

REFINE: Random Range Finder for Network Embedding

Hao Zhu

Australian National University
Data 61/CSIRO
allenhaozhu@gmail.com

Piotr Koniusz*

Data 61/CSIRO
Australian National University
name.surname@data61.csiro.au

Abstract

Network embedding approaches have recently attracted considerable interest as they learn low-dimensional vector representations of nodes. Embeddings based on the matrix factorization are effective but they are usually computationally expensive due to the eigen-decomposition step. In this paper, we propose a Random Range Finder based Network Embedding (REFINE) algorithm, which can perform embedding on one million of nodes (YouTube) within 30 seconds in a single thread. REFINE is $10\times$ faster than ProNE, which is $10-400\times$ faster than other methods such as LINE, DeepWalk, Node2Vec, GraRep, and Hope. Firstly, we formulate our network embedding approach as a skip-gram model, but with an orthogonal constraint, and we reformulate it into the matrix factorization problem. Instead of using randomized tSVD (truncated SVD) as other methods, we employ the Randomized Blocked QR decomposition to obtain the node representation fast. Moreover, we design a simple but efficient spectral filter for network enhancement to obtain higher-order information for node representation. Experimental results prove that REFINE is very efficient on datasets of different sizes (from thousand to million of nodes/edges) for node classification, while enjoying a good performance.

1 Introduction

Network embeddings have drawn a lot of interest due to their ability to produce low-dimensional representation for nodes while encapsulating the structure/properties of the network [5]. Such representations serve as latent features for off-the-shelf machine learning for a variety of tasks on graphs *e.g.*, node classification [16, 22, 34], link prediction [6], graph classification [12], relation learning [27], trajectory analysis [17, 21, 20], spatio-temporal graphs [10, 11] and network reconstruction [26].

Many network embedding approaches are based on random walks [16], matrix factorization [19] and deep learning [26]. The skip-gram model has significantly advanced network embeddings *e.g.*, DeepWalk [16], LINE [24], PTE [23], and Node2Vec [6] approaches, which sample node pairs from k -step transition matrices with different values of k , and train a skip-gram [14] model on these pairs to get node embeddings. The above methods can be unified into the closed-form Matrix Factorization (MF) framework [19] with two steps: (i) building a higher-order proximity matrix and (ii) obtaining the node embedding by using eigen-decomposition. Although MF provides an efficient way to obtain network embeddings compared with skip-gram based methods, it faces the computational and space challenges on large scale networks ($\geq 10,000$ nodes) as these two steps depend on SVD and dense proximity matrices.

Many methods use different approximation approaches to accelerate the MF based network embedding because the SVD, especially on dense matrices, limits the ability to scale up the efficient network

*The corresponding author. The code is available at: <https://github.com/allenhaozhu/REFINE>.

embedding. Some methods bypass higher-order proximity computations [31], or sparsify higher-order proximity matrices [18] *e.g.*, they use the randomized SVD [8] to achieve acceleration. Although [31, 18] are very fast network embedding approaches based on matrix factorization, they depend on SVD, which limits their use on large scale networks. Therefore, approach [32] employs random projections, a simple and powerful technique, which forms a low-dimensional embedding space for the network while preserving the original graph structure. Although network embeddings based on random projections are very efficient, the performance is worse compared to learning-based approximation methods [31, 18].

We present a new method to offer a trade-off between speed [32] and performance [31, 18]. By adding an orthogonal constraint on context vectors, we propose another Skip-Gram Network Embedding (SGNE) framework. Subsequently, we obtain node representations by a Randomized Blocked QR with power iteration (range finder), which is much faster than SVD-based methods. Finally, we present an efficient and effective spectral filter to enhance network embeddings by obtaining higher-order structural information of neighborhoods. Our contributions are:

- i. We propose a new skip-gram network embedding by adding orthogonal constraints, making it an MF problem free of SVD.
- ii. We propose a computationally more efficient than tSVD range finder based on the Randomized Blocked QR with power iteration.
- iii. We propose a compact and effective spectral filter to enhance node representations.
- iv. We validate the effectiveness/performance of REFINE on several standard benchmarks, and show it yields network embeddings on one million of nodes within 30 seconds (a single CPU thread), which is $\geq 10\times$ faster compared to the state of the art.

2 Related Work

Many off-the-shelf ML techniques leverage network embeddings, which we review below. A popular deep embedding model called DeepWalk [16] uses truncated random walks (to explore the network structure) and the skip-gram word embedding model [14] (to obtain embedding vectors of nodes). LINE [24] sets the walk length as one, and it introduces the negative sampling strategy [14] to accelerate training. Node2Vec [6] generalizes the above two methods and modifies the definition of neighborhood. The above methods are equivalent to factorizing a higher-order proximity matrix [19].

Many explicit matrix factorization methods have been used for network embeddings. GraRep [2] applies SVD to preserve higher-order proximity matrices (time complexity $\mathcal{O}(|V|^3)$). HOPE [15] uses generalized SVD to preserve the asymmetric transitivity in directed networks. Community structure (a mesoscopic structure of network) is preserved by non-negative matrix factorization in [28]. Approach [4] uses the matrix factorization with sparsification to accelerate SVD. Another approximate matrix factorization technique is used by approach [29]. AROPE [33] improves upon the above works by preserving arbitrary-order proximity simultaneously.

Many methods accelerate factorization *e.g.*, ProNE [31] avoids computing the higher-order proximity matrix by initializing embeddings with a low-order proximity matrix and applying a graph filter to improve the performance. Approach [18] uses higher-order proximity and then employs random-walk polynomial sparsification to higher-order proximity matrix. The above methods use a sparse proximity matrix with randomized tSVD which is much faster in obtaining the network embedding than SVD on a dense matrix.

3 Methodology

Below, we present our approach. Firstly, we propose a new objective function for SGNE, by imposing an orthogonal constraint on context vectors. We rewrite the objective function as the matrix factorization. Secondly, we present a range estimator based Randomized Blocked QR with power iterations to solve the matrix factorization. Finally, to enhance network embedding with higher-order proximity, we propose a simple spectral filter, approximated by the second-order Taylor expansion.

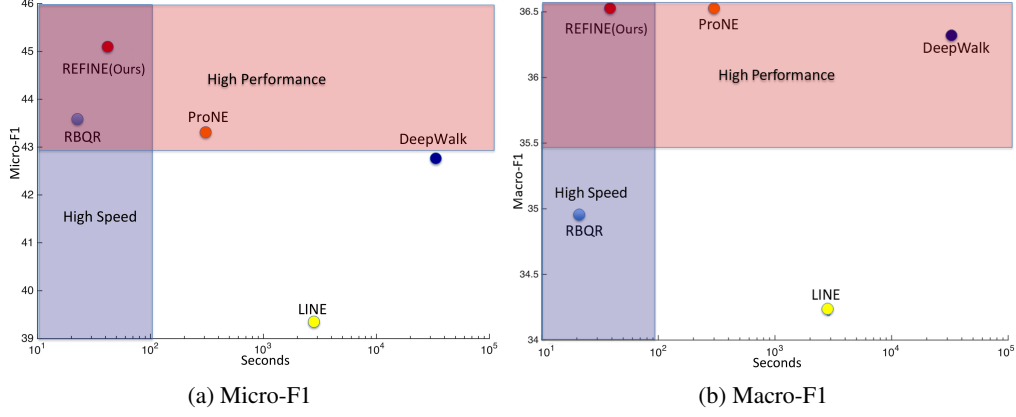


Figure 1: The Micro-F1 (and Macro-F1) vs. runtime for different methods on the YouTube dataset (90% of dataset used for training).

3.1 SGNE as Column Pivoting QR Factorization

Based on the Skip-Gram with Negative Sampling (SGNS) model[14, 31], we introduce an orthogonal constraint on context vectors, which leads to a very fast optimization procedure. Firstly, consider the objective function:

$$l = - \sum_{(i,j) \in \hat{E}} [p_{i,j} \ln \sigma(\mathbf{r}_i^\top \mathbf{c}_j) + \lambda \phi(\hat{E}, j) \ln \sigma(-\mathbf{r}_i^\top \mathbf{c}_j)], \text{ s.t. } \mathbf{C}^\top \mathbf{C} = \mathbf{I}, \quad (1)$$

where $\lambda \geq 0$ is a coefficient controlling the negative noise sample ratio, $\sigma(\cdot)$ is the sigmoid function, context vector \mathbf{c}_j is the j -th row of context matrix $\mathbf{C} \in \mathbb{R}^{n \times k}$, whereas $\mathbf{r}_i \in \mathbb{R}^k$ are k -dimensional embeddings and $p_{i,j}$ is the (i, j) -th coefficient of the degree-normalized adjacency matrix. Finally, $\phi(\hat{E}, j)$ forms the empirical context for node j given the edge set \hat{E} associated with j , given as:

$$\phi(\hat{E}, j) = \frac{\sum_{i: (i,j) \in \hat{E}} p_{i,j}}{\sum_{(i',j') \in \hat{E}} p_{i',j'}}. \quad (2)$$

Without the constraint $\mathbf{C}^\top \mathbf{C} = \mathbf{I}$, Eq. (1) has a trivial solution $c_{ij} \rightarrow \infty$. The constraint fulfills two different roles: i) it bounds context vectors to be normalized to the ℓ_2 unit norm, and (ii) it decorrelates them. Thus, \mathbf{C} can be thought of as a subspace as $k \ll n$.

A sufficient condition for minimizing the objective (1) is to let its partial derivative with respect to $\mathbf{c}_j^\top \mathbf{r}_i$ be equal zero, thus:

$$\mathbf{r}_i^\top \mathbf{c}_j = \ln \frac{p_{i,j}}{\lambda \phi(\hat{E}, j)}, \quad (i, j) \in \hat{E}. \quad (3)$$

If $M_{ij} = \mathbf{r}_i^\top \mathbf{c}_j$, \mathbf{M} would result in a dense matrix, the source of inefficiency in MF. Thus, ProNE [31] defines a matrix \mathbf{M} with entries:

$$M_{i,j} = \begin{cases} \ln \frac{p_{i,j}}{\lambda \phi(\hat{E}, j)}, & (i, j) \in \hat{E} \\ 0, & (i, j) \notin \hat{E}. \end{cases} \quad (4)$$

Thus, we obtain an approximate sparse matrix \mathbf{M} with the low-rank product of embedding matrix \mathbf{R} and the context embedding matrix \mathbf{C} . We avoid the truncated Singular Value Decomposition (tSVD), and achieve the optimal rank k factorization w.r.t. the p' norm by solving:

$$\min_{\mathbf{R}, \mathbf{C}} \|\mathbf{M} - \mathbf{M}^*\|_{p'}, \text{ s.t. } \mathbf{M}^* = \mathbf{R}\mathbf{C}^\top, \mathbf{C}^\top \mathbf{C} = \mathbf{I}, \quad (5)$$

where \mathbf{R} and \mathbf{C} are $n \times k$ matrices whose rows stand for a node embedding and context embedding, respectively. Note that this objective is different from other objectives in [31, 18, 19]. Intuitively, through SVD, the matrix \mathbf{M} can be decomposed into $\mathbf{M} = \mathbf{U}\mathbf{\Sigma}\mathbf{V}^\top$, where $\mathbf{U} = [\mathbf{u}_1, \dots, \mathbf{u}_k]$ is orthogonal and $\mathbf{\Sigma} = \text{diag}([\sigma_1, \dots, \sigma_k])$ with $\sigma_1 \geq \sigma_2 \geq \dots \geq \sigma_k$. Thus, $\mathbf{R} = \mathbf{U}\mathbf{\Sigma}$ and $\mathbf{C} = \mathbf{V}$. Thus, we avoid SVD as its computational cost is high.

3.2 Randomized Range Finder for Network Embedding

Another solution to Eq. (5) can be obtained by the Randomized Range Finder, which yields an orthonormal matrix \mathbf{C} with very few columns, such that $\|(\mathbf{I} - \mathbf{C}\mathbf{C}^\top)\mathbf{M}\|_F \leq \epsilon$ for a desired tolerance ϵ . Below, we explain the use of the randomized blocked QR factorization to obtain \mathbf{C} .

Suppose that we seek a basis for the range of matrix \mathbf{M} with an exact rank k . Draw a random vector $\boldsymbol{\omega}$, and form the product $\mathbf{y} = \mathbf{M}\boldsymbol{\omega}$. Repeat such a sampling process k times: $\mathbf{y}_i = \mathbf{M}\boldsymbol{\omega}_i, i = 1, 2, \dots, k$. Owing to the randomness, the set $\boldsymbol{\Omega} = [\boldsymbol{\omega}_1, \boldsymbol{\omega}_2, \dots, \boldsymbol{\omega}_k]$ of random vectors is likely to be in the so-called general linear position. In particular, the random vectors form a linearly independent set and no linear combination falls in the null space of \mathbf{M} . As a result, the set $\mathbf{Y} = [\mathbf{y}_1, \dots, \mathbf{y}_k]$ of sample vectors is also linearly independent, so it spans the range of \mathbf{M} . Therefore, to produce an orthonormal basis \mathbf{C} for the range of \mathbf{M} , one just needs to orthonormalize the sample vectors (QR decomposition). The theoretical performance guarantees of randomized QR decomposition based low-rank approximation are given in the Theorem 1.

The randomized QR decomposition works well for matrices whose singular values exhibit some decay, but it may produce a poor basis if the input matrix has a flat spectrum or the input matrix is very large. The power iteration is thus considered in our solution despite requiring extra q times matrix-matrix multiplications ($\mathbf{M}^q\boldsymbol{\Omega}$) as the power iteration is far more accurate if singular values of \mathbf{A} decay slowly. A heuristic we use has a nice property: if the original scheme ($q = 1$) produces a basis whose approximation error is within a factor c' of the optimum [8], the power scheme yields an approximation error within $c'^{1/q}$ of the optimum. The power iteration shrinks the approximation gap towards one exponentially fast.

Theorem 1. *Based on [7], let \mathbf{M} be a matrix ($m \times n$) of real numbers with singular values $\sigma_1 \geq \sigma_2 \geq \sigma_3 \geq \dots$. Choose a target rank $k \geq 2$ and an oversampling parameter $p \geq 2$, where $k + p = l$ and $l \leq \min\{m, n\}$. Draw coefficients of matrix $\boldsymbol{\Omega}$ of size $n \times l$ according to the Normal distribution, $\Omega_{ij} \sim \mathcal{N}(0, 1/k)$, and construct the sample matrix $\mathbf{Y}(\boldsymbol{\Omega}) = \mathbf{M}\boldsymbol{\Omega}$. Then the expected approximation error of $\text{QR}(\mathbf{Y}(\boldsymbol{\Omega})) = \mathbf{C}(\boldsymbol{\Omega})$ yields:*

$$\mathbb{E}_{\boldsymbol{\Omega}} \|\mathbf{M} - \mathbf{C}(\boldsymbol{\Omega})\mathbf{C}(\boldsymbol{\Omega})^\top \mathbf{M}\|_F \leq \min_{k+p=l} \left(1 + \frac{k}{p-1}\right)^{\frac{1}{2}} \left(\sum_{j>k} \sigma_j^2\right)^{\frac{1}{2}}. \quad (6)$$

Although the eigen-decomposition or randomized SVD can be avoided, the orthogonalization (QR decomposition) is the most computationally intensive part of the entire algorithm. Take QR decomposition as an example, the computational cost is approximately $2nk^2 - \frac{2}{3}k^3$ flops. The cost is highly dependent on k . By blocking techniques, one can improve computational efficiency to obtain cost $2nb^2 - \frac{2}{3}qb^3$ flops, where $b \ll k$.

In the proposed by us randomized blocked QR based network embedding, we use power iterations to improve the performance and speed. In the experimental section, we discuss the parameters for our techniques. The pseudocode for our algorithm is outlined in Algorithm 1. Of the parameters of the algorithm, k (target rank) is problem dependent, whereas b (block size) and q (the number of power iterations) are chosen by the user to control the quality and computational cost of the approximation. The algorithm requires the choice of b and q to satisfy $qb \geq k$, and it is not the standard blocked QR as we avoid computing the residual $\mathbf{M}' = \mathbf{M} - \mathbf{R}\mathbf{C}^\top$, which is a dense matrix leading to extra storage costs. Empirically, we found that removing this term does not influence the final performance.

Algorithm 1: Network Embedding based on the Randomized Blocked QR (RBQR).

Input: $\mathbf{M} \in \mathbb{R}^{n \times n}$, rank $k \leq n, b, q$

Output: $\mathbf{R} \in \mathbb{R}^{n \times k}$

for $i \leftarrow 1$ **to** k/b **do**

$\boldsymbol{\Omega} : \Omega_{ij} \sim \mathcal{N}(0, 1)$ ($\boldsymbol{\Omega}$ is of $n \times b$ size)

$\mathbf{C}_i = \text{QR}(\mathbf{M}^q \boldsymbol{\Omega})$

$\mathbf{C}_i = \text{QR}(\mathbf{C}_i - \sum_{j=1}^{i-1} \mathbf{C}_j \mathbf{C}_j^\top \mathbf{C}_i)$

$\mathbf{R}_i^\top = \mathbf{C}_i^\top \mathbf{M}$

end

return $\mathbf{R} = [\mathbf{R}_1, \dots, \mathbf{R}_{k/b}]$

3.3 Spectral Graph Filter for Network Embedding Enhancement

The randomized blocked QR factorization relies on the low-proximity matrix, which means that such a node embedding does not capture the relationship between distant neighbors. In social networks, this is of concern as not direct neighbors are of interest *e.g.*, representing close friends and distant acquaintances.

To improve our model, we capture such relations via the graph diffusion. Intuitively, we apply the heat to the node under consideration and then continuously diffuse the heat towards other neighbors. After a certain time, the heat distribution defines the edge weights from the starting node to other nodes. Thus, we obtain a matrix that defines a new, continuously weighted graph. We define a graph diffusion on node embedding as:

$$\mathbf{R}^* = \sum_{k'=0}^K \theta_{k'} \mathbf{T}^{k'} \mathbf{R}, \quad (7)$$

where $\theta_{k'}$ are coefficients and \mathbf{T} is the transition matrix $\mathbf{D}^{-1}\mathbf{A}$. Coefficients $\theta_{k'}$ are predefined by the specific diffusion variant we choose, such as the heat kernel [9] or Markov diffusion kernel [35]. Increasing K will utilize the information from the K -hop neighborhoods of the node at an increased computational cost. We set $K = 2$ for all datasets, as empirically we observe it is sufficient.

3.4 Computational Complexity

We compare the computational cost of the randomized QR and Algorithm 1. To this end, let $C_{\text{sppm}}, C_{\text{mm}}$ and C_{qr} denote the scaling constants for the cost of sparse matrix-matrix multiplication, matrix-matrix multiplication and a full QR factorization, respectively. The computational complexity for the randomized QR is $(q+1)C_{\text{sppm}}|E|k + C_{\text{qr}}nk^2$, where $|E|$ is the number of edges. Alg. 1 costs $(q+1)C_{\text{sppm}}|E|k + C_{\text{mm}}nk^2 + \frac{2}{k/b}C_{\text{qr}}nk^2$, where b is the block size. The blocked QR is faster if $n/b \geq 2$, whereas more power iterations (larger q) yields better results and costs more. Alg. 1 spends less time executing the full QR factorization, as expected. The computational cost of Eq. (7) is $3|E|k + 3nk$.

Table 1: Micro/Macro-F1(%) of node classification on the Blogcatalog and the PPI datasets.

Metric	ALG	BlogCatalog					PPI				
		10.0%	30.0%	50.0%	70.0%	90.0%	10.0%	30.0%	50.0%	70.0%	90.0%
Micro-F1	LINE	25.35	32.05	35.16	36.61	37.35	11.7	14.2	16	17.82	19.59
	DeepWalk	35.85	39.91	41.62	42.45	42.9	16.06	19.37	21.26	22.63	24.36
	ProNE	36.52	39.97	41.20	41.75	42.16	17.37	22.45	24.28	25.08	26.31
	RBQR	32.88	36.42	37.85	38.61	39.33	16.44	20.90	22.59	23.29	23.71
	REFINE	36.46	39.75	41.00	41.75	42.29	17.79	22.57	24.30	24.96	26.35
Macro-F1	LINE	14.38	19.11	21.36	22.25	22.62	9.5	12.15	13.82	15.35	15.92
	DeepWalk	21.16	25.59	27.58	28.47	28.66	12.89	16.71	18.25	19.48	20.36
	ProNE	18.04	22.68	24.05	24.95	25.03	12.42	17.44	19.59	20.52	20.86
	RBQR	13.74	18.25	20.21	21.30	21.56	11.10	15.69	17.97	18.89	18.71
	REFINE	17.76	22.61	24.17	25.09	25.35	12.67	17.57	19.73	20.60	20.86

4 Experiments

Below, we evaluate our method on the node classification. We study the computational efficiency, the performance and parameters.

4.1 Experiments Setting

We evaluate the algorithms on four widely-used real-world network datasets listed below:

- 1) **BlogCatalog** [30] is a network of social relationships of bloggers in the BlogCatalog website, whose labels represents interests of bloggers.

Table 2: Micro/Macro-F1(%) of node classification on the Wikipedia and Flickr datasets.

Metric	ALG	Wikipedia					Flickr				
		10.0%	30.0%	50.0%	70.0%	90.0%	1.0%	3.0%	5.0%	7.0%	9.0%
Micro-F1	LINE	41.3	48.35	51.89	53.57	54.86	25.3	28.64	30.07	31.28	32.34
	DeepWalk	42.32	47.02	48.65	49.8	50.35	32.06	35.89	37.46	38.29	38.84
	ProNE	48.61	53.96	55.63	56.52	57.43	30.77	35.14	36.69	37.54	38.12
	RBQR	45.72	49.73	50.88	51.56	52.31	30.55	34.23	35.58	36.41	37.04
	REFINE	51.17	56.15	57.50	58.34	58.84	31.17	35.22	36.74	37.66	38.32
Macro-F1	LINE	8.51	10.53	12.63	13.4	13.16	9.01	13.42	15.77	17.44	18.68
	DeepWalk	7.26	9.02	9.71	10.03	9.92	13.36	19.45	22.21	23.94	25.07
	ProNE	8.40	10.73	11.40	11.83	12.21	8.21	13.55	16.14	17.86	19.12
	RBQR	7.10	9.39	9.98	10.36	10.63	7.68	12.21	14.47	15.98	17.21
	REFINE	9.42	11.75	12.36	12.79	13.11	8.18	13.50	16.29	18.21	19.54

- 2) Protein-Protein Interactions (**PPI**) [6] is a subgraph of the PPI network for Homo Sapiens, whose labels represent biological states.
- 3) **Wikipedia** is a co-occurrence network of words appearing in the first million bytes of the Wikipedia dump. The labels are Part-of-Speech (POS) tags inferred using the Stanford POS-Tagger.
- 4) **Flickr** [25] is a network of contacts between Flickr users, whose labels represent the user’s interest group.
- 5) **YouTube** [30] is a social network between YouTube users. The labels represent groups of viewers that enjoy common video genres. The statistics of these datasets are shown in Table 3.

Experiments were conducted on a Ubuntu workstation (AMD Ryzen 2700, 64G RAM).

Baseline Algorithms. We compare the proposed algorithms with LINE(2nd) [24] (number of samples 10, number of negative samples 5, initial learning rate 0.025), DeepWalk [16] (window size 10, walk length 40, and the number of walks 80), and ProNE [31]. The proposed methods include two variants, with or without spectral filters, referred to as RBQR (no spectral filter) and REFINE (RBQR with spectral filters). The dimension of representation is set $k = 128$.

For the node classification, we follow the same experimental setting as DeepWalk. In particular, we randomly sample a portion of labeled nodes for training and use the rest for testing. For BlogCatalog, PPI and Wiki datasets, the training set size ranges from 10% to 90% of a given dataset, with 20% step. For Flickr and YouTube, the training set size ranges from 1% to 9% of the dataset, with 2% step. Using the one-vs-rest logistic regression, we repeat experiment 10 times and report the Micro-F1 and Macro-F1 performance.

4.2 Computational Efficiency

The efficiency of baselines is accelerated by using 16 CPU threads, whereas our approach and ProNE use a single CPU thread. Note that REFINE can be easily extended to a multi-threaded REFINE as QR scales gracefully for multi-threading compared with SVD. As shown in Figure 1a, DeepWalk needs more than 10K seconds to obtain the node representation on YouTube. LINE is $10\times$ faster than DeepWalk while ProNE is $10\times$ faster than LINE. Our approach with/without spectral filters is $10\times$ faster than ProNE. Table 5 reports the running times (both I/O and computation time) of our approach and the fastest baseline, ProNE (other methods are much slower). The runtime results suggest that for PPI and Wiki (small networks of 1,000+ nodes), REFINE requires less than 0.15 seconds to complete while the fastest baseline ProNE is at least $10\times$ slower. Similar speed-ups are consistently observed on BlogCatalog and Flickr, moderate-size networks (10K+ nodes), and YouTube, a relatively big

Table 3: The statistics of datasets.

Dataset	BlogCatalog	Wiki	PPI	Flickr	YouTube
#nodes	10312	4777	3890	80513	1138499
#edges	333983	184812	76584	5899882	2990443
#labels	39	40	50	195	47

Table 4: Micro-F1 scores for node classification. Asterisk ‘*’ indicates that we use 9% of a given dataset for training, other methods use 80% of a given dataset for training.

	HARP	RandNE	MILE	NetSMF	LouvainNE	Ours
Blogcatalog	0.316	0.308	0.264	0.334	0.306	0.420
Flickr	0.384	0.385	0.386	0.356	0.389	0.389*
YouTube	0.305	0.303	0.304	0.307	0.307	0.451*

Table 5: Computational efficiency (runtime in seconds).

Dataset	BlogCatalog	Wiki	PPI	Flickr	YouTube
ProNE	10.04	2.06	0.959	195	305
REFINE	0.49	0.15	0.097	13	40

network (1M+ nodes). Remarkably, our proposed method embeds the YouTube network within 40 seconds by using one thread, whereas ProNE takes 300 seconds. Remaining baselines need between half an hour and a dozen of hours. Without spectral filters, using RBQR alone (Algorithm 1 alone) to obtain the network embedding takes **26** seconds. We analyzed the time cost of our method and found that the random generation costs $\geq 15\%$ of the total time. In practice, the cost may be limited by using pools of random numbers. In this paper, we kept the cost of random generator in our comparisons.

To summarize, the single-threaded REFINE is about 8-20 \times faster than the ProNE, which is 10–400 \times faster than the 16-threaded LINE, DeepWalk, and Node2Vec.

4.3 Performance

For YouTube, Figure 1a shows that REFINE enjoys the best Micro-F1 score and the best runtime compared to three fastest baselines, given training set of 90% size of the dataset. Figure 1b shows that REFINE simultaneously achieves the best Macro-F1 score and the best runtime compared with competing methods. Thus, REFINE without spectral filtering outperforms other more advanced methods such as DeepWalk on YouTube. Tables 1 and 2 summarize the prediction performance on small- and large-scale datasets, respectively, and report the results in terms of Micro-F1 and Macro-F1 metrics. We provide the embedding results generated by the RBQB (Algorithm 1) without spectral filters as well as results with the spectral filter (REFINE). We note that REFINE, RBQR, LINE, and ProNE have all similar performance in Macro-F1 because they are based on the low-proximity matrix factorization. Benefiting from the high-order information, DeepWalk performs well in terms of the Macro-F1 score. Table 1 shows that REFINE works very well in terms of the Micro-F1 scores outperforming DeepWalk by $\sim 8\%$ on Wikipedia, and outperforming LINE for training ratios of 10% and 30% on Macro-F1. REFINE also outperforms the initial RBQR (Algorithm 1) by a large margin. We also conduct the comparison with fast network embedding methods without matrix factorization [1, 13, 3]. As shown in Table 4, our method outperforms other methods significantly. Kindly note that on Flickr and Youtube, our method only uses 9% of a given dataset for the training set, whereas other methods use 80% of a given dataset for training.

4.4 Parameter Analysis

Our approach has two parameters which control the efficiency and effectiveness, that is, power iteration number q and the block size b . As we aforementioned, the power iteration can improve the accuracy of the low-rank approximation at the cost of extra computations. The block size, given a rank k , can determine the number of blocks. Below, we analyze the computational cost of blocked QR, which is highly dependent on the block number.

Power Iteration. Figure 2 shows Micro-F1 w.r.t. the number of power iterations (ranging from 1 to 4, indicated by the four curves) on Blogcatalog. The plot shows that the larger number of iterations is, the better the results of REFINE are. REFINE with $q \geq 2$ significantly improves the performance over without power iteration ($q = 1$). The performance at $q = 3$ saturates, whereas $q = 4$ does not offer any further significant gain.

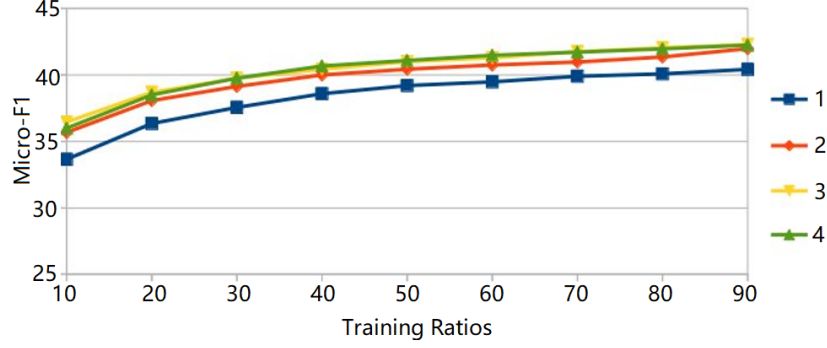


Figure 2: Micro-F1 w.r.t. the number of power iterations on Blogcatalog under different training ratios in range 10%-90%.

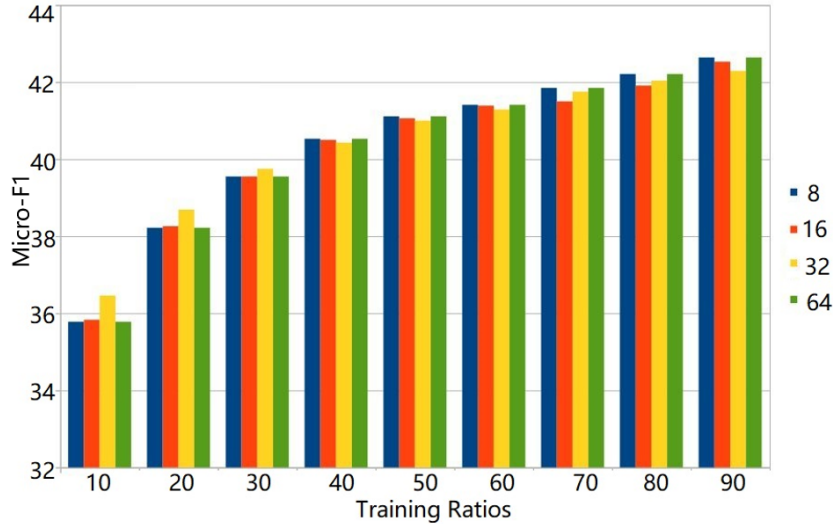


Figure 3: Micro-F1 w.r.t. different block sizes on Blogcatalog under different training ratios (10-90%).

Block Size. Figure 3 shows Macro-F1 w.r.t. different block sizes (8, 16, 32, 64, indicated by different colors) on the Blogcatalog dataset. The plot shows that there are no significant differences between these variants. For the sake of efficiency, the smaller the block size is, the lesser the computational cost is *e.g.*, the QR decomposition may account for less than 10% of runtime of the whole algorithm, which appears to be faster than the cost of random generation.

5 Conclusions

In this work, we have proposed REFINE, a fast and scalable network embedding approach. REFINE achieves a favourable computational efficiency and performance compared to recent powerful network embedding approaches, such as DeepWalk and LINE. The proposed method is 8 – 20 \times faster than the ProNE, which is 10 – 400 \times faster than the aforementioned baselines that are already accelerated by multi-threaded codes. Due to the speed and accuracy of REFINE, its use is suitable for scenarios where node embeddings are frequently updated *e.g.*, in commercial recommendation systems.

Acknowledgement

This research is supported by the Australian Government Research Training Program (RTP) scholarship.

References

- [1] Ayan Kumar Bhowmick, Koushik Meneni, Maximilien Danisch, Jean-Loup Guillaume, and Bivas Mitra. Louvainne: Hierarchical louvain method for high quality and scalable network embedding. In *Proceedings of the 13th International Conference on Web Search and Data Mining*, pages 43–51, 2020.
- [2] Shaosheng Cao, Wei Lu, and Qiongkai Xu. Grarep: Learning graph representations with global structural information. In *Proceedings of the 24th ACM international on conference on information and knowledge management*, pages 891–900, 2015.
- [3] Haochen Chen, Bryan Perozzi, Yifan Hu, and Steven Skiena. Harp: Hierarchical representation learning for networks. *Proceedings of the AAAI Conference on Artificial Intelligence*, 32(1), 2018.
- [4] Siheng Chen, Sufeng Niu, Leman Akoglu, Jelena Kovačević, and Christos Faloutsos. Fast, warped graph embedding: Unifying framework and one-click algorithm. *arXiv preprint arXiv:1702.05764*, 2017.
- [5] Peng Cui, Xiao Wang, Jian Pei, and Wenwu Zhu. A survey on network embedding. *IEEE Transactions on Knowledge and Data Engineering*, 31(5):833–852, 2018.
- [6] 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.
- [7] Nathan Halko, Per-Gunnar Martinsson, and Joel A. Tropp. Finding structure with randomness: Probabilistic algorithms for constructing approximate matrix decompositions, 2010.
- [8] Nathan Halko, Per-Gunnar Martinsson, and Joel A Tropp. Finding structure with randomness: Probabilistic algorithms for constructing approximate matrix decompositions. *SIAM review*, 53(2):217–288, 2011.
- [9] Johannes Klicpera, Stefan Weißenberger, and Stephan Günnemann. Diffusion improves graph learning. *arXiv preprint arXiv:1911.05485*, 2019.
- [10] Piotr Koniusz, Anoop Cherian, and Fatih Porikli. Tensor representations via kernel linearization for action recognition from 3d skeletons. In Bastian Leibe, Jiri Matas, Nicu Sebe, and Max Welling, editors, *Computer Vision - ECCV 2016 - 14th European Conference, Amsterdam, The Netherlands, October 11-14, 2016, Proceedings, Part IV*, volume 9908 of *Lecture Notes in Computer Science*, pages 37–53. Springer, 2016.
- [11] Piotr Koniusz, Lei Wang, and Anoop Cherian. Tensor representations for action recognition. In *IEEE Transactions on Pattern Analysis and Machine Intelligence*. IEEE, 2020.
- [12] Piotr Koniusz and Hongguang Zhang. Power normalizations in fine-grained image, few-shot image and graph classification. In *IEEE Transactions on Pattern Analysis and Machine Intelligence*. IEEE, 2020.
- [13] Jiongqian Liang, Saket Gururkar, and Srinivasan Parthasarathy. Mile: A multi-level framework for scalable graph embedding. *arXiv preprint arXiv:1802.09612*, 2018.
- [14] Tomas Mikolov, Kai Chen, Greg Corrado, and Jeffrey Dean. Efficient estimation of word representations in vector space. *arXiv preprint arXiv:1301.3781*, 2013.
- [15] Mingdong Ou, Peng Cui, Jian Pei, Ziwei Zhang, and Wenwu Zhu. Asymmetric transitivity preserving graph embedding. In *Proceedings of the 22nd ACM SIGKDD international conference on Knowledge discovery and data mining*, pages 1105–1114, 2016.
- [16] 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.
- [17] Arian Prabowo, Piotr Koniusz, Wei Shao, and Flora D. Salim. Coltrane. In *Proceedings of the 6th ACM International Conference on Systems for Energy-Efficient Buildings, Cities, and Transportation*. ACM, 2019.
- [18] Jiezhong Qiu, Yuxiao Dong, Hao Ma, Jian Li, Chi Wang, Kuansan Wang, and Jie Tang. Netsmf: Large-scale network embedding as sparse matrix factorization. In *The World Wide Web Conference*, pages 1509–1520, 2019.

- [19] Jiezhong Qiu, Yuxiao Dong, Hao Ma, Jian Li, Kuansan Wang, and Jie Tang. Network embedding as matrix factorization: Unifying deepwalk, line, pte, and node2vec. In *Proceedings of the Eleventh ACM International Conference on Web Search and Data Mining*, pages 459–467, 2018.
- [20] Wei Shao, Arian Prabowo, Sichen Zhao, Piotr Koniusz, and Flora D. Salim. Predicting flight delay with spatio-temporal trajectory convolutional network and airport situational awareness map. In *Neurocomputing*, 2021.
- [21] Wei Shao, Arian Prabowo, Sichen Zhao, Siyu Tan, Piotr Koniusz, Jeffrey Chan, Xinhong Hei, Bradley Feest, and Flora D. Salim. Flight delay prediction using airport situational awareness map. In *SIGSPATIAL '19*, page 432–435, 2019.
- [22] Ke Sun, Piotr Koniusz, and Zhen Wang. Fisher-bures adversary graph convolutional networks. In Ryan P. Adams and Vibhav Gogate, editors, *Proceedings of The 35th Uncertainty in Artificial Intelligence Conference*, volume 115 of *Proceedings of Machine Learning Research*, pages 465–475. PMLR, 22–25 Jul 2020.
- [23] Jian Tang, Meng Qu, and Qiaozhu Mei. Pte: Predictive text embedding through large-scale heterogeneous text networks. In *Proceedings of the 21th ACM SIGKDD International Conference on Knowledge Discovery and Data Mining*, pages 1165–1174, 2015.
- [24] Jian Tang, Meng Qu, Mingzhe Wang, Ming Zhang, Jun Yan, and Qiaozhu Mei. Line: Large-scale information network embedding. In *Proceedings of the 24th international conference on world wide web*, pages 1067–1077, 2015.
- [25] Lei Tang and Huan Liu. Relational learning via latent social dimensions. In *Proceedings of the 15th ACM SIGKDD international conference on Knowledge discovery and data mining*, pages 817–826, 2009.
- [26] Daixin Wang, Peng Cui, and Wenwu Zhu. Structural deep network embedding. In *Proceedings of the 22nd ACM SIGKDD international conference on Knowledge discovery and data mining*, pages 1225–1234, 2016.
- [27] Xianjing Wang, Flora D. Salim, Yongli Ren, and Piotr Koniusz. Relation embedding for personalised translation-based poi recommendation. *Lecture Notes in Computer Science*, page 53–64, 2020.
- [28] Xiao Wang, Peng Cui, Jing Wang, Jian Pei, Wenwu Zhu, and Shiqiang Yang. Community preserving network embedding. In *Thirty-first AAAI conference on artificial intelligence*, 2017.
- [29] Cheng Yang, Maosong Sun, Zhiyuan Liu, and Cunchao Tu. Fast network embedding enhancement via high order proximity approximation. In *IJCAI*, pages 3894–3900, 2017.
- [30] Reza Zafarani and Huan Liu. Social computing data repository at asu, 2009.
- [31] Jie Zhang, Yuxiao Dong, Yan Wang, Jie Tang, and Ming Ding. Prone: Fast and scalable network representation learning. In *IJCAI*, pages 4278–4284, 2019.
- [32] Ziwei Zhang, Peng Cui, Haoyang Li, Xiao Wang, and Wenwu Zhu. Billion-scale network embedding with iterative random projection. In *2018 IEEE International Conference on Data Mining (ICDM)*, pages 787–796. IEEE, 2018.
- [33] Ziwei Zhang, Peng Cui, Xiao Wang, Jian Pei, Xuanrong Yao, and Wenwu Zhu. Arbitrary-order proximity preserved network embedding. In *Proceedings of the 24th ACM SIGKDD International Conference on Knowledge Discovery & Data Mining*, pages 2778–2786, 2018.
- [34] Hao Zhu and Piotr Koniusz. Graph convolutional network with generalized factorized bilinear aggregation, 2021.
- [35] Hao Zhu and Piotr Koniusz. Simple spectral graph convolution. In *International Conference on Learning Representations*, 2021.