

Large-scale Gender/Age Prediction of Tumblr Users

Yao Zhang*, Changwei Hu[†], Yifan Hu[†], Tejaswi Kasturi[‡],
Shanmugam Ramasamy[‡], Matt Gillingham[‡], and Keith Yamamoto[‡]

*LinkedIn [†]Yahoo Research [‡]Verizon Media

Email: *yaozhang@linkedin.com [†]{changweih, yifanhu, kasturi}@verizonmedia.com

[‡]{sramasamy8, mgilling, yamamoto}@verizonmedia.com

Abstract—Tumblr, as a leading content provider and social media, attracts 371 million monthly visits, 280 million blogs and 53.3 million daily posts¹. The popularity of Tumblr provides great opportunities for advertisers to promote their products through sponsored posts. However, it is a challenging task to target specific demographic groups for ads, since Tumblr does not require user information like gender and ages during their registration. Hence, to promote ad targeting, it is essential to predict user’s demography using rich content such as posts, images and social connections. In this paper, we propose graph based and deep learning models for age and gender predictions, which take into account user activities and content features. For graph based models, we come up with two approaches, network embedding and label propagation, to generate connection features as well as directly infer user’s demography. For deep learning models, we leverage convolutional neural network (CNN) and multilayer perceptron (MLP) to prediction users’ age and gender. Experimental results on real Tumblr daily dataset, with hundreds of millions of active users and billions of following relations, demonstrate that our approaches significantly outperform the baseline model, by improving the accuracy relatively by 81% for age, and the AUC and accuracy by 5% for gender.

I. INTRODUCTION

Online social media has become a ubiquitous part of our daily life, which allows us to easily share ideas/contents with other users, discuss social events/activities, and get connected with friends. Tumblr, with over 280 million blogs² and 130 billion blog posts, is one of the most popular social media apps. The rich content including text, images, and videos, provide great opportunities for advertisers to champion their products to specific groups. In particular, Tumblr offers “native advertisement” that allows advertisers to present their sponsored posts on the users’ interface. Native advertising has gained over 3 billion paid ad impressions in 2015 since it was started in 2012 [1]. However, unlike other social networks, Tumblr does not ask for some of the user demographic information, such as gender, during registration. Even though age is a required input during registration, it is difficult to verify the accuracy of this self-declared information. This makes it a challenge to precisely target ads for a group of users with age and gender related demography profiles. To improve the performance of user specific ad targeting, it is important for Tumblr to infer users’ age and gender from rich content generated by users.

*The work was done when the author was at Yahoo Research.

¹ <http://expandedramblings.com/index.php/tumblr-user-stats-fact/>

² A Tumblr user typically corresponds to one primary blog, so in this paper users and blogs are used interchangeably.

Several attempts have been made for this task, e.g., Grbovic et al. [1] proposed a gender and interest targeting framework that leverages user-generated data. Their model has lifted user engagement of ads, however, they did not take into consideration age prediction, which is one of the key factors for targeted ads. Furthermore, their gender prediction models only use features from blog contents and user activities (the so-called “cumulative features”). We believe that in addition to the cumulative features, user interactions, represented by the Tumblr following graph, can also be used for age and gender prediction, as users who followed each other tend to have similar interests or background. It is worth mentioning that the model in [1] does use blogs a user follows as a categorical feature for a linear classifier. However, it only exploits 1-hop neighbor information of the Tumblr following graph. The rich network structure provides further indications on how users interact. Furthermore, with the development in deep learning, advanced models like convolutional neural network (CNN) and multilayer perceptron (MLP) can typically produce better performance, which the past work like [1] does not explore.

Hence, to achieve better performance for age and gender prediction, in this paper, we propose a graph based approach with deep learning techniques that leverages the multitude of information encoded by the Tumblr following graph.

The main contributions of our paper include:

- *Data*: We leverage rich cumulative features including user activities and post contents. Furthermore, we construct a large Tumblr following graph with hundreds of millions of vertices and billions of edges.
- *Graph base methods*: We apply graph embedding and label propagation techniques to (1) generate rich features from the Tumblr following graph through node embedding, which is shown to be useful to improve users demography prediction; (2) directly leverage the label propagation algorithm to boost the prediction performance.
- *Deep learning*: We apply deep learning models including CNN and MLP to further improve the performance of age and gender prediction.
- *Evaluation*: We conduct empirical studies to show that the graph based and deep learning approaches can improve the AUC and accuracy performance by 5% relatively for gender prediction and the accuracy performance by 81% for age prediction, compared to the baseline model [1], and classifiers like GBDT, XGBOOST, etc.

The rest of this paper is organized as follows. Section II describes the Tumblr data, and Section III discusses our proposed methods, including network embedding, label propagation and deep learning. Empirical studies are shown in Section IV. Finally we summarize the related work in Section V, and conclude our work with future directions in Section VI.

II. PRELIMINARY

In this section, we discuss the rich Tumblr data we use, and the labels we create.

A. Tumblr Data

We use both following graph and cumulative features.

Following Graph. In Tumblr, users can follow other users, which forms the Tumblr following graph. Formally, we define $G(V, E)$ as the following graph, where V is the set of users and E is the set of following relations. Note that even though the following graph is directed in Tumblr, we assume G is an unweighted undirected graph, as our study focuses on age/gender prediction where we assume two users with a following relation tend to share mutual interests in their background. In this paper, G consists of hundreds of millions of nodes and billions of edges.

Cumulative Features. We utilize “cumulative features” [1], which contain blog content and activities including *music* (artist), *follow* (blogs followed), *likes* (of posts), *photo captions*, *post tags*, *blog titles*, and *blog descriptions*, etc.

Note that the baseline model [1] only used cumulative features with LOGISTICREGRESSION for the Tumblr age/gender prediction. It is interesting to study how the following graph can boost the performance.

B. Label Construction

Tumblr does not ask for gender information when a user signs up. Furthermore, although it does ask for user age, no independent verification is done, and Tumblr believes that the age information is unreliable. Therefore, it is a challenge to create ground truth labels. Grbovic et al. [1] leveraged the US census data that associates people’s names with gender to create gender labels. However, their approach suffers from a natural issue that some names can be neutral (e.g., “Avery”). Instead, we use age and gender ground truth data provided internally at Yahoo, which is at a large-scale and more reliable. This golden set provides us with good quality of ground truth label information. In total there are about 3 million ground truth age and gender labels.

III. PROPOSED METHODS

In this section, we will discuss our approaches for age and gender prediction, including network embedding, label propagation and deep learning. We first enrich the cumulative features with graph features generated by network embedding and label propagation, and then apply deep learning techniques for the prediction. In addition, we also use label propagation as a direct tool for the prediction, which is shown to be efficient in the experiments.

A. Network Embedding

When applying LOGISTICREGRESSION on the cumulative features for the baseline model, we found that *follow* and *like* are the two most important features, which suggests that blogs that a user follows/reads is a good indicator on the user’s demographic. Intuitively, users tend to follow others’ activities with similar age/gender. Using *follow* features directly is akin to using words as unigram features for natural language processing. It works to some extent, but fails to uncover underlying connections among blogs. For example, if two blogs have a large overlap of followers, though they are not directly followed by each other, they are somewhat related. However, network embedding (mapping blogs to a high dimensional space based on their follow relations) can capture such relationships.

In this paper, we apply word2vec [2], [3] to generate rich network features. Word2vec is one of the most popular embedding techniques to capture syntactic and semantic relationships of entities. It leverages a multilayer perceptron (MLP) with two architectures, bag-of-words (BOW) and skip-gram. Levy and Goldberg [4] found that the skip-gram model, which uses a word to predict its neighborhood words, is equivalent to a matrix factorization with the matrix to be factored containing the word-word point-wise mutual information. With this broad understanding of word2vec in mind, the technique is applicable to tasks beyond those of traditional natural language processing tasks. It can be deployed to applications with entities and their co-occurrence patterns. Since blogs and their followers have such a pattern, we leverage word2vec for our task.

Implementation details. Network embedding has been well-studied recently. For example, DeepWalk [5] constructs “sentences” of vertices by random walks, and then applies a word2vec model to sentences to get vector representation of vertices. Vertices that appear frequently together in sentences tend to be similar. Our implementation is similar to LINE [6], where we create a sentence from a vertex and its neighbors. Specifically, each blog is combined with all the k blogs it follows. This set of $k + 1$ blogs are randomly permuted and is treated as one sentence. All together, the following graph G produces hundreds of millions of “sentences” with billions of words, all of which are embedded into 50 dimensions using word2vec in the skip-gram mode.

We use a word2vec implementation with a minimum word count of 5, which results in a dictionary size of 46 million. In order to use blog embedding successfully as features, we embed a majority of the blogs. For blogs without an embedding, we construct its embedding by averaging the embeddings of its neighbors that have an embedding. We denote the embedding algorithm as Emb .

B. Label Propagation

Emb is an efficient method to generate graph features. However, it is an unsupervised approach which does not take into account label information for training. Next we develop

semi-supervised learning algorithm based on label propagation (called `LabelProp`), which takes into account labels and multi-hop neighborhood information. `LabelProp` serves as an approach to directly predict age/gender, as well as generate additional graph features. In addition, compared to the time consuming training of `Emb` for large graphs, `LabelProp` is a faster algorithm that can be naturally implemented in a distributed environment with a parallel fashion using Spark on Hadoop.

`LabelProp` spreads existing age/gender labels over the following graph. The intuition is that users tend to follow blogs with the similar background (e.g., ladies fashion). The age/gender of such a popular blog, defined by its followers, in turn defines the age/gender of its follower whose demographic information may be unknown to us. Next we will introduce the `LabelProp` algorithm and its variants, and then discuss how to leverage it for feature generation.

The `LabelProp` Algorithm. Given the Tumblr following graph $G(V, E)$, we denote the set of labeled users as $V_\ell = \{v_1, \dots, v_\ell\}$ where $V_\ell \subseteq V$, and the set of unlabeled users is denoted as V_m where $V_m = V - V_\ell$. For $v_i \in V_\ell$, let y_i be v_i 's label. The label propagation is an iterative process: at each iteration, a node's label propagates to its 1-hop neighbors. The process starts with V_ℓ , and once a node gets labels from its neighboring nodes, it will propagate labels in the next iterations. The process ends when it converges or reach a predefined number of iterations. Note that during the process, we fix the labels of nodes in V_ℓ , i.e., only nodes without labels at the beginning will iteratively update labels.

Algorithm 1 shows the *Label Propagation* (`LabelProp`) algorithm. We iteratively update label y_i at each iteration (Line 4), where $N(i)$ is the set of neighboring nodes of v_i that have labels at previous iterations. The algorithm ends until we reach K iterations. Note that we also add a hyperparameter $\alpha \in [0, 1]$ to control contributions of neighboring labels. If α is low, we focus more on the neighboring information and vice versa.

Algorithm 1 `LabelProp`(G, V_ℓ, α, K)

Require: G, V_ℓ, α , and iteration number K

```

1:  $k = 1$ 
2: while  $k \leq K$  do
3:   for  $v_i \in V_m$  do
4:      $y_i^{(k+1)} := \alpha y_i^{(k)} + (1 - \alpha) \frac{1}{|N(i)|} \sum_{j \in N(i)} y_j^{(k)}$ 
5:   end for
6:    $k := k + 1$ 
7: end while
8: return  $y_i^{(K)}$ 

```

Implementation details. The label here can be either gender or age. For gender, we assume $y_i = 1$ as female and $y_i = 0$ as male. For age, we first split age into 7 buckets ($< 17, 18 - 24, 25 - 34, 35 - 44, 45 - 54, 55 - 64, 65+$), and use one hot encoding to create a label vector with 7 entries. We run `LabelProp` for each entry. It is natural to develop

`LabelProp` in parallel, as each node updates its label by independently querying its neighbors at each iteration. We implement `LabelProp` under the Pregel framework [7] on Spark, which is a state-of-the-art message passing parallel model for large distributed graphs.

Note that we do not run `LabelProp` until it converges, but stop at a predefined maximum number of iterations because of the scalability consideration (we may need a large amount of iterations to converge). Furthermore, as we mentioned above, the motivation of `LabelProp` is that users that have following relations have similar demography. Hence, labels propagated to a node at a large number of iteration k , may “contaminate” the prediction of u , as the propagated labels at k -th iteration can be different from the true label.

Variants of label propagation. We propose two variants of `LabelProp` for the gender prediction in order to investigate how labels from nodes with different distance make impact on the label prediction. Note that it is straightforward to generalize the analysis to the multiclass age prediction. The main idea is that instead of using a constant hyperparameter α , we adaptively change the parameter. As the value of the iteration increases, we decrease the contribution of incoming labels from neighboring nodes. Hence, labels that are far away become less important in Algorithm 1 (Line 4). If they do not affect the performance, then the label prediction results with large iterations will remain competitive compared to the case with small iterations. The followings are two alternative propagation strategies we propose:

- 1) β -strategy: $y_i^{(k+1)} := (1 - \beta^k) y_i^{(k)} + \beta^k \sum_{j \in N(i)} \frac{y_j^{(k)}}{|N(i)|}$.
- 2) γ -strategy: $y_i^{(k+1)}(m/f) := y_i^{(k)}(m/f) + \gamma \sum_{j \in N(i)} \frac{y_j^{(k)}(m/f)}{|N(i)|}$

In the first strategy, $\beta \in [0, 1]$ and k is the number of the iterations. As the iteration value increases, the impact of labels from the neighboring information decreases exponentially. In the second strategy, $\gamma \in [0, 1]$. We propagate male and female labels respectively for each node, and at the end normalize them as the final label for each node. As shown in the experiments, these two strategies can empirically demonstrate that for a node u , labels from its close neighbors are more important than long-distance labels for gender prediction.

Label Propagation as Features. In `LabelProp`, we can directly learn labels by the iteration rule (Algorithm 1 Line 4). As shown in the next experiments, it provides high quality predictions compared to baselines. However, neither does it take into account the cumulative features, nor features from `Emb`. `Emb` generates graph features in an unsupervised manner without considering existing label information. To leverage the label information to improve the prediction performance, we use `LabelProp` to general features as well.

The idea of `LabelProp` as feature generators comes from the ensemble methods: we divide the labeled data uniformly at random into N partitions. For partition i , we denote it as V_ℓ^i . We run `LabelProp` on the whole graph with the labeled data V_ℓ^i for each i respectively to get the learned label y_u^i , and

concatenate learned labels. For each node u , label propagation features are represented as a vector $\mathbf{Y}_u = [y_u^1, \dots, y_u^N]^T$. For $u \in V_\ell$, if it is in partition i (meaning we know its label y_u), instead of using y_u as the feature at dimension i , we set $y_u^i = \frac{\sum_{j \neq i} y_u^j}{N-1}$. This is because directly putting y_u as features will cause overfitting.

C. Deep Learning Models

Deep learning models have been extensively studied due to their extraordinary performance in a variety of areas including computer vision, natural language processing, robotics, etc. The baseline model used the LOGISTICREGRESSION method for age/gender prediction. In this paper, we aim to leverage deep learning models to improve users' demography prediction performance. To start with, we used multilayer perceptron (MLP), a class of feed-forward neural net with more than three layers, with embedding as features. In addition, we experimented with a convolutional neural network (CNN) [8], which has shown promising and reliable performances across a range of text classification tasks by leveraging embedding results. This is a word embedding based CNN using the text features.

As a next step we also tried RESNET [9], whose architecture has more layers (18 for our case), with each layer also connecting to 2 layers in the front. This skip connection adds the ability to train deeper models to avoid overfitting. The accuracy for this model is 1% higher than MLP with 3 hidden layers. However, RESNET took almost 5 times time to train. So we decided not to pursue this direction.

Implementation details. MLP and CNN are two advanced models which typically take long times to train over large datasets, especially for the Tumbler data with hundreds of millions of users. To speed up the training process, we sample the same number of female and male examples. Empirical studies over the whole dataset indeed show the model trained from the sampled users can still provide us with convincing performance over the baseline model (LOGISTICREGRESSION). We implemented MLP and CNN using the Keras library and Tensorflow kernels. For both MLP and CNN, we use cross entropy as the objective, and stochastic gradient descent (SGD) as the optimizer.

IV. EXPERIMENTS

We conduct our experiments using the Spark framework on Hadoop with 300 executors, each of which has 12G memory. LabelProp is implemented using the Spark GraphX library, while Emb is implemented using a multi-threaded version of word2vec.

We obtain about 3 million labeled data as discussed before. To evaluate the performance, we randomly sample 70% of the labeled data as the training set, and the rest as the testing set. We fit our model on the training set, and compute accuracy on the testing data. For gender prediction we also report AUC as it is a binary classification task.

TABLE I
PERFORMANCE OF AGE PREDICTION.

	Accuracy	Cross-Entropy Loss
CNN +MLP	0.3897	1.5066
MLP	0.3841	1.5212
CNN	0.3482	1.6157
LOGISTICREGRESSION (baseline with CumF)	0.2150	-

TABLE II
PERFORMANCE (AUC) OF LOGISTICREGRESSION ON HADOOP.

CumF	Emb	LabelProp	CumF+LP	All
0.8489	0.8413	0.8628	0.8831	0.8858

A. Results

In short, we demonstrate that CNN and MLP can achieve 0.38 of the accuracy for the multi-class (7 classes) age prediction, which outperforms the baseline model by 81% in accuracy. In addition, we show that adding features from Emb and LabelProp can greatly improve AUC for gender prediction by 5% compared to the baseline model. Finally, we explore different performance of LabelProp as iteration K and hyperparameter α change, as well as the variants of LabelProp.

Age Prediction. As discussed above, for the age prediction, we split ages into 7 buckets ($< 17, 18 - 24, 25 - 34, 35 - 44, 45 - 54, 55 - 64, 65+$) as labels, and use CumF and Emb as features. The baseline is a generalized multiclass LOGISTICREGRESSION, which uses softmax instead of logistic function as the objective, and CumF. In addition to CNN and MLP, we also combine them together (CNN +MLP) to boost the performance, by concatenating the last hidden layers of both and feeding to them to the cross entropy optimizer as the output. Table I shows the results: combining CNN and MLP produces the best results by 81.3% of performance improvement in accuracy compared to the baseline model LOGISTICREGRESSION. In addition, MLP outperforms CNN by 10%.

Gender Prediction. Table II shows the results of using different features on the Hadoop system. CumF is the result that only uses cumulative features [1], which are used in the baseline. Emb and LabelProp are the results that use the embedding and label propagation features respectively, CumF+LP is the result that combines both. All is the result that puts cumulative, node embedding and label propagation features together. First, All produces the best AUC by 5% of performance improvement. Second, it is interesting that CumF+LP also gives us competitive results: it achieves only 0.2% of the performance loss compared to All. Finally, when we separate features respectively, LabelProp outperforms CumF and Emb by 2%.

In addition to testing our model on Hadoop, we also use the following classifiers as baselines for the gender prediction on a single machine. Table III shows the quality (quantifying comparisons) of different feature integration approaches in

terms of AUC and accuracy. For `Emb`, we embed users into 50 dimensional feature space, while for `LabelProp`, we are able to obtain 5 features. First, our model MLP outperforms other classifiers by up to 10% of performance improvement in both AUC and accuracy. Second, stacking `Emb` and `LabelProp` features together can improve AUC and accuracy (up to 5% of the improvement). Compared to the baseline LOGISTICREGRESSION model with cumulative features, MLP with `Emb` and `LabelProp` features improved the AUC by 5%.

Sensitivity of LabelProp. We also evaluate the sensitivity of `LabelProp` for hyperparameter α and iteration number K . Table IV shows the results. First, when iteration number $K = 1$, `LabelProp` just queries labels from 1-hop neighbors, hence, the performance is not as good as examining multi-hop neighbors. Second, as the number of iterations goes up, the AUC first increases and then slightly decreases. This empirically confirms our intuition that labels that are close to a node u make more contributions to the labeling result of u . In addition, we notice that the convergence rate is different for different parameters. In general, it takes around $K = 50$ iterations for `LabelProp` to converge. However, in practice, since we are interested in the prediction performance, we only need to run a few iterations to obtain the best result. With regard to the parameter α , we observe that the smaller α is, the fewer iterations are needed to achieve the best AUC. For example, when $\alpha = 0.2$, we only need 2 iterations, while 10 iterations are needed for the best AUC when $\alpha = 0.8$.

Variants of LabelProp. The variants of `LabelProp` proposed in Section III are used to examine how the distance of neighbors affects results, and how many iterations are sufficient to get good prediction. Different from the `LabelProp` with the constant α , both β -strategy and γ -strategy decrease the weight on neighboring labels. As shown in Table V, both strategies get the best AUC around 4-5 iterations, which suggests that labels within 4 or 5-hops are very important for `LabelProp` to obtain good results. As the iteration number increases, the AUC almost remains the same especially for the γ -strategy, which suggests that labels with large distance make little impact on the prediction.

V. RELATED WORK

Demographic Ad Targeting. Personalized advertising is a common ads targeting strategies, which has been studied broadly [10], [11], [12]. It tries to display the most relevant ads to each individual. Demographic ad targeting is one of the personalization tactics, which effectively targets users based on their age, gender, etc. Grbovic et al. [1] first studied the gender based ad targeting in Tumblr. However, their model is based on learned gender labels, which sometimes can be unreliable. To address their issue, we provide high accuracy labels to improve gender prediction. In addition to demographic ad targeting, another type of personalized advertising used in Tumblr is interest targeting which recommends ads based on their categories [1], [13], [14].

Network Embedding. The graph embedding problem tries to generate vector representation of nodes. Previous work such as locally linear embedding [15], IsoMap [16], and spectral techniques [17], [18], [19], treats networks as matrices. They tend to be slow and do not scale to large networks. Recently, leveraging the deep learning techniques, several novel algorithms were proposed to learn feature representations of nodes [5], [20], [21], [22]. DeepWalk [5] and Node2Vec [22] extends the word2vec model [23] to networks by leveraging random walks to generate “word context”. SDNE [21] and LINE [20] learn embedding results on directed graphs that maintain the first and second order proximity of nodes. In this paper, our `Emb` implementation leverages Line to learn graph features of users.

Label Propagation. The idea of label propagation has been widely studied in the machine learning literature. Zhu et al. [24] first leveraged label propagation for a graph based semi-supervised learning algorithm, and Zhu et al. [25] further proposed an iterative algorithm from a close form solution for label propagation. After that label propagation has been applied to many domains, such as multimedia including image and video data [26] and information retrieval like relevance and keyword search [27], [28]. In addition, Rao and Yarowsky [29] proposed a parallel label propagation algorithm under the MapReduce framework. However, none of the above work studied the problem with large-scale data like ours, and implemented the label propagation algorithm on large distributed systems in practice, nor were these applied to demography classification.

VI. CONCLUSION AND FUTURE WORK

In this paper, we study the important and challenging problem of Tumblr age and gender prediction for large scale Tumblr data, which can be used to better target sponsored ads against specific demographic audiences. We propose to add the following graph information to the existing cumulative features in Tumblr to enhance the age and gender prediction performance. In particular, we leverage graph embedding and label propagation techniques to generate informative user features, and apply deep learning models including CNN and MLP to utilize these features. Experimental results demonstrate that our approaches outperforms the baseline models by relatively 81% of accuracy improvement for age, and 5% of accuracy and AUC improvement for gender.

As future work, we would like to incorporate Tumblr’s raw age signals either as a feature, or as a way to gain more trustworthy labels in the age prediction model training. In addition, the Tumblr data consists of many images. If we could leverage the images and their annotations as extra features, then this could possibly boost the performance further.

REFERENCES

- [1] M. Grbovic, V. Radosavljevic, N. Djuric, N. Bhamidipati, and A. Nagarajan, “Gender and interest targeting for sponsored post advertising at tumblr,” in *Proceedings of the 21th ACM SIGKDD International Conference on Knowledge Discovery and Data Mining*, ser. KDD ’15. New York, NY, USA: ACM, 2015, pp. 1819–1828. [Online]. Available: <http://doi.acm.org/10.1145/2783258.2788616>

TABLE III
COMPARISON OVER VARIOUS CLASSIFIERS. LEFT: AUC; RIGHT: ACCURACY.

AUC	Emb	LabelProp	Emb+LP	Accuracy	Emb	LabelProp	Emb+LP
MLP	0.883	0.865	0.899	MLP	0.811	0.804	0.830
LOGISTICREGRESSION	0.842	0.864	0.869	LOGISTICREGRESSION	0.782	0.783	0.806
XGBOOST	0.841	0.865	0.879	XGBOOST	0.777	0.804	0.812
GBDT	0.831	0.865	0.882	GBDT	0.768	0.804	0.815
LINEARSVM	0.762	0.784	0.798	LINEARSVM	0.779	0.772	0.804

TABLE IV
PERFORMANCE (AUC) OF LABELPROP AS ITERATION K AND HYPERPARAMETER α CHANGE.

$\alpha \backslash K$	1	2	3	4	5	6	7	8	9	10	15	20
0.2	0.770	0.881	0.878	0.878	0.871	0.868	0.866	0.865	0.863	0.862	0.860	0.858
0.5	0.770	0.853	0.873	0.877	0.877	0.876	0.874	0.872	0.870	0.868	0.863	0.861
0.8	0.770	0.816	0.834	0.848	0.858	0.865	0.869	0.872	0.874	0.875	0.875	0.872

TABLE V
RESULTS (AUC) OF LABELPROP VARIATIONS. K IS THE NUMBER OF ITERATIONS, $\beta = 0.8$, AND $\gamma = 0.9$

K	1	2	3	4	5	10
β -strategy	0.770	0.850	0.871	0.877	0.878	0.869
γ -strategy	0.770	0.868	0.877	0.878	0.878	0.875

- [2] T. Mikolov, K. Chen, G. Corrado, and J. Dean, "Efficient estimation of word representations in vector space," in *Proceedings of Workshop at ICLR*, 2013.
- [3] J. Pennington, R. Socher, and C. D. Manning, "Glove: Global vectors for word representation," *Proceedings of the Empirical Methods in Natural Language Processing (EMNLP 2014)*, vol. 12, 2014.
- [4] O. Levy and Y. Goldberg, "Neural word embedding as implicit matrix factorization," in *Advances in Neural Information Processing Systems 27: Annual Conference on Neural Information Processing Systems 2014, December 8-13 2014, Montreal, Quebec, Canada*, 2014, pp. 2177–2185.
- [5] B. Perozzi, R. Al-Rfou, and S. Skiena, "Deepwalk: Online learning of social representations," in *Proceedings of the 20th ACM SIGKDD international conference on Knowledge discovery and data mining*. ACM, 2014, pp. 701–710.
- [6] J. Tang, M. Qu, M. Wang, M. Zhang, J. Yan, and Q. Mei, "Line: Large-scale information network embedding," in *Proceedings of the 24th International Conference on World Wide Web*, ser. WWW '15. Republic and Canton of Geneva, Switzerland: International World Wide Web Conferences Steering Committee, 2015, pp. 1067–1077. [Online]. Available: <https://doi.org/10.1145/2736277.2741093>
- [7] G. Malewicz, M. H. Austern, A. J. Bik, J. C. Dehnert, I. Horn, N. Leiser, and G. Czajkowski, "Pregel: a system for large-scale graph processing," in *Proceedings of the 2010 ACM SIGMOD International Conference on Management of data*. ACM, 2010, pp. 135–146.
- [8] Y. Kim, "Convolutional neural networks for sentence classification," 08 2014.
- [9] K. He, X. Zhang, S. Ren, and J. Sun, "Deep residual learning for image recognition," in *The IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, June 2016.
- [10] D. Essex, "Matchmaker, matchmaker," *Communications of the ACM*, vol. 52, no. 5, pp. 16–17, 2009.
- [11] A. Z. Broder, "Computational advertising and recommender systems," in *Proceedings of the 2008 ACM conference on Recommender systems*. ACM, 2008, pp. 1–2.
- [12] A. Majumder and N. Shrivastava, "Know your personalization: learning topic level personalization in online services," in *Proceedings of the 22nd international conference on World Wide Web*. ACM, 2013, pp. 873–884.
- [13] D. Shin, S. Cetintas, and K.-C. Lee, "Recommending tumblr blogs to follow with inductive matrix completion," in *RecSys Posters*, 2014.
- [14] N. Barbieri, F. Bonchi, and G. Manco, "Who to follow and why: link prediction with explanations," in *Proceedings of the 20th ACM SIGKDD international conference on Knowledge discovery and data mining*. ACM, 2014, pp. 1266–1275.
- [15] S. T. Roweis and L. K. Saul, "Nonlinear dimensionality reduction by locally linear embedding," *science*, vol. 290, no. 5500, pp. 2323–2326, 2000.
- [16] J. B. Tenenbaum, V. De Silva, and J. C. Langford, "A global geometric framework for nonlinear dimensionality reduction," *science*, vol. 290, no. 5500, pp. 2319–2323, 2000.
- [17] F. R. Chung, *Spectral graph theory*. American Mathematical Soc., 1997, vol. 92.
- [18] F. R. Bach and M. I. Jordan, "Learning spectral clustering," in *NIPS*, vol. 16, 2003.
- [19] M. Belkin and P. Niyogi, "Laplacian eigenmaps and spectral techniques for embedding and clustering," in *NIPS*, vol. 14, no. 14, 2001, pp. 585–591.
- [20] J. Tang, M. Qu, M. Wang, M. Zhang, J. Yan, and Q. Mei, "Line: Large-scale information network embedding," in *Proceedings of the 24th International Conference on World Wide Web*. ACM, 2015, pp. 1067–1077.
- [21] D. Wang, P. Cui, and W. Zhu, "Structural deep network embedding," in *Proceedings of the 22nd ACM SIGKDD International Conference on Knowledge Discovery and Data Mining*. ACM, 2016, pp. 1225–1234.
- [22] A. Grover and J. Leskovec, "node2vec: Scalable feature learning for networks," in *Proceedings of the 22nd ACM SIGKDD International Conference on Knowledge Discovery and Data Mining*. ACM, 2016, pp. 855–864.
- [23] T. Mikolov, I. Sutskever, K. Chen, G. S. Corrado, and J. Dean, "Distributed representations of words and phrases and their compositionality," in *Advances in neural information processing systems*, 2013, pp. 3111–3119.
- [24] X. Zhu, Z. Ghahramani, and J. D. Lafferty, "Semi-supervised learning using gaussian fields and harmonic functions," in *Proceedings of the 20th International conference on Machine learning (ICML-03)*, 2003, pp. 912–919.
- [25] D. Zhou, O. Bousquet, T. N. Lal, J. Weston, and B. Schölkopf, "Learning with local and global consistency," in *Advances in neural information processing systems*, 2004, pp. 321–328.
- [26] W.-Y. Lee, L.-C. Hsieh, G.-L. Wu, and W. Hsu, "Graph-based semi-supervised learning with multi-modality propagation for large-scale image datasets," *Journal of Visual Communication and Image Representation*, vol. 24, no. 3, pp. 295–302, 2013.
- [27] J. He, M. Li, H.-J. Zhang, H. Tong, and C. Zhang, "Manifold-ranking based image retrieval," in *Proceedings of the 12th annual ACM international conference on Multimedia*. ACM, 2004, pp. 9–16.
- [28] H. Tong, J. He, M. Li, C. Zhang, and W.-Y. Ma, "Graph based multi-modality learning," in *Proceedings of the 13th annual ACM international conference on Multimedia*. ACM, 2005, pp. 862–871.
- [29] D. Rao and D. Yarowsky, "Ranking and semi-supervised classification on large scale graphs using map-reduce," in *Proceedings of the 2009 Workshop on Graph-based Methods for Natural Language Processing*. Association for Computational Linguistics, 2009, pp. 58–65.