difficult maintenance cycle.

Recent works in the domain of Heterogeneous Information Network (HIN) have proposed to better model such real-life systems with multiple heterogeneous sources of side-information, by constructing heterogeneous graphs and typically use certain metapaths side-information [6] or powerful attention/LSTM based models [31] to extract meaningful information from the heterogeneous graphs. However, these methods still, in general, model a user or item as fixed static entity with a few attributes, so do not capture challenges with many real-world systems where the users/items themselves can be "rich", i.e., are associated with multiple static and dynamic entities (see below).

In this paper, we propose a general formulation for the recommendation problem with rich users and items. In particular, we model the real-world scenarios by considering two types of entities: (a) *static* and (b) *dynamic*. In a given time window, static entities are (mostly) fixed while dynamic entities are generated and modified at a high rate. For example, consider the citation prediction problem [34], where the goal is to predict the set of papers that should be cited by a given paper. Here, the authors and the conferences represent the static entities because their attributes/behavior is fixed in a considerably large time window. On the other hand, the content of the paper (e.g. title/abstract) represents the dynamic entities, as they are generated at a considerably high frequency.

Given this notion of static and dynamic entities, we formalize both the *user* – an entity that seeks out the recommendation – and the *item* – an entity to be recommended – as a collection of static entities as well as a dynamic component. For example, consider an author who is seeking citation recommendations for her paper in preparation, which she plans to submit to a particular conference venue. Here, the user (i.e. the paper) is composed of static entities such as author and conference

venue while its dynamic component is the paper title or abstract. Similarly, the research papers to be recommended are composed of their authors (static), conference venues (static) where they were published, and their title or abstract (dynamic). We call this representation a ‘rich’ representation of users and items. Notice that many items or users may have the same static entities but a different dynamic component.

Standard systems log a multitude of behavioral information about static entities using their past interaction with each other. These interactions, relationships, and engagements can be naturally represented using graphs. One example is the author-conference graph, where an edge denotes the number of papers an author has published in a particular conference. Typical systems would have multiple such graphs among the entities.

Using the above mentioned “rich” representation for users/items along with their behavior captured by multiple graphs, we formulate the recommendation problem as that of finding relevance of an item for a user. The training data consists of labeled (user, item) tuples and multiple bi-partite graphs between static entities. For this formulation, we propose a novel architecture, Multiple Entity based Deep REcommendation System (MEDRES), that combines the multiple entities and graphs associated with users and items via graph embedding techniques to embed each entity in a vector space after accumulating information from multiple graphs for that entity. We then concatenate the entity embeddings with the dynamic component to obtain richer representations of users as well as items. These representations are then fed to a multi-layer perceptron (MLP) to infer how relevant an item is for a user.

Now, in general, MEDRES can be combined with any standard graph embedding technique like Graph Convolution Networks (GCN) [15], Graph Recurrent Neural Networks (GRNN) [33], Node2Vec [9] to process the given multiple graphs and produce representations for static entities. GCN, in particular, has been shown to be one of the most successful techniques in this domain and allows for efficient end-to-end training of the architecture. However, GCN is designed for semi-supervised settings where a node in the graph can correspond to only one labeled data point. So GCN’s goal is to do label propagation over the graph. In our case, a node can correspond to multiple labeled points, and the total volume of labeled data is substantial, so the goal of graph-processing is to extract powerful features by training with the labeled data. To this end, we propose a novel AL-GCN block, that builds upon GCN but adds residual connections, and also allows for learning non-linear functions of edge-weights to extract powerful features. We observe empirically that AL-GCN indeed improves the performance of MEDRES significantly, and due to its generality as a Graph Neural Network block, can be used in various other graph embedding applications.

One crucial advantage of MEDRES is that it is an end-to-end architecture that can be trained for the required metric. Thus, we may optimize MEDRES for standard metrics like AUC, partial-AUC (pAUC), Precision@k (prec@k) [5, 22, 17]. However, existing metrics struggle with two well-known characteristics of real-world recommendation systems: a) the number of items (k) to be recommended per user is bounded, b) a large variance in the number of positives across users (or nodes in the graph). Due to a combination of the above-mentioned issues, existing metrics can fail to capture the “correct” solution and can paint a misleading picture (see Section 3.3). In this paper, we design a new metric –partial-AUC+precision@k (pAp@k)– that combines pAUC and prec@k, and is aware of the above-mentioned characteristics. In fact, the metric is used in *multiple production systems* deployed by MSTeams. In this paper, we formalize the metric and provide an optimization algorithm for the same that we use to train our MEDRES architecture.

Finally, we study the wide-applicability of our formulation and effectiveness of our algorithm

by modeling two problems in MSTeams: message recommendation and Team/Group recommendation. We observe that MEDRES provides impressive gains over the production models for the message recommendation system deployed in MSTeams. We also apply our MEDRES architecture to two publicly available datasets: a) citation dataset, b) Flickr dataset. For both these datasets, we show how challenging recommendation problems can easily be instantiated in our framework, and can provide up to 5 – 6% more accurate recommendations than standard baselines.

In summary, following are the key contributions of the paper:

- We propose a new formulation for the modern recommendation systems via "rich" users and items, that effectively captures dynamically generated content, static entities, and multiple relationships among the entities (Section 2).
- We propose MEDRES, a multiple Graph-CNN based architecture to exploit the above formulation (Section 3). MEDRES' backbone architecture critically uses AL-GCN, a novel graph embedding technique.
- We propose a new metric $pAp@k$, that captures subtleties in several real-world recommendation systems. We also provide an algorithm to optimize the proposed metric (Section 3). The $pAp@k$ metric is used for evaluation in multiple real-world recommendation systems.
- We present empirical validation of our proposed solutions on two benchmark datasets. We also demonstrate our technique on two real-world recommendation systems for MSTeams¹ where our method significantly outperforms strong production-grade models, and is being deployed in the product.

1.1 Related Works

Heterogeneous Information Networks (HINs) is a popular approach to model multiple heterogeneous entities in the system, and are used in several representation learning problems [35, 37, 7]. The core idea of HINs is the emphasis on the graph's heterogeneous entities, but the user-item pairs are considered as static nodes in the graph, thus do not allow for rich users/items with multiple static entities or dynamic component. For example, in the formulation of HINs for citation networks, if citations are required for a new paper, ideally the entire model needs retraining. Our formulation, on the other hand, allows for a more flexible and richer aspect of user-item pairs by allowing items to be dynamic and users and related entities to be static in the system. For instance, in the case of a new paper, we model it as a dynamic item represented by the paper's content vector in our formulation. In another related study[28], the authors also use GCNs to learn embeddings from multiple graphs to compute relevance between user-item pairs, but their approach is also restricted to a static set of items, contrary to the proposed framework which can model items dynamically. Thus, HIN and related works are complementary to our approach, where their architecture can be used as the graph embeddings extraction block in our novel framework.

Graph Representation Learning (GRL) has been widely studied in the literature, and several GRL approaches have been proposed such as Graph Convolutional Networks (GCN) [15], Graph Recurrent Neural Networks (GRNN) [33], node2Vec [9], DeepWalk [24], NGM [3], etc. but all these methods take graph weights as given, which could be problematic because of noise in the real world graphs. GAM [27] is an exception to this, but the method proposed in GAM still uses the same weighing scheme, with augmentation of the propagated weights. AL-GCN on the other hand

¹MSTeams is an enterprise chat+meeting product with millions of active users

learns the graph weights in an end-to-end fashion, thereby extracting task-specific signals from the graph.

Finally, there exists a vast literature on metrics in recommendation system. AUC, pAUC, prec@k, and NDCG@k are the most well-known and popular metrics in this space [22, 11, 30, 19, 4]. However, these metrics do not address the bipartite ranking problem with varied user engagement levels and limitations on the number of recommended items, a problem at the heart of several modern recommendation systems. The pAp@k metric is designed to address the above-mentioned problem and is in use by multiple production systems at MS Teams (see Section 3).

2 Problem Formulation

Given a user/document \mathcal{U} and an item \mathcal{I} , the goal is to find relevance score of \mathcal{I} for \mathcal{U} . Both \mathcal{U} and \mathcal{I} are complex entities with several static sub-entities and a dynamic component. Let $\mathcal{E}^{(1)}, \dots, \mathcal{E}^{(E)}$ be given *entity types* for some $E \in \mathbb{Z}_+$. An entity type can be a set of authors, groups, venues etc. $\mathcal{E}_U^{(i)}$ denotes the i -th entity type for user. For example, if $\mathcal{E}^{(1)}$ denotes the set of groups, then $\mathcal{E}_U^{(1)}$ denotes the group of \mathcal{U} . With these notations in place, \mathcal{U} is now defined as:

$$\mathcal{U} = (\mathcal{E}_U^{(a_1)}, \dots, \mathcal{E}_U^{(a_{k_1})}, \zeta(\mathcal{U})), \quad (1)$$

where each \mathcal{U} is associated with k_1 entity types: $\mathcal{E}^{(a_i)}$, $i \in [k_1]$, and $[k] = \{1, \dots, k\}$ denotes the index set for $k \in \mathbb{Z}_+$. $\zeta(\mathcal{U})$ is the *dynamic* part of \mathcal{U} . For example, for a research paper \mathcal{U} being written for conference C by author A , the entities are conference $\mathcal{E}_U^{conf.} = C$ and author $\mathcal{E}_U^{auth.} = A$. The *dynamic* part is the text of the paper, i.e., $\zeta(\mathcal{U}) \in \mathbb{R}^{D_U}$ can be say D_U -dimensional word2vec embedding [20] of the paper’s abstract and title. Similarly, an item \mathcal{I} is defined as:

$$\mathcal{I} = (\mathcal{E}_I^{(b_1)}, \dots, \mathcal{E}_I^{(b_{k_2})}, \nu(\mathcal{I})), \quad (2)$$

where $\nu(\mathcal{I}) \in \mathbb{R}^{D_I}$ denotes the *dynamic* part of the item with D_I being the dimension of the dynamic component. For instance, a citation item has two static entities (authors and conference venue) and it’s dynamic part is the word-embedding of paper title. The goal is to recommend citations i.e. a set of items for a new research paper \mathcal{U} .

Additionally, a multitude of behavioral information about static entities can be collected using their past interactions. These interactions are expressed via (multiple) bipartite graphs. For example, we may have an author-conference graph where an edge represents the number of times an author has published at a conference. We denote these bipartite graphs between entities $\mathcal{E}^{(a)}$, $\mathcal{E}^{(b)}$ as: $G^{a,b} = (V^{a,b}, A^{a,b})$, where $V^{a,b}$ is the set of nodes and $A^{a,b} \in \mathbb{R}^{|\mathcal{E}^{(a)}| \times |\mathcal{E}^{(b)}|}$ is the adjacency matrix of the graph. The rows and columns in $A^{a,b}$ are associated with entity instances of type $\mathcal{E}^{(a)}$ and $\mathcal{E}^{(b)}$, respectively. For any two entity types $\mathcal{E}^{(a)}$ and $\mathcal{E}^{(b)}$, we may have multiple bipartite graphs denoting different interactions. Thus, let $\mathcal{G}^{a,b} = \{G^{a,b,1}, \dots, G^{a,b,|\mathcal{G}^{a,b}|}\}$ be the set of graphs between entities $\mathcal{E}^{(a)}$ and $\mathcal{E}^{(b)}$. Furthermore, let the set of all graphs be represented by $\mathcal{G} = \{G^1, \dots, G^{|\mathcal{G}|}\}$. In addition to such graphs, we are provided the following dataset for training:

$$\mathcal{D} = \{(\mathcal{U}_1, \mathcal{I}_1, y_1), \dots, (\mathcal{U}_n, \mathcal{I}_n, y_n)\}, \quad (3)$$

where $y_i \in \{0, 1\}$ is the label of the i -th (user, item) pair and denotes the binary relevance of item \mathcal{I}_i for user \mathcal{U}_i . Given all this information, the goal is to find a scoring function that computes

relevance of \mathcal{I} for \mathcal{U} . That is, given \mathcal{D} , entities $\mathcal{E}^{(1)}, \dots, \mathcal{E}^{(E)}$, graphs \mathcal{G} , the goal is to find a scoring function $s(\mathcal{U}, \mathcal{I})$ that works best for a certain metric $\mathcal{M}(\mathcal{D})$. For example, for binary labels y_i 's, the goal would be to find a score s s.t. $s(\mathcal{U}_i, \mathcal{I}_i) > s(\mathcal{U}_j, \mathcal{I}_j)$ if $y_i > y_j$.

Note that most of the existing recommender system formulations can be expressed as a special case of this general formulation. For example, typical collaborative filtering [16] is defined by $\mathcal{U} = (\mathcal{E}^{user}), \mathcal{I} = (\mathcal{E}^{item})$. That is, there are only two entity types: the set of users and set of items. Furthermore, users/items do not have any dynamic content. Berg et al. [2] extended this model by considering the ratings graph $G^{\mathcal{U}, \mathcal{I}}$ between users and items. That is, edge $(\mathcal{U}, \mathcal{I})$ exists in graph if \mathcal{U} has given \mathcal{I} a fixed rating. This model is also a special case of our formulation where we have only two types of entities \mathcal{U}, \mathcal{I} , and the only graphs are the ratings graphs.

Similarly, inductive matrix factorization [10] models the problem as: $\mathcal{U} = (\zeta(\mathcal{U})), \mathcal{I} = (\mathcal{E}^{item})$, where the set of entities contains only the item entity itself. $\zeta(\mathcal{U})$ is the dynamic user/document feature of \mathcal{U} . Recently, several papers [25, 21, 37, 36, 31] have studied the matrix completion problem with graph based side-information and several attributes. Typically, these works model the users/items as: $\mathcal{U} = (\mathcal{E}^{user}), \mathcal{I} = (\mathcal{E}^{item})$ with graphs among users, items, and their attributes. Finally, the standard content filtering approaches [1, 13, 18] do not exploit graph structure and directly models the $(\mathcal{U}, \mathcal{I})$ pair as: $\phi(\mathcal{U}, \mathcal{I})$, i.e., via hand-coded features for \mathcal{U}, \mathcal{I} pair. This shows that most of the existing recommendation problems are a special case of our general framework. Moreover, unlike standard formulations, our general framework leads to a real-world recommendation system that captures much more complex entities and stitches together various sources of information.

3 Method

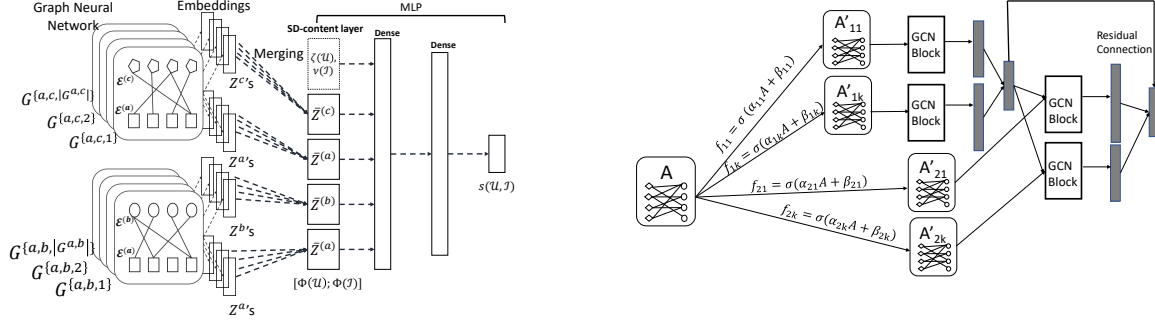
We first discuss the proposed architecture MEDRES, which stitches several types of interactions amongst multiple entities to recommend \mathcal{I} s for \mathcal{U} s followed by AL-GCN which extends GCN by using non-linear adjacency matrix. We then discuss the new metric, pAp@ k , which appropriately measures the effectiveness of a recommendation system in a classification+ranking scenario such as ours followed by an optimization method for the proposed metric. Lastly, we discuss the training of MEDRES for optimizing pAp@ k metric.

3.1 MEDRES Architecture

In this section, we propose MEDRES, an end-to-end architecture, to learn multiple embeddings for all entity instances of every entity type. We then concatenate all these embeddings along with dynamic content to obtain rich $[\mathcal{U}, \mathcal{I}]$ representations and feed them into a multi-layer perceptron (MLP) network to obtain the final relevance score.

Unlike existing methods [18], which use feature extraction from graphs and a separate classification model, we learn all the parameters simultaneously for training our end-to-end architecture while optimizing the proposed metric pAp@ k . Please see Figure 1a for a block representation of the framework.

Let us suppose that we are given a bipartite graph $G^{a,b}$ between two entities $\mathcal{E}^{(a)}$ and $\mathcal{E}^{(b)}$. For simplicity of exposition, let us represent the graph as $G = \{V^{(a)} \cup V^{(b)}, A\}$, where $V^{(a)}$ and $V^{(b)}$ correspond to entity instances of entity types $\mathcal{E}^{(a)}$ and $\mathcal{E}^{(b)}$, respectively, and A is the adjacency



(a) MEDRES: Embeddings are generated using any graph neural network technique. These embeddings are merged together to obtain multiple *Emb-merge* embeddings. Dynamic content is augmented with all the *Emb-merge* embeddings in the *S-D Content* layer, which is then passed to an MLP for binary classification.

(b) AL-GCN: First, the adjacency matrix is transformed by applying k functions at each layer, then these transformed matrices are used by the GCN Block to propagate embeddings, and then results are aggregated at the output of each layer. In the end, residual connections are established between outputs of all layers.

Figure 1

matrix between the nodes defined as:

$$A = \begin{pmatrix} 0 & A^{a,b} \\ (A^{a,b})^T & 0 \end{pmatrix}.$$

Let $F^{(a)} \in \mathbb{R}^{|V^{(a)}| \times p}$ and $F^{(b)} \in \mathbb{R}^{|V^{(b)}| \times q}$ be the feature matrices of entity instances associated with $\mathcal{E}^{(a)}$ and $\mathcal{E}^{(b)}$, respectively². Then, the embeddings for a and b are computed as:

$$\begin{bmatrix} Z^{(a)} \\ Z^{(b)} \end{bmatrix} = f_e(A, F) \quad (4)$$

where f_e is an embedding function (such as Graph Convolution Networks [15]) that takes as input the adjacency matrix A and the feature matrix F of $G^{a,b}$ and outputs $Z^{(a)} \in \mathbb{R}^{|V^{(a)}| \times D_{CN}}$, $Z^{(b)} \in \mathbb{R}^{|V^{(b)}| \times D_{CN}}$, the embeddings of entities $\mathcal{E}^{(a)}$ and $\mathcal{E}^{(b)}$. Here, F is given by: $F = [F^{(a)} \ 0; 0 \ F^{(b)}]$.

Now, as we discussed in Section 2, there can be multiple graphs between entity types $\mathcal{E}^{(a)}$ and $\mathcal{E}^{(b)}$. For example, interactions such as ‘number of papers published’ and ‘number of papers cited’ in a conference by an author may be represented as two different graphs. We merge all the embeddings of an entity type learnt by f_e independently from all the interaction graphs using a fully connected layer. We refer to this layer as the *Emb-merge layer*. Formally, let $\mathcal{G}^{a,b}$ be the set of graphs between entities $\mathcal{E}^{(a)}$ and $\mathcal{E}^{(b)}$. Then, we first concatenate all embeddings i.e. $Z^{(a)}$ ’s (4) of entity instances in $\mathcal{E}^{(a)}$ emerging from the graphs in $\mathcal{G}^{a,b}$ and merge them using a fully connected layer. The same is done for entity instances in $\mathcal{E}^{(b)}$ as described in the following formulation:

$$\bar{Z}^{a,b} = \begin{bmatrix} \bar{Z}^{(a)} \\ \bar{Z}^{(b)} \end{bmatrix} = RELU \left(\begin{bmatrix} Z_1^{(a)} \dots Z_{|G^{a,b}|}^{(a)} \\ Z_1^{(b)} \dots Z_{|G^{a,b}|}^{(b)} \end{bmatrix} \begin{bmatrix} W^{(a)} \\ W^{(b)} \end{bmatrix} + \begin{bmatrix} c^{(a)} \\ c^{(b)} \end{bmatrix} \right), \quad (5)$$

where $\bar{Z}^{a,b} \in \mathbb{R}^{(|V^a|+|V^b|) \times D_{MN}}$ with D_{MN} being the *Emb-merge layer*’s embedding dimensions, $c^{(a)}, c^{(b)} \in \mathbb{R}^{D_{MN}}$ and $W^{(a)}, W^{(b)} \in \mathbb{R}^{(|G^{a,b}| * D_{CN}) \times D_{MN}}$ are the biases and weights of the *Emb-merge*

²In the absence of feature vectors, we can select the one-hot encoding of nodes for $\mathcal{E}^{(a)}$.

layer, respectively. Notice that we always have the flexibility to choose different dimensions D_{MN} for different entity types. $\bar{Z}^{(a)}$ represents an embedding of entity instances in $\mathcal{E}^{(a)}$ with regards to $\mathcal{E}^{(b)}$. For example, at this stage, we have embedding of authors due to all the interactions with conferences.

After utilizing all the static information (graphs) to get entity embeddings, now we augment the dynamic content as features. That is, we represent every user/document \mathcal{U} as:

$$\Phi(\mathcal{U}) = [\bar{Z}_{\mathcal{U}}^{1,2}, \bar{Z}_{\mathcal{U}}^{1,3}, \dots, \bar{Z}_{\mathcal{U}}^{E,E}, \zeta(\mathcal{U})], \quad (6)$$

where $Z_{\mathcal{U}}^{i,j}$ is the *Emb-merge* embedding of \mathcal{U} from all the graphs between entities $\mathcal{E}^{(i)}$ and $\mathcal{E}^{(j)}$. If \mathcal{U} is not of type $\mathcal{E}^{(i)}$ and $\mathcal{E}^{(j)}$, then $Z_{\mathcal{U}}^{i,j} = 0$. We represent each item \mathcal{I} similarly.

We call the concatenation of user-item representations $[\Phi(\mathcal{U}); \Phi(\mathcal{I})]$ as the *S-D-content* embedding, where S-D stands for ‘static’-‘dynamic’. We then apply an MLP layer on the final *S-D-content* as follows:

$$s(\mathcal{U}, \mathcal{I}) = \sigma(\text{RELU}(\text{RELU}([\Phi(\mathcal{U}); \Phi(\mathcal{I})]M_1)M_2)M_3), \quad (7)$$

where M_i ’s are the weight matrices and σ is the sigmoid function (for binary classification problem).

3.2 Adjacency Learnt GCN (AL-GCN)

MEDRES can be combined with any standard graph embedding technique like GCN, GRNN, Node2Vec to process multiple graphs, but now we propose AL-GCN, a new graph embedding technique which performs better than these techniques, when combined with MEDRES. AL-GCN is a novel GCN style graph embedding technique that provides more powerful representations than the standard GCN methods (Section 6). Now, a typical GCN layer is defined as:

$$X^{l+1} = \sigma(AX^lW^l), \quad (8)$$

where X^l is l^{th} layer embedding of GCN, W^l is l^{th} layer weight matrix, A is the Adjacency Matrix, and σ is a non-linear function. Let us consider a multiple-graph setting, where graphs $A^1 \dots A^m$ are given between entities $\mathcal{E}^{(a)}$ and $\mathcal{E}^{(b)}$. In AL-GCN, we apply a function on the adjacency matrix before using it for propagating embeddings over the graph. That is, AL-GCN requires functions $f_p : \mathbb{R} \rightarrow \mathbb{R}$ that are applied entrywise and acts on the aggregation of all graphs, i.e, $\hat{A}_{ij} = [A_{ij}^1, A_{ij}^2 \dots A_{ij}^m]$:

$$\tilde{X}^{l+1} = [\sigma(f_1(\hat{A})\tilde{X}^lW^l), \dots, \sigma(f_k(\hat{A})\tilde{X}^lW^l)], \text{ where } f_p(\hat{A}_{ij}) = \sigma(\alpha_p\hat{A}_{ij} + \beta_p), \quad (9)$$

\tilde{X}^l is the l^{th} layer embedding, W^l is the weight matrix, σ is a non-linear function, $\alpha_p \in \mathbb{R}^m$ and $\beta_p \in \mathbb{R}$. Note that, W^l, α_p, β_p are all learnable parameters in this model.

We further add **residual connections** in MEDRES to make it more robust to overfitting and to allow richer embeddings of the entities to the MLP layers for classification. To add the residual connections, we use the output of the AL-GCN filter at each hop and concatenate the outputs together. After concatenation, the final embeddings looks like:

$$Z = \tilde{X}^1 \oplus \tilde{X}^2 \oplus \dots \oplus \tilde{X}^L, \quad (10)$$

where Z is the final embeddings obtained from AL-GCN, \tilde{X}^l is the output of l^{th} layer AL-GCN as defined in Eq (9) and \oplus makes residual connections between output of these layers. Adding these residual connections allows the network to effectively encode information from multiple hops of a node’s neighborhood. Through our experiments, we observe that adding these residual connections from multiple hops of AL-GCN always outperforms taking embeddings from only the last hop of the architecture. See Fig 1b for a visual illustration of the AL-GCN method.

3.3 The Performance Metric – pAp@k

An important component of any strong recommendation system is the metric that is used to evaluate different methods, and it should capture and reflect the key nuances of the system. Modern systems typically use prec@k, pAUC, AUC, NDCG@k style metrics for this purpose [5, 22, 11, 30, 19, 17].

Now, metrics like NDCG@k do not apply to bipartite ranking style problems where the goal is to rank positives above negatives, and gain/loss for ranking a positive above a negative should not change by their position in the list. For bipartite ranking, pAUC and prec@k style measures are the most popular [5, 11]. However, both of them struggle when two key characteristics common to several recommendation systems are in effect jointly: (a) amount of engagement, i.e., fraction of positives, per user can have high variance, and (b) number of items to be recommended are limited,

For example, let $k = 3$ and let there be 8 points with 2 positives. Let there be two methods that produce rankings r_1 and r_2 of the given points, and let f_1 and f_2 be the true labels corresponding to the ranked order. For example, let $f_1 : \{0, 0, 1, 1, 0, 0, 0, 0\}$ and $f_2 : \{1, 0, 0, 1, 0, 0, 0, 0\}$. The value for prec@k is 0.33 for both f_1 and f_2 , but it is clear that f_2 is a better ordering as a positive is at the top of the list in that case. Now, for a similar example with $k = 3$ and 3 positives, let certain methods produce ranked orders f_3 and f_4 where $f_3 = \{1, 1, 0, 0, 0, 0, 0, 1\}$ and $f_4 = \{0, 1, 1, 1, 0, 0, 0, 0\}$. Here for $k = 3$, f_3 is clearly better but pAUC for threshold $3/6$ is 0.66 for both the rankings.

Naturally, the issue is that in the first case, a classification metric like prec@k does not care for the ranking, while in the second case pAUC is not much aware of top- k recommendations.

Next, we present pAp@k, a novel metric which attempts to mitigate the above mentioned concerns with the existing methods. We first define pAp@k assuming only one user for which we have a given set of labeled points $(x_i, y_i), \dots, (x_n, y_n)$, where $x_i \in \mathcal{X}$ and $y_i \in \{0, 1\}$ and a scoring function $s : \mathcal{X} \rightarrow \mathbf{R}$. Let n_+ be the number of positives (i.e. $y_i = 1$) in the data and n_- be the set of negatives ($y_i = 0$). Let S^+ be the set of top- β positives ordered by scoring function s , where $\beta = \min(n_+, k)$ and k is the number of items to be recommended. Similarly, let S^- be the set of top- k negatives ordered by s . Then,

$$\text{pAp@k}(s) = \frac{1}{\beta k} \sum_{x_i \in S^+} \sum_{x_j \in S^-} \mathbb{1}[s(x_i) \geq s(x_j)], \quad (11)$$

where $\mathbb{1}$ is the indicator function. Note that the metric considers pairwise comparisons between top negative items and top positive items and computes how many times a relevant item has secured higher score than an irrelevant item. Thus, pAp@k rewards presence of every relevant item in top- k scored elements and penalizes high scores of negative items. Also, note that the metric essentially behaves like pAUC for $n_+ \ll k$ and like prec@k for $n_+ \gg k$. Since, n_+ can

be significantly different for various users (and datasets), pAp@k provides a more nuanced and comparative evaluation, due to which it is used as the key-performance-indicator by MSTeams’ production system.

We now define the micro and macro versions of pAp@k to handle multiple users:

$$\text{Micro-pAp@k}(s) = \frac{1}{P} \sum_p \text{pAp@k}_p(s), \quad \text{Macro-pAp@k}(s) = \frac{\sum_p \sum_{x_i \in S_p^+} \sum_{x_j \in S_p^-} \mathbf{1}[s(x_i) \geq s(x_j)]}{P \cdot k \cdot \sum_p \beta_p}, \quad (12)$$

where $\{\mathcal{U}_p\}_{p=1}^P$ are the set of users. We define S_p^+ , S_p^- , β_p for each user p as mentioned above. Finally, pAp@k_p is pAp@k computed for user \mathcal{U}_p . In this paper, we focus on the micro-pAp@k since it is more accommodative of varied engagement level across users and indicative of average of performance per user. For the rest of the paper, unless specified pAp@k refers to micro-pAp@k.

3.4 Training: optimizing pAp@k

MEDRES is trained to approximately optimize regularized pAp@k i.e. the goal is to find parameters Θ of the MEDRES architecture to minimize: $\mathcal{L} = (1 - \text{pAp@k}(\Theta)) + \tau \|\Theta\|_2^2$, where the first term is the pAp@k loss function and τ is the regularizer. Note that the loss function is discontinuous and hard to optimize in general. However, we propose a simple iterative procedure to refine the Θ parameter. In each iteration, we create a pool of datapoints where the pool contains the set of top- β_p positives and top- k negatives for each user p according to the current scoring function. Here, $\beta_p = \min((n_+)_p, k)$ and $(n_+)_p$ is the number of positives for user p . We then optimize the standard cross-entropy loss on the selected points using Adam [14] optimizer. We then update the score function, and again compute the set of top positives and negatives, and iterate till convergence. See Algorithm 1 for a pseudo-code of the method. Note that *Select-top(s,k,label)* function selects top k points of given *label* when sorted by the score s . Also, the subroutine in line 8 in Algorithm 1 is fast since it runs on a few selected points.



Figure 2: The current Production System deploys a three layer architecture with layers– L1: user weaning, L2: message classification and L3: personalized delivery model. Our model is relevant for the $L2$ layer and is evaluated against the production-grade models for the $L2$ layer.

Table 1: Details of MS Teams Message Recommendation Dataset

Statistic	Training Data	Val Data	Test Data
Number of points	3416632	379626	1499083
% of positives	4.40	4.87	4.80
Users	901	898	898
Authors	677	509	593
Channels	476	367	391

Algorithm 1 Training MEDRES to optimize $pAp@k$

Input: TrainingData, Iterations, k , P : no. of users

Output: MEDRES model

```

1: procedure ITERATIVE TRAINING
2:    $S \leftarrow$  TrainingData
3:   model  $\leftarrow$  Model_Init()
4:    $(n_+)_p =$  no. of positives for user  $p$ ,  $1 \leq p \leq P$ 
5:    $\beta_p = \min((n_+)_p, k)$ ,  $1 \leq p \leq P$ 
6:    $i \leftarrow 0$ 
7:   while  $i \leq$  Iterations do
8:     model.train(S)
9:      $scores_p \leftarrow$  model.score(TrainingData,  $p$ ),
            $1 \leq p \leq P$ 
10:     $S \leftarrow \cup_p \{ \text{Select-top}(scores_p, k, '-') \cup$ 
            $\text{Select-top}(scores_p, \beta_p, '+') \}$ 
11:     $i \leftarrow i+1$ 
12:  end while
13: end procedure

```

During our ablation studies (Table 3) we observe that the iterative procedure indeed improves $pAp@k$ metric significantly over the iterations of the proposed method. We leave further theoretical investigation into convergence and consistency guarantees for future work.

4 Experimental Setup

We now discuss the setup we use to evaluate our methods in the next three sections. We refer to the AL-GCN-MEDRES setup with $pAp@k$ optimization as MEDRES in the further sections.

Baselines: For most of our experiments, we use the standard content-filtering baseline (CoF) based on LightGBM [12] and non-linear collaborative filtering Collab as the key baselines. We later introduce other baselines that are scenario-specific. For CoF, we first compute joint features between users and items by using one-hop graph features. That is, for each graph $G^{d,a} = (V^{d,a}, A^{d,a})$ between entity $\mathcal{E}^{(d)}$ of the user \mathcal{U} , and $\mathcal{E}^{(a)}$ of the item \mathcal{I} , we use edge $A^{d,a}(\mathcal{U}, \mathcal{I})$ as the feature. We then collect all the features along with $\nu(\mathcal{I})$ and $\zeta(\mathcal{U})$ (see (1), (2)) to form data

point $\{\phi(\mathcal{U}_i, \mathcal{I}_i), y_i\}$. For Collab baseline, as is standard, we learn features for each entity directly from the given labelled data. That is, intuitively, Collab baseline is the same as MEDRES architecture but rather than using the graphs to regularize/extract features, it learns features from scratch to fit the given labels.

Note that we focus on CoF and Collab as baselines, because they are the most popular and widely applicable techniques used routinely in real-world systems with multiple source of information. Most of the other recommendation system techniques like GCNN-MC[21], HIN [7] do not easily apply to the problems that we study. For example, GCNN-MC [21, 2] in it’s original form handles only user and item graphs and their labels, hence if we have other entities or a dynamic component in users/items, it will not be able to exploit that information faithfully. Similar observations hold for other methods like inductive matrix completion [10] and matrix completion with graph side information [25].

To demonstrate the effectiveness of the AL-GCN block, we also provide a comparison with the MEDRES-GCN architecture where the graph embedding is computed using the state-of-the-art GCN method [15]. We also evaluate the effect of using pAp@k optimization in Section 5. Finally, for the team recommendation problem (Section 6), which is a link-prediction problem with graph side information, we use the Laplacian smoothing method by [25]. Furthermore, in this case, MEDRES-GCN reduces to the standard matrix-completion with GCN formulation [21, 15].

Metrics: For all the benchmarks, we use pAp@k as the primary evaluation metric. In addition, we use prec@k metric for the Teams Recommendation scenario, which is essentially a link-prediction problem, and hence prec@k is a standard metric. For the remaining benchmarks and scenarios, we also compute AUC as it is the default evaluation metric for bipartite ranking style problems.

Hyperparameters: For CoF, we sweep on the following hyper-parameters: learning rate $\{10^{-3}, 10^{-2}\}$, number of trees $\{50, \dots, 300\}$, and number of leaves $\{10, 15, \dots, 30\}$. For MEDRES-GCN, MEDRES and Collab, the cross-validation is done on embedding size $\{2^5, 2^6, 2^7\}$, regularization parameter $\{10^{-4}, 10^{-5}\}$, nodes in fully connected layers $\{2^7, \dots, 2^{10}\}$, batch size $\{1, 3, 5\} \times 10^3$, and learning rate $\{10^{-3}, 10^{-4}\}$. As stated earlier in the paper, we report only Micro-pAp@k (Eq. 12) as it is more relevant metric with a large number of heterogeneous users, and refer to it as pAp@k throughout this section. For the Laplacian smoothing approach [25], we tune on the following hyperparameters: learning rate $\{10^{-3}, 10^{-4}\}$, regularization parameter $\{10^{-3}, 10^{-2}, 10^{-1}\}$ and embedding size $\{2^5, 2^6\}$.

5 Case Study 1: MS Teams Message Recommendation

Microsoft Teams (MSTeams) is a fast-growing enterprise conversation platform. It is composed of teams, users, and channels, where users are a part of a team, and users post on channels to start a conversation. Each user has access to all channels of a team that she/he is a part of, thus privy to a large number of messages posted per day, hence prioritizing messages is a significant challenge and creates an information overload [8]. This problem can be addressed by an accurate message recommendation system that recommends the most relevant messages to users.

The message recommendation problem can be stated as learning a function $s(\mathcal{U}, \mathcal{I})$, which takes as input a user \mathcal{U} and a message \mathcal{I} , and predicts whether the user would be interested in this message or not. That is, the problem is a special case of our formulation in Section 2 with \mathcal{U} being

any user of MSTeams, item \mathcal{I} being any message in a channel, and label y if the \mathcal{U} engages with the message. The term *engagement* can be defined in various ways, here we define a user to have engaged with a message if the user likes the message or replies to the message.

Overall, the system has three entity types: $(\mathcal{E}^U, \mathcal{E}^A, \mathcal{E}^C)$, i.e., the users, the authors of messages, and the channels in which messages are posted. The user \mathcal{U} is defined by only the user entity type. Message \mathcal{I} , on the other hand, is composed of the author of the message, channel of the message as well as the content of the message. That is, $\mathcal{I} = (\mathcal{E}_I^A, \mathcal{E}_I^C, \nu(\mathcal{I}))$ where \mathcal{E}^A is the set of authors, \mathcal{E}^C is the set of channels, and $\nu(\mathcal{I})$ is a vector embedding of the message content. In addition, we are given 10 graphs between the three entities: $\mathcal{E}^U, \mathcal{E}^A, \mathcal{E}^C$.

1. User-Author Graph: This is a graph between users and authors with directional edges such as the number of times a user liked an author’s post, and vice versa, and number of times a user replied to author’s post, etc.
2. User-Channel Graph: This graph models user channel interaction by capturing edges such as how many times a user visited the channel in the past month, how many messages a user posted on the channel, etc.

Production System Architecture: MSTeams has a daily budget of k messages to be recommended per user. Furthermore, messages are recommended as they arrive, so the position of the message is irrelevant, and hence NDCG style metrics are not suitable. Instead, the problem requires bipartite ranking with limited items and highly varied engagement-per-user. Hence, pAp@k metric is a good fit and is used in the production system. Now, the classifier $s(\mathcal{U}, \mathcal{I})$ needs to be invoked for every message \mathcal{I} and the set of users who can access the message. This can be computationally infeasible due to a large volume of messages. The current production system avoids this issue by deploying a three layer architecture with layers– L1: user weaning, L2: message classification and L3: personalized delivery model. Our model is relevant for the $L2$ layer and is evaluated against the production-grade CoF model.

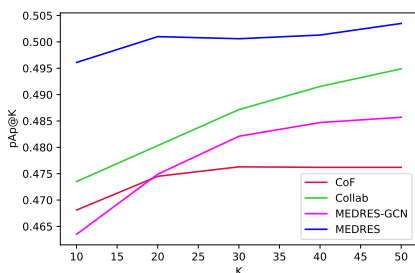


Figure 3: pAp@k accuracy obtained by various methods. MEDRES outperforms baselines by significant margin.

#Graphs = 1			
k	CoF	MEDRES-GCN	MEDRES
10	0.447	0.468	0.478
20	0.448	0.474	0.481
30	0.452	0.478	0.485
40	0.451	0.478	0.488
50	0.452	0.479	0.492

Table 2: Accuracy of MEDRES and CoF baseline with only 2 graphs

Empirical Results: We use a subset of the production dataset for evaluating our MEDRES method against various baselines (Section 4); Table 1 summarizes the dataset. Recall that we define the label of a (user, message) pair to be positive, i.e., $y = 1$, if the user likes/replies to the message. Figure 3 reports pAp@k (12) accuracy for various methods. MEDRES is significantly better than both CoF and Collab baselines across various k values. Furthermore, MEDRES with AL-GCN block is also 3% better than the MEDRES-GCN method. Table 6 shows the AUC measurements,

Model	K=10	K=20	K=30	K=40	K=50
CoF (Prod Model)	0.468	0.475	0.476	0.476	0.476
MEDRES-GCN Baseline	0.464	0.475	0.482	0.485	0.486
MEDRES-GCN + pAp@50 opt.	0.467	0.481	0.488	0.490	0.495
MEDRES 1 Layer(Concatenation)	0.482	0.491	0.496	0.500	0.503
MEDRES 1 Layer(Sum)	0.476	0.490	0.496	0.499	0.502
MEDRES 2 Layers (Concatenation)	0.495	0.499	0.499	0.499	0.501
MEDRES 2 Layers (Sum)	0.465	0.475	0.480	0.485	0.489
MEDRES 2 Layers (Concatenation) + pAp@50 opt.	0.4961	0.5010	0.501	0.501	0.504

Table 3: Ablation results for impact of pAp@k optimization, aggregation method, number of layers, and AL-GCN.

and we observe a similar trend. MEDRES performs the best followed by MEDRES-GCN and then followed by Collab and CoF.

Next, we evaluate MEDRES when only two graphs are available, instead of the 10 graphs used above. That is, 1 User-Author graph and 1 User-Channel graph. Figure 2 presents pAp@k accuracy for this setting. We observe that the performance of MEDRES and MEDRES-GCN is similar to the case when all 10 graphs are available. Thus, MEDRES is able to extract powerful features from 2 graphs while the accuracy of the production baseline (CoF) suffers.

5.1 Ablation Study

We compare the performance of various flavors of MEDRES, including the impact of the AL-GCN block against MEDRES-GCN, the impact of pAp@k optimization, and the impact of the aggregation function (concatenation and sum). Figure 3 summarizes our results. Firstly, we observe that MEDRES has superior performance to MEDRES-GCN for all K values (rows 2, 3 and 6), demonstrating that the AL-GCN block learns powerful features. Secondly, the proposed pAp@k optimization method improves accuracy for both MEDRES-GCN and MEDRES method by about .5%. Thirdly, we observe that the concatenation method always performs better than the sum of functions because the summation of embeddings could lead to a loss of signal which is otherwise available in the concatenation aggregation. Finally, we also compare 1 layer and 2 layers of AL-GCN and show that for concatenation aggregation method, the results improve as we move to 2 layers. This shows that 2 hops of neighborhood aggregation leads to more accurate embeddings. For the sum aggregation method, we observe a different trend i.e. a decrease in performance with 2 layers because for 2 layers, the summation of embeddings is possibly leading to a massive loss of relevant signal as compared to a 1 layer MEDRES.

6 Case Study 2: Team Recommendation (MS Teams)

In this section, we discuss the second real-world use case of our proposed method, which is that of recommending Teams or groups (in MS Teams) that a user can join. This is a popular feature in MS Teams to increase user engagement. The problem is essentially a link prediction problem, but several users have joined limited Teams, so one needs side-information about users. This is

	Train	Val	Test
No. of points	8.7M	2.9M	2.9M
No. of positives	33K	8K	8K
Users	1K	300	300
Teams	9K	9K	9K

Table 4: Details of MS Teams Team Recommendation Dataset

provided in the form of a user-user graph which reflects how the user interacts with other users on certain O365 products.

Formally, the problem is that we have two entities: Users \mathcal{E}^U , and Teams \mathcal{E}^T and the user-item pair is also the same: Users \mathcal{U} and Teams \mathcal{T} . The only graph is between User-User entities, given by: \mathcal{G}^{UU} . So the problem naturally fits in the MEDRES architecture, although it is a significantly simpler problem from a modeling standpoint. Since it is a link prediction problem, we can also use standard techniques like KNN (which is similar to the production model), Inductive Matrix Completion (IMC) where the user features are generated using PCA of the user-user graph (referred to as PCA), the Laplacian smoothed Graph Information method by [25], and finally the GCN based matrix completion technique [15, 21] which reduces to MEDRES-GCN. CoF baseline does not directly apply here as we cannot compute Teams feature from existing graphs (only graph given is \mathcal{G}^{UU}), and Collab baseline does not generalize to new users or users with little Team enrollment, so instead we use the PCA based IMC method mentioned above.

Experiments: We perform our experiments on a subset of the production dataset, as described in Table 4.

For evaluation, we measure both pAp@k and prec@k, which are also used by the production system. Fig. 4 (a) shows the prec@k of the various approaches for different values of k. It can be seen that MEDRES outperforms all the related approaches and performs 3-4% better than the KNN based production baseline for all $k > 3$; and the performance gap increases with an increase in k , which is the practically relevant setting. MEDRES also performs 2% better than MEDRES-GCN and approx. 4% better than PCA for all k values. Fig. 4 (b) shows a similar trend for pAp@k, with the only exception being relatively poor performance of MEDRES-GCN as it might not be extracting strong features for relatively disconnected users on the graph; MEDRES with AL-GCN block is the most accurate method on this dataset.

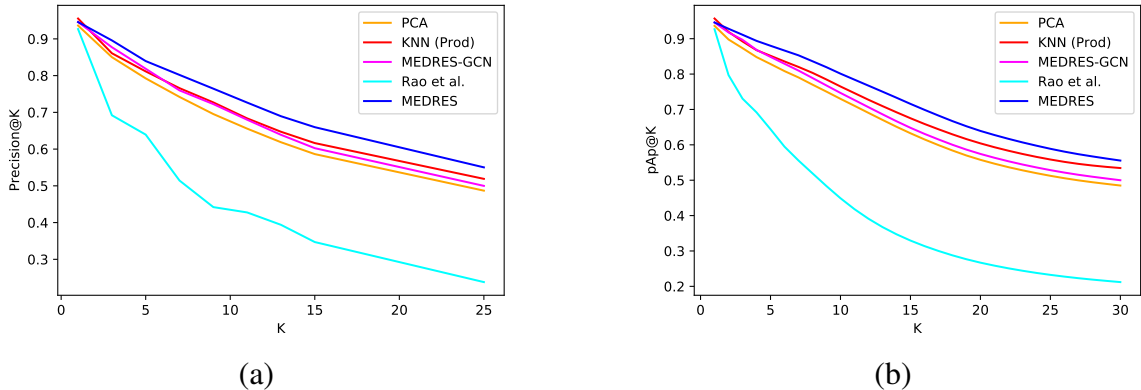


Figure 4: (a), (b): prec@k and pAp@k accuracy for Teams Recommendation problem. Clearly, MEDRES outperforms existing baselines significantly on both the metrics.

	Citation Dataset			Flickr Dataset		
	Train	Val	Test	Train	Val	Test
No. of Points	238k	84k	85k	260k	20k	59k
% of positives	9.99	9.40	9.48	12.87	12.24	12.72
Entities (Count of Entity)	Users(1602);Authors(1147); Conferences(951)			Users(839);Groups(464); Entities(311);Categ(1471)		

Table 5: Details of benchmark datasets

7 Benchmark Dataset Experiments

In this section, we apply our general formulation and MEDRES method on two publicly available datasets of citation networks and Flickr³ and show the wide applicability of our proposed method. We further make the citation dataset and source code available for reproducibility along with the stated hyper-parameters in Section 4.

Citation network problem: Here, we study the problem of recommending relevant papers for a given research paper, i.e., recommending relevant citations for a paper. In particular, the goal is to recommend an existing paper from an author A_j , appearing in conference C_j to a new research paper authored by U_j . That is, the user entity (\mathcal{U}) is a new research paper and item \mathcal{I} , to be recommended, is a citation. There are three entity types: user/new-author \mathcal{E}^U , authors \mathcal{E}^A , and conference \mathcal{E}^C venues. In addition, the paper-title is taken as the dynamic component of the citation. So, $\mathcal{U} = (\mathcal{E}_U^U, \zeta(\mathcal{U}))$ where $\zeta(\mathcal{U})$ is the GloVe embedding of the paper-title [23]. Similarly, $\mathcal{I} = (\mathcal{E}_I^A, \mathcal{E}_I^C, \nu(\mathcal{I}))$ where $\nu(\mathcal{I})$ is the GloVe embedding of the title of the paper-to-be-cited. For citation networks, we consider 5 graphs: 1) User-Published-Conference graph: Number of times a user published in a conference; 2) User-Cited-Conference graph: Number of times a user cited a conference; 3) User-CoAuthor-Author graph: Number of times, a user and author were coauthors in a paper; 4) User-Cited-Author graph: Number of times a user cited an author; and 5) Author-Cited-User graph: Number of times the user (as an author) was cited by another author. We extract data from the Citation-network V1⁴ – a publicly available citation dataset [29]. For creating entity graphs and training data, we used the citation data from 2000-2003, while the test set was created using citation data from 2004-2005. We consider every author and citation pair for a given paper as positive data points. For generating the negative class points, we use the following method: if a user U_i has cited any paper of author A_j , then all the papers of A_j which are not cited by U_i are considered as negative data points. We further sample this data such that a user should have at least 20 data points in the training set, and the author being referred should also have at least 20 data points in the training data. We further remove all rows in the test set where either the user, author, or conference did not occur in the training data. For the paper-titles, we use 50-dimensional GloVe embedding. Table 5 provides key statistics of the training and the test dataset.

Flickr dataset: Here, we consider the task of recommending relevant Flickr groups for users to post their images. Flickr has various groups based on different themes where each group can be associated with several categories as well as tags. The original dataset denotes tags as entities, but we avoid that term to avoid confusion. We want to recommend groups to photos, so we refer to a particular group as item \mathcal{I} and a photo as a user/document \mathcal{U} . \mathcal{U} is represented by the specific user

³www.flickr.com

⁴<https://aminer.org/citation>

Dataset	CoF	Collab	MEDRES GCN	MEDRES
Message Reco.	0.8901	0.8973	0.9026	0.9029
Citation Reco.	0.7645	0.8198	0.8105	0.8156
Group Reco.	0.7535	0.7404	0.7651	0.7874

Table 6: AUC ROC Results

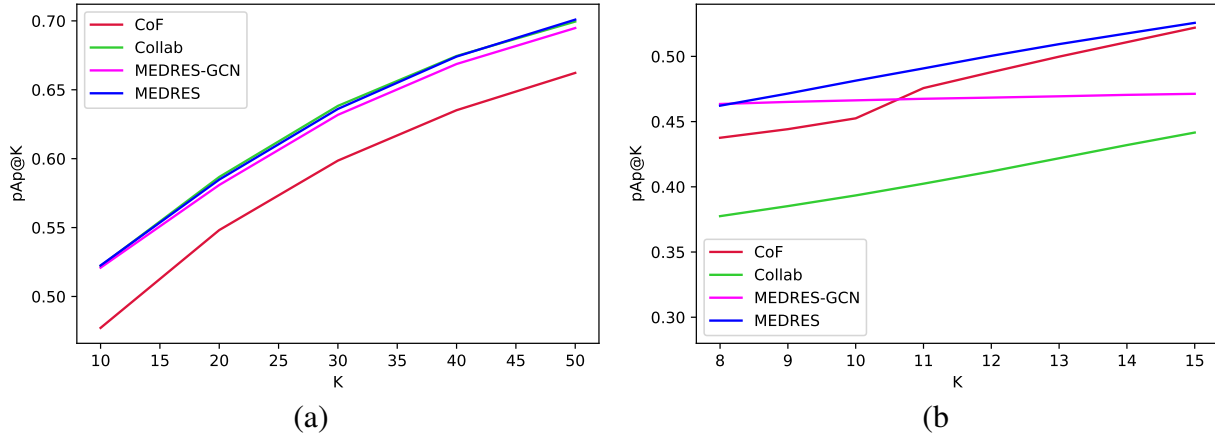


Figure 5: (a): pAp@k accuracy on Citation networks dataset, (b): pAp@k accuracy on Flickr dataset. For Flickr, MEDRES is 4-5% more accurate than baselines.

entity $\mathcal{E}^{(U)}$ and the content of the photograph itself, i.e., $\mathcal{U} = (\mathcal{E}^{(U)}, \zeta(\mathcal{U}))$; we embed images into 4096 dimensional features using VGG19 [26] and then do 100 dimensional PCA reduction of the same to consume as image features. Each item is just the group entity itself, i.e., $\mathcal{I} = (\mathcal{E}^{(G)})$. In addition, we have the following engagement graphs: 1) User-PublishedIn-Group: Number of posts in a group by a user; 2) User-PublishedIn-Category: Number of posts by a user in a category; and 3) User-PublishedIn-Entity: Number of posts by a user in groups containing the entity type under consideration. We use the publicly available Flickr dataset [32]. For each photo posted by a user in a particular group, we create a positive class data point. We further create negative class data points for a user and photo pair from all groups which belong to the same category as the original group of the photo. We then randomly sample the dataset to consider only 10% of the data-points for generating train and test datasets. Table 5 provides statistics for the dataset.

Empirical Results: We applied MEDRES with the same hyperparameters, as discussed in the experimental setting section. Figure 5(c) shows pAp@k obtained by various methods with varying k for the citation recommendation problem. For recommending $k = 10$ articles, MEDRES-GCN is (\geq)4.37% more accurate than the CoF baselines, and similar trend holds for varying values of k , showing that MEDRES, even with vanilla GCN can outperform content filtering model. We further observe an increase of 0.5% increase for pAp@k when using AL-GCN over GCN and Collab. To understand the gain of MEDRES over Collab, we computed results when only 50% data is used in training. While for 100% data the performance of Collab was close to MEDRES, for 50% data, we saw a drop in pAp@k of 3.2% Collab compared to only a 1.5% drop in MEDRES, with pAp@50

for Collab at 0.6658 and while for MEDRES it is 0.6870. Table 6 shows the AUC values where we observe that the Collab baseline has a slightly higher value than MEDRES. This results from the issue that different users have different numbers of data points as stated in Section 3.3. To further understand it, we computed the average of per user AUC (micro-AUC). For MEDRES, its value is 0.7427 and for the Collab baseline, it is 0.7426 which is a reverse trend compared to AUC. This clearly indicates need for more fine-grained metric like pAp@k. Moreover, the pAp@k trend remains clear across datasets.

Finally, we observe a similar trend for the Flickr group problem (Figure 5(d)) where MEDRES performs 2.5% better than MEDRES-GCN, 3% better than CoF and approximately 9% better than Collab at $k = 10$. Thus demonstrating that MEDRES can extract signals which are otherwise not available in CoF and Collab baselines.

8 Conclusion and Future Work

We considered the problem of constructing a recommendation system with "rich" users/items, with multiple entities expressed through multiple engagement/relationship graphs. We showed that our framework encapsulates several existing formulations, and used it to develop a novel MEDRES architecture that pools together signals from the graphs via Graph Neural Networks (GNN). We proposed AL-GCN, a novel GNN block that can be used with MEDRES architecture to learn graph weights and extract more powerful task-specific signals. Our training method can optimize MEDRES for the proposed pAp@k metric. Finally, we used our framework to model several different recommendation problems, including two real-world tasks, and demonstrated effectiveness of MEDRES against different baselines.

Due to the generality of our framework, it can be applied with several other embedding techniques like Heterogeneous Information Networks (HIN) that can be richer in some settings than GCN and AL-GCN. Furthermore, further exploring optimization of pAp@k metric and studying its theoretical properties is also an exciting research direction.

References

- [1] Hyung Jun Ahn. A new similarity measure for collaborative filtering to alleviate the new user cold-starting problem. *Information Sciences*, 178(1):37–51, 2008.
- [2] Rianne van den Berg, Thomas N Kipf, and Max Welling. Graph convolutional matrix completion. *arXiv preprint arXiv:1706.02263*, 2017.
- [3] Thang D Bui, Sujith Ravi, and Vivek Ramavajjala. Neural graph learning: Training neural networks using graphs. In *Proceedings of the Eleventh ACM International Conference on Web Search and Data Mining*, pages 64–71, 2018.
- [4] Stéphan J Cléménçon and Nicolas Vayatis. Empirical performance maximization for linear rank statistics. In *Advances in neural information processing systems*, pages 305–312, 2009.
- [5] Corinna Cortes and Mehryar Mohri. Auc optimization vs. error rate minimization. In *Advances in Neural Information Processing Systems*, pages 313–320, 2004.

- [6] Yuxiao Dong, Nitesh V Chawla, and Ananthram Swami. metapath2vec: Scalable representation learning for heterogeneous networks. In *Proceedings of the 23rd ACM SIGKDD international conference on knowledge discovery and data mining*, pages 135–144, 2017.
- [7] Tao-yang Fu, Wang-Chien Lee, and Zhen Lei. Hin2vec: Explore meta-paths in heterogeneous information networks for representation learning. In *Proceedings of the 2017 ACM on Conference on Information and Knowledge Management*, pages 1797–1806, 2017.
- [8] Manuel Gomez-Rodriguez, Krishna P Gummadi, and Bernhard Schoelkopf. Quantifying information overload in social media and its impact on social contagions. In *ICWSM*, pages 170–179, 2014.
- [9] Aditya Grover and Jure Leskovec. node2vec: Scalable feature learning for networks. In *Proceedings of the 22nd ACM SIGKDD international conference on Knowledge discovery and data mining*, pages 855–864, 2016.
- [10] Prateek Jain and Inderjit S Dhillon. Provable inductive matrix completion. *arXiv preprint arXiv:1306.0626*, 2013.
- [11] Purushottam Kar, Harikrishna Narasimhan, and Prateek Jain. Surrogate functions for maximizing precision at the top. In *International Conference on Machine Learning*, pages 189–198, 2015.
- [12] Guolin Ke, Qi Meng, Thomas Finley, Taifeng Wang, Wei Chen, Weidong Ma, Qiwei Ye, and Tie-Yan Liu. Lightgbm: A highly efficient gradient boosting decision tree. In *Advances in neural information processing systems*, pages 3146–3154, 2017.
- [13] Heung-Nam Kim, Inay Ha, Kee-Sung Lee, Geun-Sik Jo, and Abdulmotaleb El-Saddik. Collaborative user modeling for enhanced content filtering in recommender systems. *Decision Support Systems*, 51(4):772–781, 2011.
- [14] Diederik P Kingma and Jimmy Ba. Adam: A method for stochastic optimization. *arXiv preprint arXiv:1412.6980*, 2014.
- [15] Thomas N Kipf and Max Welling. Semi-supervised classification with graph convolutional networks. *arXiv preprint arXiv:1609.02907*, 2016.
- [16] Yehuda Koren, Robert Bell, and Chris Volinsky. Matrix factorization techniques for recommender systems. *Computer*, (8):30–37, 2009.
- [17] Quoc Le and Alexander Smola. Direct optimization of ranking measures. *arXiv preprint arXiv:0704.3359*, 2007.
- [18] Xin Li and Hsinchun Chen. Recommendation as link prediction in bipartite graphs: A graph kernel-based machine learning approach. *Decision Support Systems*, 54(2):880–890, 2013.
- [19] Tie-Yan Liu. Learning to rank for information retrieval. *Foundations and trends in information retrieval*, 3(3):225–331, 2009.

- [20] Tomas Mikolov, Ilya Sutskever, Kai Chen, Greg S Corrado, and Jeff Dean. Distributed representations of words and phrases and their compositionality. In *Advances in Neural Information Processing Systems*, pages 3111–3119, 2013.
- [21] Federico Monti, Davide Boscaini, Jonathan Masci, Emanuele Rodola, Jan Svoboda, and Michael M Bronstein. Geometric deep learning on graphs and manifolds using mixture model cnns. In *Computer Vision and Pattern Recognition*, volume 1, page 3, 2017.
- [22] Harikrishna Narasimhan and Shivani Agarwal. Support vector algorithms for optimizing the partial area under the roc curve. *Neural computation*, 29(7):1919–1963, 2017.
- [23] Jeffrey Pennington, Richard Socher, and Christopher D Manning. Glove: Global vectors for word representation. In *Proceedings of the 2014 conference on empirical methods in natural language processing (EMNLP)*, pages 1532–1543, 2014.
- [24] Bryan Perozzi, Rami Al-Rfou, and Steven Skiena. Deepwalk: Online learning of social representations. In *Proceedings of the 20th ACM SIGKDD international conference on Knowledge discovery and data mining*, pages 701–710, 2014.
- [25] Nikhil Rao, Hsiang-Fu Yu, Pradeep K Ravikumar, and Inderjit S Dhillon. Collaborative filtering with graph information: Consistency and scalable methods. In *Advances in Neural Information Processing Systems*, pages 2107–2115, 2015.
- [26] Karen Simonyan and Andrew Zisserman. Very deep convolutional networks for large-scale image recognition. *arXiv preprint arXiv:1409.1556*, 2014.
- [27] Otilia Stretcu, Krishnamurthy Viswanathan, Dana Movshovitz-Attias, Emmanouil Platanios, Sujith Ravi, and Andrew Tomkins. Graph agreement models for semi-supervised learning. In *Advances in Neural Information Processing Systems*, pages 8710–8720, 2019.
- [28] Jianing Sun and Yingxue Zhang. Multi-graph convolutional neural networks for representation learning in recommendation.
- [29] Jie Tang, Jing Zhang, Limin Yao, Juanzi Li, Li Zhang, and Zhong Su. Arnetminer: extraction and mining of academic social networks. In *Proceedings of the 14th ACM SIGKDD international conference on Knowledge discovery and data mining*, pages 990–998, 2008.
- [30] Hamed Valizadegan, Rong Jin, Ruofei Zhang, and Jianchang Mao. Learning to rank by optimizing ndcg measure. In *Advances in neural information processing systems*, pages 1883–1891, 2009.
- [31] Xiao Wang, Houye Ji, Chuan Shi, Bai Wang, Yanfang Ye, Peng Cui, and Philip S Yu. Heterogeneous graph attention network. In *The World Wide Web Conference*, pages 2022–2032, 2019.
- [32] Yueyang Wang, Yuanfang Xia, Siliang Tang, Fei Wu, and Yueting Zhuang. Flickr group recommendation with auxiliary information in heterogeneous information networks. *Multimedia Systems*, 23(6):703–712, 2017.

- [33] Jiaxuan You, Rex Ying, Xiang Ren, William L Hamilton, and Jure Leskovec. Graphrnn: Generating realistic graphs with deep auto-regressive models. In *Proceedings of the 35th International Conference on Machine Learning*, 2018.
- [34] Xiao Yu, Quanquan Gu, Mianwei Zhou, and Jiawei Han. Citation prediction in heterogeneous bibliographic networks. In *SIAM International Conference on Data Mining*, pages 1119–1130. SIAM, 2012.
- [35] Xiao Yu, Xiang Ren, Yizhou Sun, Bradley Sturt, Urvashi Khandelwal, Quanquan Gu, Brandon Norick, and Jiawei Han. Recommendation in heterogeneous information networks with implicit user feedback. In *Proceedings of the 7th ACM conference on Recommender systems*, pages 347–350, 2013.
- [36] Chuxu Zhang, Dongjin Song, Chao Huang, Ananthram Swami, and Nitesh V Chawla. Heterogeneous graph neural network. In *Proceedings of the 25th ACM SIGKDD International Conference on Knowledge Discovery & Data Mining*, pages 793–803, 2019.
- [37] Huan Zhao, Quanming Yao, Jianda Li, Yangqiu Song, and Dik Lun Lee. Meta-graph based recommendation fusion over heterogeneous information networks. In *Proceedings of the 23rd ACM SIGKDD International Conference on Knowledge Discovery and Data Mining*, pages 635–644, 2017.