

Advances in Collaborative Filtering and Ranking

By

LIWEI WU

B.S. City University of Hong Kong 2014

M.S. University of California, Davis 2016 & 2018

DISSERTATION

Submitted in partial satisfaction of the requirements for the degree of

DOCTOR OF PHILOSOPHY

in

Statistics

in the

OFFICE OF GRADUATE STUDIES

of the

UNIVERSITY OF CALIFORNIA

DAVIS

Approved:

James Sharpnack, Chair

Cho-Jui Hsieh, Chair

Thomas Lee

Committee in Charge

2020

Copyright © 2020 by

Liwei Wu

All rights reserved.

*To my loving parents and
adorable girlfriend . . .*

CONTENTS

List of Figures	vii
List of Tables	ix
Abstract	xiii
Acknowledgments	xiv
1 Introduction	1
1.1 Overview	1
1.2 Contributions and Outline of This Thesis	3
2 Collaborative Filtering with Graph Encoding	9
2.1 Introduction	9
2.2 Related Work	10
2.3 Methodology	11
2.3.1 Bloom Filter	12
2.3.2 Graph DNA Encoding Via Bloom Filters	13
2.4 Collaborative Filtering with Graph DNA	14
2.4.1 Graph Regularized Matrix Factorization	15
2.4.2 Co-Factorization with Graph Information	16
2.4.3 Graph Convolutional Matrix Completion	16
2.5 Experiments	17
2.5.1 Simulation Study	17
2.5.2 Graph Regularized Matrix Factorization for Explicit Feedback	19
2.5.3 Co-Factorization with Graph for Explicit Feedback	19
2.5.4 Graph Regularized Weighted Matrix Factorization for Implicit Feedback	22
2.5.5 Graph Convolutional Matrix Factorization	22
2.6 Conclusion	24

3	Large-scale Pairwise Collaborative Ranking in Near-Linear Time	25
3.1	Introduction	25
3.2	Related Work	27
3.3	Problem Formulation	28
3.4	Proposed Algorithms	30
3.4.1	Motivation and Overview	30
3.4.2	Primal-CR: the proposed algorithm for pairwise input data	31
3.4.3	Primal-CR++: the proposed algorithm for rating data	40
3.4.4	Parallelization	42
3.5	Experiments	43
3.5.1	Compare single thread versions using the same subsamples	45
3.5.2	Compare parallel versions	47
3.5.3	Performance using Full Training Data	47
3.6	Conclusions	48
4	SQL-Rank: A Listwise Approach to Collaborative Ranking	50
4.1	Introduction	50
4.1.1	A brief history of collaborative ranking	50
4.1.2	Contributions of this work	52
4.2	Methodology	53
4.2.1	Permutation probability	53
4.2.2	Deriving objective function for SQL-Rank	55
4.2.3	Non-convex implementation	57
4.3	Experiments	59
4.3.1	Implicit Feedback	59
4.3.2	Explicit Feedback	62
4.3.3	Training speed	64
4.3.4	Effectiveness of Stochastic Queuing (SQ)	64
4.3.5	Effectiveness of using the Full List	65
4.4	Conclusions	66

5	Stochastic Shared Embeddings: Data-driven Regularization of Embedding Layers	67
5.1	Introduction	67
5.2	Related Work	68
5.3	Stochastic Shared Embeddings	69
5.3.1	General SSE with Knowledge Graphs: SSE-Graph	72
5.3.2	Simplified SSE with Complete Graph: SSE-SE	73
5.3.3	Theoretical Guarantees	74
5.4	Experiments	78
5.4.1	Neural Networks with One Hidden Layer (Matrix Factorization and BPR)	78
5.4.2	Transformer Encoder Model for Sequential Recommendation	79
5.4.3	Neural Machine Translation	80
5.4.4	BERT for Sentiment Classification	80
5.4.5	Speed and convergence comparisons.	82
5.5	Conclusion	82
6	SSE-PT: Sequential Recommendation Via Personalized Transformer	83
6.1	Introduction	83
6.2	Related Work	84
6.2.1	Session-based and Sequential Recommendation	84
6.2.2	Regularization Techniques	85
6.3	Methodology	85
6.3.1	Sequential Recommendation	85
6.3.2	Personalized Transformer Architecture	86
6.3.3	Handling Long Sequences: SSE-PT++	88
6.4	Experiments	89
6.4.1	Attention Mechanism Visualization	94
6.4.2	Training Speed	94
6.4.3	Ablation Study	95

6.5	Conclusion	96
7	Conclusions	97
7.1	Summary of Contributions	97
7.2	Future Work	98
8	Appendix	100
8.1	Appendix to Chapter 2	100
8.1.1	Simulation Study	100
8.1.2	Metrics	101
8.1.3	Graph Regularized Weighted Matrix Factorization for Implicit feed-back	103
8.1.4	Reproducibility	103
8.1.5	Code	105
8.2	Appendix to Chapter 4	105
8.3	Appendix to Chapter 5	107
8.3.1	Neural Networks with One Hidden Layer	107
8.3.2	Neural Machine Translation	110
8.3.3	BERT	112
8.3.4	Proofs	113
8.4	Appendix to Chapter 6	117

LIST OF FIGURES

2.1	Illustration of Algorithm 1: the graph DNA encoding procedure. The curly brackets at each node indicate the nodes encoded at a particular step. At $d = 0$ each node’s Bloom filter only encodes itself, and multi-hop neighbors are included as d increases.	14
2.2	Illustration of our proposed DNA encoding method (DNA-3), with the corresponding bipartite graph representation.	15
2.3	Compare Training Speed of GRMF, with and without Graph DNA.	22
3.1	Comparing Primal-CR, Primal-CR++ and Collrank, MovieLens1m data, 200 ratings/user, rank 100, lambda = 5000	41
3.2	Comparing Primal-CR, Primal-CR++ and Collrank, MovieLens10m data, 500 ratings/user, rank 100, lambda = 7000	42
3.3	Comparing Primal-CR, Primal-CR++ and Collrank, Netflix data, 200 ratings/user, rank 100, lambda = 10000	42
3.4	Comparing parallel version of Primal-CR and Collrank, MovieLens10m data, 500 ratings/user, rank 100, lambda = 7000	43
3.5	Speedup of Primal-CR, MovieLens10m data, 500 ratings/user, rank 100, lambda = 7000	43
3.6	Varies number of ratings per user in training data, MovieLens1m data, rank 100	46
3.7	Varies number of ratings per user in training data, Netflix data, rank 100	47
4.1	Demonstration of Stochastic Queuing Process—the rating matrix R (left) generates multiple possible rankings Π ’s (right), $\Pi \in \mathcal{S}(R, \Omega)$ by breaking ties randomly.	55
4.2	Training time of implicit feedback methods.	64

5.1	SSE-Graph described in Algorithm 9 and Figure 5.2 can be viewed as adding exponentially many distinct reordering layers above the embedding layer. A modified backpropagation procedure in Algorithm 9 is used to train exponentially many such neural networks at the same time.	71
5.2	Illustration of how SSE-Graph algorithm in Figure 5.1 works for a simple neural network.	71
5.3	Projecting 50-dimensional embeddings obtained by training a simple neural network without SSE (Left), and with SSE-Graph (Center) , SSE-SE (Right) into 3D space using PCA.	74
5.4	Compare Training Speed of Simple Neural Networks with One Hidden Layer, i.e. MF and BPR, with and without SSE-SE.	82
6.1	Illustration of our proposed SSE-PT model	87
6.2	Illustration of how SASRec (Left) and SSE-PT (Right) differs on utilizing the Engagement History of A Random User in Movielens1M Dataset. . .	93
6.3	Illustration of the speed of SSE-PT	95
8.1	Compare Training Speed of GRMF, with and without Graph DNA.	101
8.2	Comparing implicit feedback methods.	110
8.3	Effectiveness of Stochastic Queuing Process.	111
8.4	Effectiveness of using full lists.	112
8.5	Simulation of a bound on $\rho_{L,n}$ for the movielens1M dataset. Throughout the simulation, L is replaced with ℓ (which will bound $\rho_{L,n}$ by Jensen’s inequality). The SSE probability parameter dictates the probability of transitioning. When this is 0 (box plot on the right), the distribution is that of the samples from the standard Rademacher complexity (without the sup and expectation). As we increase the transition probability, the values for $\rho_{L,n}$ get smaller.	114

LIST OF TABLES

1.1	Summary of Different Fundamental Approaches Before This Dissertation.	3
1.2	Summary of Different Fundamental Approaches After This Dissertation.	3
2.1	Comparison of Graph Regularized Matrix Factorization Variants for Explicit Feedback on Synthetic, Douban and Flixster data. We use rank $r = 10$. RGG is the Relative Graph Gain defined in (2.3).	20
2.2	Graph DNA (Algorithm 1) Encoding Speed. We set number $c = 500$ and implement Graph DNA using single-core python. We can scale up linearly in terms of depth d for a fixed c	20
2.3	Comparison of GRWMF Variants for Implicit Feedback on Douban and Flixster datasets. P stands for precision and N stands for NDCG. We use rank $r = 10$ and all results are in %.	20
2.4	Comparison of GCN Methods for Explicit Feedback on Douban, Flixster and Yahoo Music datasets (3000 by 3000 as in [8, 74]). All the methods except GC-MC utilize side graph information.	21
2.5	Comparison of GRMF Methods of different ranks for Explicit Feedback on Flixster Dataset.	24
3.1	Datasets used for experiments	46
3.2	Scability of Primal-CR and Collrank on Movielens10m	46
4.1	Comparing implicit feedback methods on various datasets.	61
4.2	Comparing explicit feedback methods on various datasets.	63
4.3	Effectiveness of Stochastic Queuing Process.	65
4.4	Comparing different k on Movielens1m data set using 50 training data per user.	65

5.1	Compare SSE-Graph and SSE-SE against ALS-MF with Graph Laplacian Regularization. The p_u and p_i are the SSE probabilities for user and item embedding tables respectively, as in (5.5). Definitions of ρ_u and ρ_i can be found in (5.3). Movielens10m does not have user graphs.	77
5.2	SSE-SE outperforms Dropout for Neural Networks with One Hidden Layer such as Matrix Factorization Algorithm regardless of dimensionality we use. p_s is the SSE probability for both user and item embedding tables and p_d is the dropout probability.	77
5.3	SSE-SE outperforms dropout for Neural Networks with One Hidden Layer such as Bayesian Personalized Ranking Algorithm regardless of dimensionality we use. We report the metric precision for top k recommendations as $P@k$	77
5.4	SSE-SE has two tuning parameters: probability p_x to replace embeddings associated with input x_i and probability p_y to replace embeddings associated with output y_i . We use the dropout probability of 0.1, weight decay of $1e^{-5}$, and learning rate of $1e^{-3}$ for all experiments.	79
5.5	Our proposed SSE-SE helps the Transformer achieve better BLEU scores on English-to-German in 10 out of 11 newstest data between 2008 and 2018.	80
5.6	Our proposed SSE-SE applied in the pre-training stage on our crawled IMDB data improves the generalization ability of pre-trained IMDB model and helps the BERT-Base model outperform current SOTA results on the IMDB Sentiment Task after fine-tuning.	81
5.7	SSE-SE pre-trained BERT-Base models on IMDB datasets turn out working better on the new unseen SST-2 Task as well.	82

6.1	Comparing various state-of-the-art temporal collaborative ranking algorithms on various datasets. The (A) to (E) are non-deep-learning methods, the (F) to (K) are deep-learning methods and the (L) to (O) are our variants. We did not report SSE-PT++ results for beauty, games and steam, as the input sequence lengths are very short (see Table 8.7), so there is no need for SSE-PT++.	91
6.2	Comparing SASRec, SSE-PT and SSE-PT++ on Movielens1M Dataset while varying the maximum length allowed and dimension of embeddings.	92
6.3	Comparing Different Regularizations for SSE-PT on Movielens1M Dataset. NO REG stands for no regularization. PS stands for parameter sharing across all users while PS(AGE) means PS is used within each age group. SASRec is added to last row after all SSE-PT results as a baseline. . . .	92
8.1	Compare Bloom filters of different depths and sizes on Synthesis Dataset. Note that the number of bits of Bloom filter is decided by Bloom filter’s maximum capacity and tolerable error rate (i.e. false positive error, we use 0.2 as default).	102
8.2	Compare nnz of different methods on Douban and Flixster datasets. GRMF_ G^4 and GRMF_DNA-2 are using the same 4-hop information in the graph but in different ways. Note that we do not exclude potential overlapping among columns.	103
8.3	Compare Matrix Factorization for Explicit Feedback on Synthesis Dataset. The synthesis dataset has 10,000 users and 2,000 items with user friendship graph of size $10,000 \times 10,000$. Note that the graph only contains at most 6-hop valid information. GRMF_ G^6 means GRMF with $G + \alpha \cdot G^2 + \beta \cdot G^3 + \gamma \cdot G^4 + \epsilon \cdot G^5 + \omega \cdot G^6$. GRMF_DNA- d means depth d is used. . . .	106
8.4	Compare Matrix Factorization methods for Explicit Feedback on Douban and Flixster data. We use rank $r = 10$	106
8.5	Compare Co-factor Methods for Explicit Feedback on Douban and Flixster Datasets. We use rank $r = 10$ for both methods.	107

8.6	Compare Weighted Matrix Factorization with Graph for Implicit Feedback on Douban and Flixster Datasets. We use rank $r = 10$ for both methods and all metric results are in %.	107
8.7	Description of Datasets Used in Evaluations.	118
8.8	Comparing our SSE-PT, SSE-PT++ with SASRec on Movielen1M dataset. We use number of negatives $C = 100$, dropout probability of 0.2 and learning rate of $1e^{-3}$ for all experiments while varying others. p_u, p_i, p_o are SSE probabilities for user embedding, input item embedding and output item embedding respectively.	119
8.9	Comparing our SSE-PT with SASRec on Movielens10M dataset. Unlike Table 8.8, we use the number of negatives $C = 500$ instead of 100 as $C = 100$ is too easy for this dataset and it gets too difficult to tell the differences between different methods: Hit Ratio@10 approaches 1.	120
8.10	Comparing Different SSE probability for user embeddings for SSE-PT on Movielens1M Dataset. Embedding hidden units of 50 for users and 100 for items, attention blocks of 2, SSE probability of 0.01 for item embeddings, dropout probability of 0.2 and max length of 200 are used.	121
8.11	Comparing Different Sampling Probability, p_s , of SSE-PT++ on Movielens1M Dataset. Hyper-parameters the same as Table 8.10, except that the max length T allowed is set 100 instead of 200 to show effects of sampling sequences.	121
8.12	Comparing Different Number of Blocks for SSE-PT while Keeping The Rest Fixed on Movielens1M and Movielens10M Datasets.	122
8.13	Varying number of negatives C in evaluation on Movielens1M dataset. Other hyper-parameters are fixed for a fair comparison.	122

ABSTRACT

Advances in Collaborative Filtering and Ranking

In this dissertation, we cover some recent advances in collaborative filtering and ranking. In chapter 1, we give a brief introduction of the history and the current landscape of collaborative filtering and ranking; chapter 2 we first talk about pointwise collaborative filtering problem with graph information, and how our proposed new method can encode very deep graph information which helps four existing graph collaborative filtering algorithms; chapter 3 is on the pairwise approach for collaborative ranking and how we speed up the algorithm to near-linear time complexity; chapter 4 is on the new listwise approach for collaborative ranking and how the listwise approach is a better choice of loss for both explicit and implicit feedback over pointwise and pairwise loss; chapter 5 is about the new regularization technique Stochastic Shared Embeddings (SSE) we proposed for embedding layers and how it is both theoretically sound and empirically effectively for 6 different tasks across recommendation and natural language processing; chapter 6 is how we introduce personalization for the state-of-the-art sequential recommendation model with the help of SSE, which plays an important role in preventing our personalized model from overfitting to the training data; chapter 7, we summarize what we have achieved so far and predict what the future directions can be; chapter 8 is the appendix to all the chapters.

ACKNOWLEDGMENTS

Being a PhD student at the Statistics Department of the University of California, Davis is an extremely fun experience. Choosing Davis among many other schools back in 2014 is definitely one of the best decisions I have ever made in my life. Choosing my two advisors is another great decision I made in my life: at Davis, I was fortunate to meet my 2 wonderful advisors Cho (who has now moved to UCLA) and James, who have both become fathers in recent years. It is an interesting fact that they both came to Davis one year later than me, so technically one can argue I am more senior at Davis than my advisors. During my study, I was given a lot of freedom to follow my heart and passion to explore diverse directions that I am most curious about and work on my own projects with lots of encouragements, and appropriate level of guidance whenever I need it. I have enjoyed every single bit of the past five and half years of my study.

In the process, I was fortunate to obtain 2 master degrees in Statistics and Computer Science in 2016 and 2018 respectively. I got to take or sit-in courses from many wonderful professors from both departments, including Ethan, Alexander, Prabir, Fushing, Jiming, Thomas, Hans, Debashis, James, Jie, Wolfgang, Duncan, and Jane-Ling from Statistics Department, Cho, Ian, Vladimir, Yong Jae, Yu, Matthew, Nitta, Sam from Computer Science Department.

I was blessed to do 4 distinct internships at 4 different companies, namely AT&T Labs at San Ramon, Facebook at Melno Park and LinkedIn at Mountain View, and Dataminr at midtown Manhattan, during which I have encountered many helpful and inspirational people, especially my colleagues, mentors and supervisors. Moreover, I got the precious opportunity to present my work at top conferences including KDD'17 at Halifax, ICML'18 at Stockholm and NeurIPS'19 at Vancouver and attend different conferences around the world. I also got to teach a lot at Davis: I would have taught as a teaching assistant for 16 quarters.

I'm truly grateful for the professors that wrote me strong recommendation letters while I was in college. Without the help and encouragements, my journey of PhD study would not have started. Thank you, Betsy, Jun, Rick and Xiang. It is an interesting fact that

except Rick has retired, all the rest have now got tenure. Congratulations and how time flies!

I also want to thank my collaborators including Hsiang-Fu and Nikhil from Amazon, Mahdi from Rutgers, Shengli, Joel and Alex from Dataminr, and of course most importantly my girlfriend, collaborator and best friend Shuqing. It has been a great pleasure to work with Shuqing, travel together to different conferences, and present our works. She also helped to proof-read the paper. Her accompany during past few years made my PhD study enjoyable. Actually SQL-Rank algorithm presented in chapter 4 is also named after her.

Last but not the least, I want to mourn for those over 1000 people who have passed away because of the novel coronavirus outbreak in China (including Dr Wenliang Li and other heroes I don't know their names), and pray for families still stuck in Wuhan, including my parents and grandmother. Born and raised up in Wuhan till age of 18, I sincerely wish we can recover from coronavirus soon together as humankind in 2020. I cannot imagine how much pain those who lost their parents, partner and loved ones have been bearing. Maybe this reminds us to spend every day just like the last day on earth. Ask yourself, if today were the last day, what is the most important thing that you want to do? Then just follow your heart and I believe everyone can achieve a meaningful life.

Chapter 1

Introduction

1.1 Overview

Nowadays in online retail and online content delivery applications, it is commonplace to have embedded recommendation systems algorithms that recommend items to users based on previous user behaviors and ratings. The field of recommender systems has gained more and more popularity ever since the famous Netflix competition [7], in which competitors utilize user ratings to predict ratings for each user-movie pair and the final winner takes home 1 million dollars. During the competition, 2 distinct approaches stand out: one being the Restricted Boltzmann Machines [93] and the other being matrix factorization [61, 73]. The combination of both approaches work well during the competition, but due to the ease of training and inference, matrix factorization approaches have dominated the collaborative filtering field before the widespread adoption of deep learning methods [95]. Collaborative filtering refers to making automatic predictions (filtering) about the interests of a user by collecting preferences or taste information from many users (collaborating). Usually this does not require approaching the recommendation problem as a ranking problem but rather pretends it is a regression or classification problem. Of course, there may be some loss incurred when using a regression or classification loss for ultimately what is a ranking problem. Because at the end of day, the ordering is the most important thing, which is directly associated with the recommender system performance. Collaborative ranking approaches mitigate the concerns by using a ranking loss. The ranking loss can be either

pairwise or listwise. The difference among pointwise, pairwise and listwise are rooted in the distinct interpretations of the same data points. Pointwise approaches assume each user-item rating datum is independent; pairwise approaches assume that pairwise comparisons for two items by the same user are independent; the listwise approaches view the list of item preferences as a whole and treats different users' list as independent data points. In a strict definition, the Collaborative Filtering approaches refers to the pointwise approach while the Collaborative Ranking refers to the pairwise and listwise approaches. In a loose definition, Collaborative Ranking sometimes is viewed as a sub-field of Collaborative Filtering.

In another classification of these approaches, recommender systems can be divided into those designed for explicit feedback, such as ratings [61], and those for implicit feedback, based on user engagement [49]. Recently, implicit feedback datasets, such as user clicks of web pages, check-in's of restaurants, likes of posts, listening behavior of music, watching history and purchase history, are increasingly prevalent. Unlike explicit feedback, implicit feedback datasets can be obtained without users noticing or active participation. All collaborative filtering and collaborative ranking approaches can be divided into these 6 (3 by 2) categories.

So far, we have not touched the features used in recommender systems. Most open datasets do not have good user-side features due to privacy concerns. Even for item side, most datasets including Netflix, Movielens datasets do not have features prepared. However, feature information is crucial for better recommendation and ranking performances besides the fundamental approaches and data. There are many useful features but among them, probably the most challenging ones are including graphs that encode item or user relationships [8, 89] and temporal information [43, 56]. The ways of constructing graphs can vary. One way to define graph over users is to exploit the friendship relationships among friends. But there are many ways to define such graphs. If there is more than one type of relationship, then we call it knowledge graphs instead of graphs. Knowledge graphs widely exist: for example, between 2 movies, they could be starred by the same actors or they could be the same genre of movies. In the field of collaborative filtering

Table 1.1. Summary of Different Fundamental Approaches Before This Dissertation.

	Explicit Feedback	Implicit Feedback
Point-wise Approach	[61, 73]	[49, 78]
Pair-wise Approach	[81]	[91]
List-wise Approach	[50, 121]	

Table 1.2. Summary of Different Fundamental Approaches After This Dissertation.

	Explicit Feedback	Implicit Feedback
Point-wise Approach	[61, 73]	[49, 78]
Pair-wise Approach	[81] [115]	[91]
List-wise Approach	[50, 121]	[116]

with graphs, [8, 89] find using graph indeed improves recommendation performances. As to temporal information, it is motivated by the fact that user-item interactions do not happen at one time and then stay static. Instead, usually they occur in a temporal order. This means standard train/validation/test splits may not be realistic. Because during inference, we want to predict future interactions only based on historical interactions. In sequential recommendation [56], train/validation/test are split in temporal ordering so they do not overlap. It has been shown in such setting, temporal ordering plays a large role in final ranking results [43, 56].

1.2 Contributions and Outline of This Thesis

This dissertation summarizes several published and under-review works that advance the field of collaborative filtering and ranking. Our first main contribution is that we fill out the void in Table 1.1. Before this dissertation (the chapter 4 [116]), no one had successfully applied the listwise approach to the implicit feedback setting due to the difficulty of the problem, because in implicit feedback contain only 1’s and 0’s without different level of ratings. A second contribution is that we speed up the previous best pairwise approaches for explicit feedback significantly. We achieved near-linear time complexity, which allows us to scale up to the full Netflix dataset without subsampling [115]. We not only contribute to the fundamental approaches but also to extended approaches that utilize extra information, including graph and temporal ordering information. On one hand, we propose a novel

way to encode long range graph interactions without require any training using bloom filters as backbone [119]. On the other hand, with the help of a new embedding-layer regularization called Stochastic Shared Embeddings (SSE) [117], we can also introduce personalization for the state-of-the-art sequential recommendation model and achieve much better ranking performance with our personalized model [118], where personalization is crucial for the success of recommender systems unlike most natural language tasks. This new regularization not only helps existing collaborative filtering and collaborative ranking algorithms but also benefits methods in natural language processing in fields like machine translation and sentiment analysis [117].

The outline of this thesis is as follows: chapter 2 we first talk about pointwise collaborative filtering problem with graph information, and how our proposed new method can encode very deep graph information which helps four existing graph collaborative filtering algorithms; chapter 3 is on the pairwise approach for collaborative ranking and how we speed up the algorithm to near-linear time complexity; chapter 4 is on the new listwise approach for collaborative ranking and how the listwise approach is a better choice of loss for both explicit and implicit feedback over pointwise and pairwise loss; chapter 5 is about the new regularization technique Stochastic Shared Embeddings (SSE) we proposed for embedding layers and how it is both theoretically sound and empirically effectively for 6 different tasks across recommendation and natural language processing; chapter 6 is how we introduce personalization for the state-of-the-art sequential recommendation model with the help of SSE, which plays an important role in preventing our personalized model from overfitting to the training data.

Summary Chapter 2:

In this chapter, we consider recommender systems with side information in the form of graphs. Existing collaborative filtering algorithms mainly utilize only immediate neighborhood information and do not efficiently take advantage of deeper neighborhoods beyond 1-2 hops. The main issue with exploiting deeper graph information is the rapidly growing time and space complexity when incorporating information from these neighborhoods. In this chapter, we propose using Graph DNA, a novel Deep Neighborhood Aware graph

encoding algorithm, for exploiting multi-hop neighborhood information. DNA encoding computes approximate deep neighborhood information in linear time using Bloom filters, and results in a per-node encoding whose dimension is logarithmic in the number of nodes in the graph. It can be used in conjunction with both feature-based and graph-regularization-based collaborative filtering algorithms. Graph DNA has the advantages of being memory and time efficient and providing additional regularization when compared to directly using higher order graph information. Code is open-sourced at <https://github.com/wuliwei9278/Graph-DNA>. This work is going to be published at the 23rd International Conference on Artificial Intelligence and Statistics (AISTATS 2020).

Summary Chapter 3:

In this chapter, we consider the Collaborative Ranking (CR) problem for recommendation systems. Given a set of pairwise preferences between items for each user, collaborative ranking can be used to rank un-rated items for each user, and this ranking can be naturally used for recommendation. It is observed that collaborative ranking algorithms usually achieve better performance since they directly minimize the ranking loss; however, they are rarely used in practice due to the poor scalability. All the existing CR algorithms have time complexity at least $O(|\Omega|r)$ per iteration, where r is the target rank and $|\Omega|$ is number of pairs which grows quadratically with number of ratings per user. For example, the Netflix data contains totally 20 billion rating pairs, and at this scale all the current algorithms have to work with significant subsampling, resulting in poor prediction on testing data.

In this chapter, we propose a new collaborative ranking algorithm called Primal-CR that reduces the time complexity to $O(|\Omega| + d_1\bar{d}_2r)$, where d_1 is number of users and \bar{d}_2 is the averaged number of items rated by a user. Note that $d_1\bar{d}_2$ is strictly smaller and often much smaller than $|\Omega|$.

Furthermore, by exploiting the fact that most data is in the form of numerical ratings instead of pairwise comparisons, we propose Primal-CR++ with $O(d_1\bar{d}_2(r + \log \bar{d}_2))$ time complexity. Both algorithms have better theoretical time complexity than existing approaches and also outperform existing approaches in terms of NDCG and pairwise

error on real data sets. To the best of our knowledge, this is the first collaborative ranking algorithm capable of working on the full Netflix dataset using all the 20 billion rating pairs, and this leads to a model with much better recommendation compared with previous models trained on subsamples. Finally, compared with classical matrix factorization algorithm which also requires $O(d_1 \bar{d}_2 r)$ time, our algorithm has almost the same efficiency while making much better recommendations since we consider the ranking loss. Code is open-sourced at <https://github.com/wuliwei9278/ml-1m> (Julia version) and <https://github.com/wuliwei9278/primalCR> (C++ version). This work has been published at the 23rd ACM SIGKDD International Conference on Knowledge Discovery and Data Mining (KDD 2017).

Summary Chapter 4:

In this chapter, we propose a listwise approach for constructing user-specific rankings in recommendation systems in a collaborative fashion. We contrast the listwise approach to previous pointwise and pairwise approaches, which are based on treating either each rating or each pairwise comparison as an independent instance respectively. By extending the work of [15], we cast listwise collaborative ranking as maximum likelihood under a permutation model which applies probability mass to permutations based on a low rank latent score matrix. We present a novel algorithm called SQL-Rank, which can accommodate ties and missing data and can run in linear time. We develop a theoretical framework for analyzing listwise ranking methods based on a novel representation theory for the permutation model. Applying this framework to collaborative ranking, we derive asymptotic statistical rates as the number of users and items grow together. We conclude by demonstrating that our SQL-Rank method often outperforms current state-of-the-art algorithms for implicit feedback such as Weighted-MF and BPR and achieve favorable results when compared to explicit feedback algorithms such as matrix factorization and collaborative ranking. Code is open-sourced at <https://github.com/wuliwei9278/SQL-Rank>. This work has been published at the Thirty-fifth International Conference on Machine Learning (ICML 2018).

Summary Chapter 5:

In deep neural nets, lower level embedding layers account for a large portion of the total number of parameters. Tikhonov regularization, graph-based regularization, and hard parameter sharing are approaches that introduce explicit biases into training in a hope to reduce statistical complexity. Alternatively, we propose stochastic shared embeddings (SSE), a data-driven approach to regularizing embedding layers, which stochastically transitions between embeddings during stochastic gradient descent (SGD). Because SSE integrates seamlessly with existing SGD algorithms, it can be used with only minor modifications when training large scale neural networks. We develop two versions of SSE: SSE-Graph using knowledge graphs of embeddings; SSE-SE using no prior information. We provide theoretical guarantees for our method and show its empirical effectiveness on 6 distinct tasks, from simple neural networks with one hidden layer in recommender systems, to the transformer and BERT in natural languages. We find that when used along with widely-used regularization methods such as weight decay and dropout, our proposed SSE can further reduce overfitting, which often leads to more favorable generalization results. Code is open-sourced at <https://github.com/wuliwei9278/SSE>. This work has been published at the Thirty-third Annual Conference on Neural Information Processing Systems (NeurIPS 2019).

Summary Chapter 6:

Temporal information is crucial for recommendation problems because user preferences are naturally dynamic in the real world. Recent advances in deep learning, especially the discovery of various attention mechanisms and newer architectures in addition to widely used RNN and CNN in natural language processing, have allowed for better use of the temporal ordering of items that each user has engaged with. In particular, the SASRec model, inspired by the popular Transformer model in natural languages processing, has achieved state-of-the-art results. However, SASRec, just like the original Transformer model, is inherently an un-personalized model and does not include personalized user embeddings. To overcome this limitation, we propose a Personalized Transformer (SSE-PT) model, outperforming SASRec by almost 5% in terms of NDCG@10 on 5 real-world

datasets. Furthermore, after examining some random users' engagement history, we find our model not only more interpretable but also able to focus on recent engagement patterns for each user. Moreover, our SSE-PT model with a slight modification, which we call SSE-PT++, can handle extremely long sequences and outperform SASRec in ranking results with comparable training speed, striking a balance between performance and speed requirements. Our novel application of the Stochastic Shared Embeddings (SSE) regularization is essential to the success of personalization. Code and data are open-sourced at <https://github.com/wuliwei9278/SSE-PT>. This work is currently still under review.

Chapter 2

Collaborative Filtering with Graph Encoding

2.1 Introduction

Recommendation systems are increasingly prevalent due to content delivery platforms, e-commerce websites, and mobile apps [98]. Classical collaborative filtering algorithms use matrix factorization to identify latent features that describe the user preferences and item meta-topics from partially observed ratings [61]. In addition to rating information, many real-world recommendation datasets also have a wealth of side information in the form of graphs, and incorporating this information often leads to performance gains [67, 89, 128]. However, each of these only utilizes the immediate neighborhood information of each node in the side information graph. More recently, [8] incorporated graph information when learning features with a Graph Convolution Network (GCN) based recommendation algorithm. GCNs [58] constitute flexible methods for incorporating graph structure beyond first-order neighborhoods, but their training complexity typically scales rapidly with the depth, even with sub-sampling techniques [18]. Intuitively, exploiting higher-order neighborhood information could benefit the generalization performance, especially when the graph is sparse, which is usually the case in practice. The main caveat of exploiting higher-order graph information is the high computational and memory cost when computing higher-order neighbors since the number of t -hop neighbors typically grows exponentially with t .

We aim to utilize higher order graph information without introducing much computational and memory overhead. We propose a Graph Deep Neighborhood Aware (Graph DNA) encoding, which approximately captures the higher-order neighborhood information of each node via Bloom filters [9]. Bloom filters encode neighborhood sets as c dimensional 0/1 vectors, where $c = O(\log n)$ for a graph with n nodes, which approximately preserves membership information. This encoding can then be combined with both graph regularized or feature based collaborative filtering algorithms, with little computational and memory overhead. In addition to computational speedups, we find that Graph DNA achieves better performance over competitors. We show that our Graph DNA encoding can be used with several collaborative filtering algorithms: graph-regularized matrix factorization with explicit and implicit feedback [89, 128], co-factoring [67], and GCN-based recommendation systems [74]. In some cases, using information from deeper neighborhoods (like 4th order) yields a 15x increase in performance, with graph DNA encoding yielding a 6x speedup compared to directly using the 4th power of the graph adjacency matrix.

2.2 Related Work

Matrix factorization has been used extensively in recommendation systems with both explicit [61] and implicit [49] feedback. Such methods compute low dimensional user and item representations; their inner product approximates the observed (or to be predicted) entry in the target matrix. To incorporate graph side information in these systems, [89, 128] used a graph Laplacian based regularization framework that forces a pair of node representations to be similar if they are connected via an edge in the graph. In [126], this was extended to the implicit feedback setting. [67] proposed a method that incorporates first-order information of the rating bipartite graph into the model by considering item co-occurrences. More recently, GC-MC [8] used a GCN approach performing convolutions on the main bipartite graph by treating the first-order side graph information as features, and [74] proposed combining GCNs and RNNs for the same task.

Methods that use higher order graph information are typically based on taking random walks on the graphs [31]. [52] extended this method to include graph side information in the model. Finally, the PageRank [76] algorithm can be seen as computing the steady

state distribution of a Markov network, and similar methods for recommender systems was proposed in [1, 122].

For a complete list of related works of representation learning on graphs, we refer the interested user to [36]. For the collaborative filtering setting, [8, 74] use Graph Convolutional Neural Networks (GCN) [23], but with some modifications. Standard GCN methods without substantial modifications cannot be directly applied to collaborative filtering rating datasets, including well-known approaches like GCN [58] and GraphSage [35], because they are intended to solve semi-supervised classification problem over graphs with nodes' features. PinSage [123] is the GraphSage extension to non-personalized graph-based recommendation algorithm but is not meant for collaborative filtering problems. GC-MC [8] extends GCN to collaborative filtering, albeit it is less scalable than [123]. Our Graph DNA scheme can be used to obtain graph features in these extensions. In contrast to the above-mentioned methods involving GCNs, we do not use the data driven loss function to train our graph encoder. This property makes our graph DNA suitable for both transductive as well as inductive problems.

Bloom filters have been used in Machine Learning for multi-label classification [21], and for hashing deep neural network models representations [22, 37, 99]. However, to the best of our knowledge, until now, they have not been used to encode graphs, nor has this encoding been applied to recommender systems. So it would be interesting to extend our work to other recommender systems settings, such as [118] and [117].

2.3 Methodology

We consider the recommender system problem with a partially observed rating matrix R and a Graph that encodes side information G . In this section, we will introduce the Graph DNA algorithm for encoding deep neighborhood information in G . In the next section, we will show how this encoded information can be applied to various graph based recommender systems.

2.3.1 Bloom Filter

The Bloom filter [9] is a probabilistic data structure designed to represent a set of elements. Thanks to its space-efficiency and simplicity, Bloom filters are applied in many real-world applications such as database systems [10,16]. A Bloom filter \mathcal{B} consists of k independent hash functions $h_t(x) \rightarrow \{1, \dots, c\}$. The Bloom filter \mathcal{B} of size c can be represented as a length c bit-array \mathbf{b} . More details about Bloom filters can be found in [12]. Here we highlight a few desirable properties of Bloom filters essential to our graph DNA encoding:

1. Space efficiency: classic Bloom filters use $1.44 \log_2(1/\epsilon)$ of space per inserted key, where ϵ is the false positive rate associated with this Bloom filter.
2. Support for the union operation of two Bloom filters: the Bloom filter for the union of two sets can be obtained by performing bitwise ‘OR’ operations on the underlying bit-arrays of the two Bloom filters.
3. Size of the Bloom filter can be approximated by the number of nonzeros in the underlying bit array: in particular, given a Bloom filter representation $\mathcal{B}(A)$ of a set A : the number of elements of A can be estimated as $|A| \approx -\frac{c}{k} \log\left(1 - \frac{\text{nnz}(\mathbf{b})}{c}\right)$, where $\text{nnz}(\mathbf{b})$ is the number of non-zero elements in array \mathbf{b} . As a result, the number of common nonzero bits of $\mathcal{B}(A_1)$ and $\mathcal{B}(A_2)$ can be used as a proxy for $|A_1 \cap A_2|$.

Algorithm 1. Graph DNA Encoding with Bloom Filters

Input: G : a graph of n nodes, c : the length of codes, k : the number of hash functions, d : the number of iterations, θ : tuning parameter to control the number of elements hashed.

Output: $B \in \{0, 1\}^{n \times c}$: a boolean matrix to denote the bipartite relationship between n nodes and c bits.

- $\mathcal{H} \leftarrow \{\mathbf{h}_t(\cdot) : t = 1, \dots, k\}$ ▷ Pick k hash functions
 - **for** $i = 1, \dots, n$: ▷ *GraphBloom* Initialization
 - $\mathcal{B}^0[\mathbf{i}] \leftarrow \text{BloomFilter}(c, \mathcal{H})$
 - $\mathcal{B}^0[\mathbf{i}].\text{add}(i)$
 - **for** $s = 1, \dots, d$: ▷ d times neighborhood propagations
 - **for** $i = 1, \dots, n$:
 - * **for** $j \in \mathcal{N}_1(i)$: ▷ degree-1 neighbors
 - **if** $|\mathcal{B}^s[\mathbf{i}]| > \theta$: break;
 - $\mathcal{B}^s[\mathbf{i}].\text{union}(\mathcal{B}^{s-1}[\mathbf{j}])$
 - $B_{ij} \leftarrow \mathcal{B}^d[\mathbf{i}].\mathbf{b}[\mathbf{j}] \ \forall (i, j) \in [n] \times [c]$
-

2.3.2 Graph DNA Encoding Via Bloom Filters

Now we introduce our Graph DNA encoding. The main idea is to encode the deep (multi-hop) neighborhood aware embedding for each node in the graph approximately using the Bloom filter, which helps avoid performing computationally expensive graph adjacency matrix multiplications. In Graph DNA, we have Bloom filters $\mathcal{B}[\mathbf{i}], i = 1, \dots, n$ for the n graph nodes. All the Bloom filters $\mathcal{B}[\mathbf{i}]$ share the same k hash functions. The role of $\mathcal{B}[\mathbf{i}]$ is to store the deep neighborhood information of the i -th node. Taking advantage of the union operations of Bloom filters, one node’s neighborhood information can be propagated to its neighbors in an iterative manner using gossip algorithms [97]. Initially, each $\mathcal{B}[\mathbf{i}]$ contains only the node itself. At the s -th iteration, $\mathcal{B}[\mathbf{i}]$ is updated by taking union with node i ’s immediate neighbors’ Bloom filters $\mathcal{B}[\mathbf{j}]$. By induction, we see that after the d iterations, $\mathcal{B}[\mathbf{i}]$ represents $\mathcal{N}_d(i) := \{j : \text{distance}_G(i, j) \leq d\}$, where $\text{distance}_G(i, j)$

is the shortest path distance between nodes i and j in G . As the last step, we stack array representations of all Bloom filters and form a sparse matrix $B \in \{0, 1\}^{n \times c}$, where the i -th row of B is the bit representation of $\mathcal{B}[i]$. As a practical measure, to prevent over-saturation of Bloom filters for popular nodes in the graph, we add a hyper-parameter θ to control the max saturation level allowed for Bloom filters. This would also prevent hub nodes dominating in graph DNA encoding. The pseudo-code for the proposed encoding algorithm is given in Algorithm 1. We use graph DNA- d to denote our obtained graph encoding after applying Algorithm 1 with s looping from 1 to d . We also give a simple example to illustrate how the graph DNA is encoded into Bloom filter representations in Figure 2.1. Our usage of Bloom filters is very different from previous works in [86, 96, 101], which use Bloom filter for standard hashing and is unrelated to graph encoding.

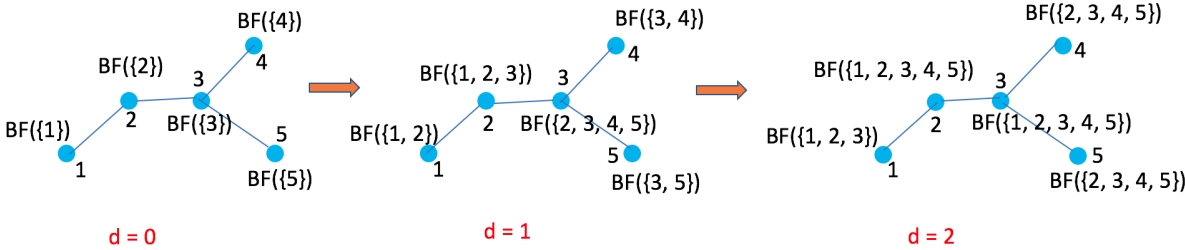


Figure 2.1. Illustration of Algorithm 1: the graph DNA encoding procedure. The curly brackets at each node indicate the nodes encoded at a particular step. At $d = 0$ each node’s Bloom filter only encodes itself, and multi-hop neighbors are included as d increases.

2.4 Collaborative Filtering with Graph DNA

Suppose we are given the sparse rating matrix $R \in \mathbb{R}^{n \times m}$ with n users and m items, and a graph $G \in \mathbb{R}^{n \times n}$ encoding relationships between users. For simplicity, we do not assume a graph on the m items, though including it should be straightforward.

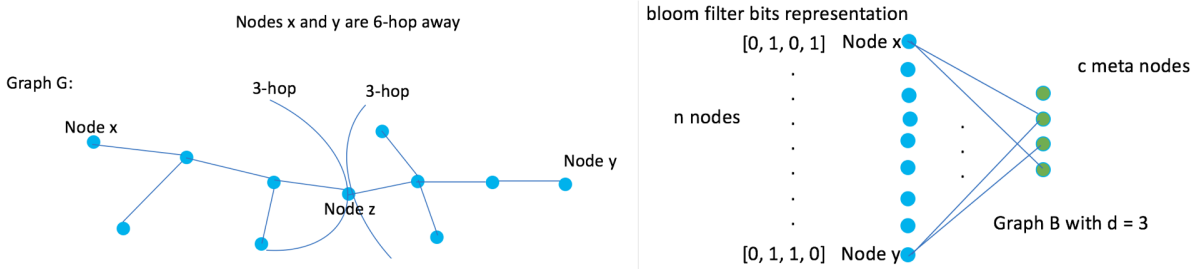


Figure 2.2. Illustration of our proposed DNA encoding method (DNA-3), with the corresponding bipartite graph representation.

2.4.1 Graph Regularized Matrix Factorization

Explicit Feedback : The objective function of Graph Regularized Matrix Factorization (GRMF) [13, 89, 128] is:

$$\min_{U, V} \sum_{(i, j) \in \Omega} (R_{i, j} - u_i^\top v_j)^2 + \frac{\lambda}{2} (\|U\|_F^2 + \|V\|_F^2) + \mu \text{tr}(U^\top \text{Lap}(G)U) \quad (2.1)$$

where $U \in \mathbb{R}^{n \times r}$, $V \in \mathbb{R}^{m \times r}$ are the embeddings associated with users and items respectively, n is the number of users and m is the number of items, $R \in \mathbb{R}^{n \times m}$ is the sparse rating matrix, $\text{tr}(\cdot)$ is the trace operator, λ, μ are tuning coefficients, and $\text{Lap}(\cdot)$ is the graph Laplacian operator.

The last term is called graph regularization, which tries to enforce similar nodes (measured by edge weights in G) to have similar embeddings. One naive way [14] to extend this to higher-order graph regularization is to replace the graph G with $\sum_{i=1}^K w_i \cdot G^i$ and then use the graph Laplacian of $\sum_{i=1}^K w_i \cdot G^i$ to replace G in (2.1). Computing G^i for even small i is computationally infeasible for most real-world applications, and we will soon lose the sparsity of the graph, leading to memory issues. Sampling or thresholding could mitigate the problem but suffers from performance degradation.

In contrast, our graph DNA obtained from Algorithm 1 does not suffer from any of these issues. The space complexity of our method is only of order $O(n \log n)$ for a graph with n nodes, instead of $O(n^2)$. The reduced number of non-zero elements using graph

DNA leads to a significant speed-up in many cases.

We can easily use graph DNA in GRMF as follows: we treat the c bits as c new pseudo-nodes and add them to the original graph G . We then have $n + c$ nodes in a modified graph \dot{G} :

$$\dot{G} = \begin{bmatrix} G \in \mathbb{R}^{n \times n} & B \in \mathbb{R}^{n \times c} \\ B^\top \in \mathbb{R}^{c \times n} & \mathbf{0} \in \mathbb{R}^{c \times c} \end{bmatrix}. \quad (2.2)$$

To account for the c new nodes, we expand $U \in \mathbb{R}^{n \times r}$ to $\dot{U} \in \mathbb{R}^{(n+c) \times r}$ by appending parameters for the meta-nodes. The objective function for GRMF with Graph DNA will be the same as (2.1) except replacing U and G with \dot{U} and \dot{G} . At the prediction stage, we discard the meta-node embeddings.

Implicit Feedback : For implicit feedback data, when R is a 0/1 matrix, weighted matrix factorization is a widely used algorithm [48, 49]. The only difference is that the loss function in (2.1) is replaced by $\sum_{(i,j):R_{ij}=1} (R_{ij} - u_i^\top v_j)^2 + \sum_{(i,j):R_{ij}=0} \rho (R_{ij} - u_i^\top v_j)^2$ where $\rho < 1$ is a hyper-parameter reflecting the confidence of zero entries. In this case, we can apply the Graph DNA encoding as before trivially.

2.4.2 Co-Factorization with Graph Information

Co-Factorization of Rating and Graph Information (Co-Factor) [67, 103] is ideologically very different from GRMF and GRWMF, because it does not use graph information as regularization term. Instead it treats the graph adjacency matrix as another rating matrix, sharing one-sided latent factors with the original rating matrix. Co-Factor minimizes the following objective function: $\min_{U,V} \sum_{(i,j) \in \Omega_R} (R_{i,j} - u_i^\top v_j)^2 + \frac{\lambda}{2} (\|U\|_F^2 + \|V\|_F^2 + \|V'\|_F^2) + \sum_{(i,j) \in \Omega_G} (G_{i,j} - u_i^\top v'_j)^2$, where $U \in \mathbb{R}^{n \times r}, V \in \mathbb{R}^{m \times r}, V' \in \mathbb{R}^{n \times r}$. We can extend Co-Factor to incorporate our DNA-d by replacing G with B in the equation above, where $B \in \mathbb{R}^{n \times c}$ is the Bloom filter bipartite graph adjacency matrix of n real-user nodes and c pseudo-user nodes, similar to B as in (2.2). We call the extension Co-Factor_DNA-d.

2.4.3 Graph Convolutional Matrix Completion

Graph Convolutional Matrix Completion (GC-MC) is a graph convolutional network (GCN) based geometric matrix completion method [8]. In [8], the rating matrix R is

treated as an adjacency matrix in GCN while side information G is treated as feature matrix for nodes — each user has an n -dimensional 0/1 feature that corresponds to a column of G . The GCN model then performs convolutions of these features on the bipartite rating graph. Convolutions of these features are performed on the bipartite rating graph. We find in our experiments that using these one-hot encodings of the graph as feature is an inferior choice both in terms of performance and speed. To capture higher order side graph information, it is better to use $G + \alpha G^2$ for some constant α and this alternate choice usually gives smaller generalization error than the original GC-MC method. However, it is hard to explicitly calculate $G + \alpha G^2$ and store the entire matrix for a large graph for the same reason described in Section 2.4.1. Again, we can use graph DNA to efficiently encode and store the higher order information before feeding it into GC-MC. We show in our experiments that this outperforms current state-of-the-art GCN methods [8, 74] as well as GC-MC with graph encoding methods that require training, such as Node2vec [32] and Deepwalk [84]. Our encoding scheme does not require training and therefore is a lot faster than previous encoding methods. More details are discussed in the experiment section 2.5.3.

2.5 Experiments

We show that our Graph DNA encoding technique can improve the performance of 4 popular graph-based recommendation algorithms: graph-regularized matrix factorization, co-factorization, weighted matrix factorization, and GCN-based graph convolution matrix factorization. All experiments except GCN are conducted on a server with Intel Xeon E5-2699 v3 @ 2.30GHz CPU and 256G RAM. The GCN experiments are conducted on Google Cloud with Nvidia V100 GPU.

2.5.1 Simulation Study

We first simulate a user/item rating dataset with user graph as side information, generate its graph DNA, and use it on a downstream task: matrix factorization.

We randomly generate user and item embeddings from standard Gaussian distributions, and construct an Erdős-Rényi Random graph of users. User embeddings are generated

using Algorithm 11 in Appendix: at each propagation step, each user’s embedding is updated by an average of its current embedding and its neighbors’ embeddings. Based on user and item embeddings after $T = 3$ iterations of propagation, we generate the underlying ratings for each user-item pairs according to the inner product of their embeddings, and then sample a small portion of the dense rating matrix as training and test sets.

We implement our graph DNA encoding algorithm in python using a scalable python library [3] to generate Bloom filter matrix B . We adapt the GRMF C++ code to solve the objective function of GRMF_DNA-K with our Bloom filter enhanced graph \dot{G} . We compare the following variants:

1. MF: classical matrix factorization only with ℓ_2 regularization without graph information.
2. GRMF_ G^d : GRMF with ℓ_2 regularization and using G, G^2, \dots, G^d [14].
3. GRMF_DNA- d : GRMF with ℓ_2 but using our proposed graph DNA- d .

We report the prediction performance with Root Mean Squared Error (RMSE) on test data. All results are reported on the test set, with all relevant hyperparameters tuned on a held-out validation set. To accurately measure how large the relative gain is from using deeper information, we introduce a new metric called Relative Graph Gain (RGG) for using information X , which is defined as:

$$\text{RGG}(X)\% = \left(\frac{\text{RMSE without Graph} - \text{RMSE with } X}{\text{RMSE without Graph} - \text{RMSE with } G} - 1 \right) \times 100, \tag{2.3}$$

where RMSE is measured for the same method with different graph information. This metric would be 0 if only first order graph information is utilized and is only defined when the denominator is positive.

In Table 2.1, we can easily see that using a deeper neighborhood helps the recommendation performances on this synthetic dataset. Graph DNA-3’s gain is 166% larger than that of using first-order graph G . We can see an increase in performance gain for an increase in depth d when $d \leq 3$. This is expected because we set $T = 3$ during our creation of this dataset.

2.5.2 Graph Regularized Matrix Factorization for Explicit Feedback

Next, we show that graph DNA can improve the performance of GRMF for explicit feedback. We conduct experiments on two real datasets: Douban [70] and Flixster [127]. Both datasets contain explicit feedback with ratings from 1 to 5. There are 129,490 users, 58,541 items in Douban. There are 147,612 users, 48,794 items in Flixster. Both datasets have a graph defined on the respective sets of users.

We pre-processed Douban and Flixster following the same procedure in [89, 115]. The experimental setups and comparisons are almost identical to the synthetic data experiment (see details in section 2.5.1). Due to the exponentially growing non-zero elements in the graph as we go deeper (see Table 8.2), we are unable to run full GRMF_ G^4 and GRMF_ G^5 for these datasets. In fact, GRMF_ G^3 itself is too slow so we thresholded G^3 by only considering entries whose values are equal to or larger than 4. For the Bloom filter, we set a false positive rate of 0.1 and use capacity of 500 for Bloom filters, resulting in $c = 4,796$.

We can see from Table 2.1 that deeper graph information always helps. For Douban, graph DNA-3 is most effective, giving a relative graph gain of 82.79% compared to only 2% gain when using G^2 or G^3 naively. Interestingly for Flixster, using G^2 is better than using G^3 . However, Graph DNA-3 and DNA-4 yield 10x and 15x performance improvements respectively, lending credence to the implicit regularization property of graph DNA. For a fixed size Bloom filter, the computational complexity of graph DNA scales linearly with depth d , as compared to exponentially for GRMF_ G^d . We measure the speed in Table 2.2. The memory cost is only a fraction of n^2 after hashing. Such low memory and computational complexity allow us to scale to larger d , compared to baseline methods.

2.5.3 Co-Factorization with Graph for Explicit Feedback

We show our graph DNA can improve Co-Factor [67, 103] as well. The results are in Table 2.1. We find that applying DNA-3 to the Co-Factor method improves performance on both the datasets, more so for Flixster. This is consistent with our observations for GRMF in Table 2.1: deep graph information is more helpful for Flixster than Douban. Applying Graph DNA to Co-Factor is detailed in the Appendix.

Table 2.1. Comparison of Graph Regularized Matrix Factorization Variants for Explicit Feedback on Synthetic, Douban and Flixster data. We use rank $r = 10$. RGG is the Relative Graph Gain defined in (2.3).

Dataset	Synthetic		Douban		Flixster	
	RMSE ($\times 10^{-1}$)	% RGG	RMSE ($\times 10^{-1}$)	% RGG	RMSE ($\times 10^{-1}$)	% RGG
MF	2.9971	-	7.3107	-	8.8111	-
GRMF_G	2.7823	0	7.2398	0	8.8049	0
GRMF_G ²	2.6543	59.5903	7.2381	2.3977	8.7849	322.5806
GRMF_G ³	2.5687	99.4413	7.2432	-4.7954	8.7932	188.7097
GRMF_G ⁴	2.5562	105.2607	-	-	-	-
GRMF_G ⁵	2.4853	138.2682	-	-	-	-
GRMF_G ⁶	2.4852	138.3147	-	-	-	-
GRMF_DNA-1	2.4303	163.8734	7.2191	29.1960	8.8013	58.0645
GRMF_DNA-2	2.4510	154.2365	7.2359	5.5007	8.8007	67.7419
GRMF_DNA-3	2.4247	166.4804	7.1811	82.7927	8.7383	1074.1935
GRMF_DNA-4	2.4466	156.2849	7.1971	60.2257	8.7122	1495.1613
Co-Factor_G	-	-	7.2743	0	8.7957	0
Co-Factor_DNA-3	-	-	7.2623	32.9670	8.7354	391.5584

Table 2.2. Graph DNA (Algorithm 1) Encoding Speed. We set number $c = 500$ and implement Graph DNA using single-core python. We can scale up linearly in terms of depth d for a fixed c .

Dataset	Graph Statistics		Graph DNA Encoding Time (secs)			
	Number of Nodes	Graph Density	DNA-1	DNA-2	DNA-3	DNA-4
Douban	129,490	0.0102%	132.2717	266.3740	403.9747	580.1547
Flixster	147,612	0.0117%	157.3103	317.7706	482.0360	686.8048

Table 2.3. Comparison of GRWMF Variants for Implicit Feedback on Douban and Flixster datasets. P stands for precision and N stands for NDCG. We use rank $r = 10$ and all results are in %.

Dataset	Methods	MAP	HLU	P@1	P@5	N@1	N@5
Douban	GRWMF_G	8.340	13.033	14.944	10.371	14.944	12.564
	GRWMF_DNA-3	8.400	13.110	14.991	10.397	14.991	12.619
Flixster	GRWMF_G	10.889	14.909	12.303	7.9927	12.303	12.734
	GRWMF_DNA-3	11.612	15.687	12.644	8.1583	12.644	13.399

Table 2.4. Comparison of GCN Methods for Explicit Feedback on Douban, Flixster and Yahoo Music datasets (3000 by 3000 as in [8, 74]). All the methods except GC-MC utilize side graph information.

Dataset	Methods	Test RMSE ($\times 10^{-1}$)	% RGG
Douban	SRGCNN (reported by [8])	-	-
	GC-MC	7.3109 \pm 0.0150	-
	GC-MC_G	7.3698 \pm 0.0737	N/A
	GC-MC_G ²	7.3123 \pm 0.0139	N/A
	GC-MC_Node2vec	7.3666 \pm 0.0218	N/A
	GC-MC_Deepwalk	7.3394 \pm 0.0343	N/A
	GC-MC_DNA-2	7.3117 \pm 0.0129	N/A
Flixster	SRGCNN (reported by [8])	9.2600	-
	GC-MC	9.2614 \pm 0.0578	-
	GC-MC_G	9.2374 \pm 0.1045	0
	GC-MC_G ²	8.9344 \pm 0.0333	1262.4999
	GC-MC_Node2vec	12.0370 \pm 1.9474	N/A
	GC-MC_Deepwalk	9.0507 \pm 0.1692	777.9167
	GC-MC_DNA-2	8.9536 \pm 0.0770	1182.4999
Yahoo Music	SRGCNN (reported by [8])	-	-
	GC-MC	22.6697 \pm 0.3530	-
	GC-MC_G	21.3672 \pm 0.4190	0
	GC-MC_G ²	20.2189 \pm 0.8664	88.1612
	GC-MC_Node2vec	19.8901 \pm 0.7948	113.4050
	GC-MC_Deepwalk	20.1603 \pm 0.9342	92.6603
	GC-MC_DNA-2	19.3879 \pm 0.2874	151.9616

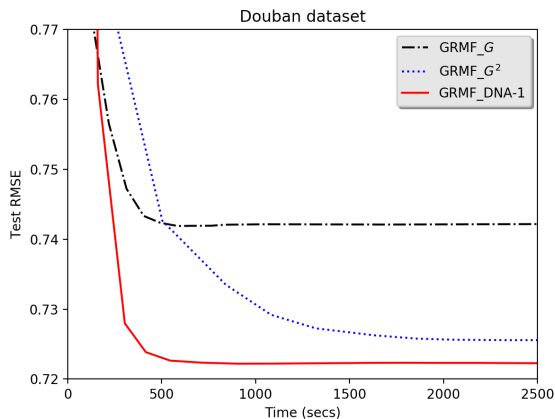


Figure 2.3. Compare Training Speed of GRMF, with and without Graph DNA.

2.5.4 Graph Regularized Weighted Matrix Factorization for Implicit Feedback

We follow the same procedure as in [116] to set ratings of 4 and above to 1, and the rest to 0. We compare the baseline graph based weighted matrix factorization [48, 49] with our proposed weighted matrix factorization with DNA-3. We do not compare with Bayesian personalized ranking [91] and the recently proposed SQL-rank [116] as they cannot easily utilize graph information.

The results are summarized in Table 2.3 with experimental details in the Appendix. Again, using DNA-3 achieves better prediction results over the baseline in terms of every single metric on both Douban and Flixster datasets.

2.5.5 Graph Convolutional Matrix Factorization

We can use graph DNA instead to efficiently encode and store the higher order information before feeding it into GC-MC.

We use the same split of three real-world datasets and follow the exact procedures as in [8, 74]. We tuned hyperparameters using a validation dataset and obtain the best test results found within 200 epochs using optimal parameters. We repeated the experiments 6 times and report the mean and standard deviation of test RMSE. After some tuning, we use the capacity of 10 Bloom filters for Douban and 60 for Flixster, as the latter has a much denser second-order graph. With a false positive rate of 0.1, this implies that we use

96-bits Bloom filters for Douban and 960 bits for Flixster. We use the resulting bloom filter bitarrays as the node features, and pass that as the input to GC-MC. Using Graph DNA-2, the input feature dimensions are thus reduced from 3000 to 96 and 960, which leads to a significant speed-up. The original GC-MC method did not scale up well beyond 3000 by 3000 rating matrices with the user and the item side graphs as it requires using normalized adjacency matrix as user/item features. PinSage [123], while scalable, does not utilize the user/item side graphs. Furthermore, it is not feasible to have $O(n)$ dimensional features for the nodes, where n is the number of nodes in side graphs. In contrast, our method only requires $O(\log(n))$ dimensional features. We can see from Table 2.4 that we outperform both GCN-based methods [8] and [74] in terms of performance by a large margin.

Note that another potential way to improve over GC-MC is to use other graph encoding schemes like Node2Vec [32] and DeepWalk [84] to encode the user-user graph into node features. One clear drawback is that those graph embedding methods are time-consuming. Using the official Node2vec implementation, excluding reading and writing, it takes 416.13 seconds to encode the 3K by 3K subsampled Yahoo-Music item graph and obtain resulting 760-d node embeddings. For our method, it only takes 7.55 seconds to obtain the same 760-d features. Similarly, it takes over 15 mins to run the official C++ codes for DeepWalk [84] using the same parameters as Node2Vec to encode the graph. In fact, fast encoding via hashing and bitwise-or that does not require training is one of the main advantages of our method.

Furthermore, even without considering the time overhead, we found our graph DNA encoding outperforms Node2Vec and DeepWalk in terms of test RMSE. Details can be found in Table 2.4. This could be due to that encoding higher-order information is more important for graph-regularized recommendation tasks, and graph DNA is a better and more direct way to encode higher order information compared with Node2Vec and DeepWalk.

Speed Comparisons Next, we compare the speed-ups obtained by graph DNA- d with GRMF G^d (a naive way to encode higher order information by computing powers of G).

Table 2.5. Comparison of GRMF Methods of different ranks for Explicit Feedback on Flixster Dataset.

Rank	methods	test RMSE ($\times 10^{-1}$)	% gain
10	GRMF_ G^2	8.7849	-
	GRMF_DNA-3	8.7383	0.8262
20	GRMF_ G^2	8.9179	-
	GRMF_DNA-3	8.7565	1.8098
30	GRMF_ G^2	9.0865	-
	GRMF_DNA-3	8.9255	1.7719

Figure 3 suggests that graph DNA-1 (which encodes hop-2 information) scales better than directly computing G^2 in GRMF.

Exploring Effects of Rank Finally, we investigate whether the proposed DNA coding can achieve consistent improvements when varying the rank in the GRMF algorithm. In Table 2.5, we compare the proposed GRMF_DNA-3 with GRMF_ G^2 , which achieves the best RMSE without using DNA coding in the previous tables. The results clearly show that the improvement of the proposed DNA coding is consistent over different ranks and works even better when rank is larger.

2.6 Conclusion

In this chapter, we proposed Graph DNA, a deep neighborhood aware encoding scheme for collaborative filtering with graph information. We make use of Bloom filters to incorporate higher order graph information, without the need to explicitly minimize a loss function. The resulting encoding is extremely space and computationally efficient, and lends itself well to multiple algorithms that make use of graph information, including Graph Convolutional Networks. Experiments show that Graph DNA encoding outperforms several baseline methods on multiple datasets in both speed and performance.

Chapter 3

Large-scale Pairwise Collaborative Ranking in Near-Linear Time

3.1 Introduction

In online retail and online content delivery applications, it is commonplace to have embedded recommendation systems—algorithms that recommend items to users based on previous user behaviors and ratings. Online retail companies develop sophisticated recommendation systems based on purchase behavior, item context, and shifting trends. The Netflix prize [7], in which competitors utilize user ratings to recommend movies, accelerated research in recommendation systems. While the winning submissions agglomerated several existing methods, one essential methodology, latent factor models, emerged as a critical component. The latent factor model means that the approximated rating for user i and item j is given by $u_i^\top v_j$ where u_i, v_j are k -dimensional vectors. One interpretation is that there are k latent topics and the approximated rating can be reconstructed as a combination of factor weights. By minimizing the square error loss of this reconstruction we arrive at the incomplete SVD,

$$\min_{U,V} \sum_{i,j \in \Omega} (R_{i,j} - u_i^\top v_j)^2 + \frac{\lambda}{2} (\|U\|_F^2 + \|V\|_F^2), \quad (3.1)$$

where Ω contains sampled indices of the rating matrix, R .

Often the performance of recommendation systems is not measured by the quality of rating prediction, but rather the ranking of the items that the system returns for

a given user. The task of finding a ranking based on ratings or relative rankings is called Collaborative Ranking. Recommendation systems can be trained with ratings, that may be passively or actively collected, or by relative rankings, in which a user is asked to rank a number of items. A simple way to unify the framework is to convert the ratings into rankings by making pairwise comparisons of ratings. Specifically, the algorithm takes as input the pairwise comparisons, $Y_{i,j,k}$ for each user i and item pairs j, k . This approach confers several advantages. Users may have different standards for their ratings, some users are more generous with their ratings than others. This is known as the calibration drawback, and to deal with this we must make a departure from standard matrix factorization methods. Because we focus on ranking and not predicting ratings, we can expect improved performance when recommending the top items. Our goal in this chapter is to provide a collaborative ranking algorithm that can scale to the size of the full Netflix dataset, a heretofore open problem.

The existing collaborative ranking algorithms, (for a summary see section 3.2), are limited by the number of observed ratings per user in the training data and cannot scale to massive datasets, therefore, making the recommendation results less accurate and less useful in practice. This motivates our algorithm, which can make use of the entire Netflix dataset without sub-sampling. Our contribution can be summarized below:

- For input data in the form of pairwise preference comparisons, we propose a new algorithm Primal-CR that alternatively minimizes latent factors using Newton’s method in the primal space. By carefully designing the computation of gradient and Hessian vector product, our algorithm reduces the sample complexity per iteration to $O(|\Omega| + d_1\bar{d}_2r)$, while the state-of-the-art approach [81] have $O(|\Omega|r)$ complexity. Here $|\Omega|$ (total number of pairs), is much larger than $d_1\bar{d}_2$ (d_1 is number of users and \bar{d}_2 is averaged number of items rated by a user). For the Netflix problem, $|\Omega| = 2 \times 10^{10}$ while $d_1\bar{d}_2 = 10^8$.
- For input data in the form of ratings, we can further exploit the structure to speedup the gradient and Hessian computation. The resulting algorithm, Primal-CR++, can further reduce the time complexity to $O(d_1\bar{d}_2(r + \log \bar{d}_2))$ per iteration. In this

setting, our algorithm has time complexity near-linear to the input size, and have comparable speed with classical matrix factorization model that takes $O(d_1 \bar{d}_2 r)$ time, while we can achieve much better recommendation by minimizing the ranking loss. We show that our algorithms outperform existing algorithms on real world datasets and can be easily parallelized.

3.2 Related Work

Collaborative filtering methodologies are summarized in [95] (see [24] for an early work). Among them, matrix factorization [61] has been widely used due to the success in the Netflix Prize. Many algorithms have been developed based on matrix factorization [19, 48, 90, 91, 102], and many scalable algorithms have been developed [29, 61]. However, they are not suitable for ranking top items for a user due to the fact that their goal is to minimize the mean-square error (MSE) instead of ranking loss. In fact, MSE is not a good metric for recommendation when we want to recommend the top K items to a user. This has been pointed out in several papers [5] which argue normalized discounted cumulative gain (NDCG) should be used instead of MSE, and our experimental results also confirm this finding by showing that minimizing the ranking loss results in better precision and NDCG compared with the traditional matrix factorization approach that is targeting squared error.

Ranking is a well studied problem, and there has been a long line of research focuses on learning one ranking function, which is called Learning to Rank. For example, RankSVM [53] is a well-known pair-wise model, and an efficient solver has been proposed in [17] for solving rankSVM. [15] is a list-wise model implemented using neural networks. Another class of point-wise models fit the ratings explicitly but has the issue of calibration drawback (see [34]).

The collaborative ranking (CR) problem is essentially trying to learn multiple rankings together, and several models and algorithms have been proposed in literature. The Cofrank algorithm [114], which tailors maximum margin matrix factorization [105] for collaborative ranking, is a point-wise model for CR, and is regarded as the performance benchmark for

this task. If the ratings are 1-bit, a weighting scheme is proposed to improve the usual point-wise Matrix Factorization approach [78]. List-wise models for Learning to Rank can also be extended to many rankings setting, [100]. However it is still quite similar to a point-wise approach since they only consider the top-1 probabilities.

For pairwise models in collaborative ranking, it is well known that they do not encounter the calibration drawback as do point-wise models, but they are computationally intensive and cannot scale well to large data sets [100]. The scalability problem for pairwise models is mainly due to the fact that their time complexity is at least proportional to $|\Omega|$, the number of pairwise preference comparisons, which grows quadratically with number of rated items for each user. Recently, [81] proposed a new Collrank algorithm, and they showed that Collrank has better precision and NDCG as well as being much faster compared with other CR methods on real world datasets, including Bayesian Personalized Ranking (BPR) [91]. Unfortunately their scalability is still constrained by number of pairs, so they can only run on subsamples for large datasets, such as Netflix. In this chapter, our algorithm Primal-CR and Primal-CR++ also belong to the family of pairwise models, but due to cleverly re-arranging the computation, we are able to have much better time complexity than existing ones, and as a result our algorithm can scale to very large datasets.

There are many other algorithms proposed for many rankings setting but none of these mentioned below can scale up to the extent of using all the ratings in the full Netflix data. There are a few using Bayesian frameworks to model the problem [91], [79], [111], the last of which requires many specified parameters. Another one proposed retargeted matrix factorization to get ranking by monotonically transforming the ratings [62]. [33] proposes a similar model without making generative assumptions on ratings besides assuming low-rank and correctness of the ranking order.

3.3 Problem Formulation

We first formally define the collaborative ranking problem using the example of item recommender system. Assume we have d_1 users and d_2 items, the input data is given in the form of “for user i , item j is preferred over item k ” and thus can be represented by

a set of tuples (i, j, k) . We use Ω to denote the set of observed tuples, and the observed pairwise preferences are denoted as $\{Y_{ijk} \mid (i, j, k) \in \Omega\}$, where $Y_{ijk} = 1$ denotes that item j is preferred over item k for a particular user i and $Y_{ijk} = -1$ to denote that item k is preferred over item j for user i .

The goal of collaborative ranking is to rank all the unseen items for each user i based on these partial observations, which can be done by fitting a scoring matrix $X \in \mathbb{R}^{d_1 \times d_2}$. If the scoring matrix has $X_{ij} > X_{ik}$, it implies that item j is preferred over item k by the particular user i and therefore we should give higher rank for item j than item k . After we estimate the scoring matrix X by solving the optimization problem described below, we can then recommend top k items for any particular user.

The Collaborative Ranking Model referred to in this chapter is the one proposed recently in [81]. It belongs to the family of pairwise models for collaborative ranking because it uses pairwise training losses [5]. The model is given as

$$\min_X \sum_{(i,j,k) \in \Omega} L(Y_{ijk}(X_{ij} - X_{ik})) + \lambda \|X\|_*, \quad (3.2)$$

where

$L(\cdot)$ is the loss function, $\|X\|_*$ is the nuclear norm regularization defined by the sum of all the singular value of the matrix X , and λ is a regularization parameter. The ranking loss defined in the first term of (3.2) penalizes the pairs when $Y_{ijk} = 1$ but $X_{ij} - X_{ik}$ is positive but small, and penalizes even more when the difference is negative. The second term in the loss function is based on the assumption that there are only a small number of latent factors contributing to the users' preferences which is analogous to the idea behind incomplete SVD for matrix factorization mentioned in the introduction. In general we can use any loss function, but since

L_2 -hinge loss defined as

$$L(a) = \max(0, 1 - a)^2 \quad (3.3)$$

gives the best performance in practice [81] and enjoys many nice properties, such as smoothness and differentiable, we will focus on L_2 -hinge loss in this chapter. In fact, our first algorithm Primal-CR can be applied to any loss function, while Primal-CR++ can

only be applied to L_2 -hinge loss.

Despite the advantage of the objective function in equation (3.2) being convex, it is still not feasible for large-scale problems since d_1 and d_2 can be very large so that the scoring matrix X cannot be stored in memory, not to mention how to solve it. Therefore, in practice people usually transform (3.2) to a non-convex form by replacing $X = UV^T$, and in that case since $\|X\|_* = \min_{X=UV^T} \frac{1}{2}(\|U\|_F^2 + \|V\|_F^2)$ [105], problem (3.2) can be reformulated as

$$\min_{U,V} \sum_{(i,j,k) \in \Omega} L(Y_{ijk} \cdot u_i^T(v_j - v_k)) + \frac{\lambda}{2}(\|U\|_F^2 + \|V\|_F^2), \quad (3.4)$$

We use u_i and v_j denote columns of U and V respectively. Note that [81] also solves the non-convex form (3.4) in their experiments, and in the rest of the paper we will propose a faster algorithm for solving (3.4).

3.4 Proposed Algorithms

3.4.1 Motivation and Overview

Although collaborative ranking assumes that input data is given in the form of pairwise comparisons, in reality almost all the datasets (Netflix, Yahoo-Music, MovieLens, etc) contain user ratings to items in the form of $\{R_{ij} \mid (i, j) \in \bar{\Omega}\}$, where $\bar{\Omega}$ is the subset of observed user-item pairs. Therefore, in practice we have to transform the rating-based data into pair-wise comparisons by generating all the item pairs rated by the same user:

$$\Omega = \{(i, j, k) \mid j, k \in \bar{\Omega}_i\}, \quad (3.5)$$

where $\bar{\Omega}_i := \{j \mid (i, j) \in \bar{\Omega}\}$ is the set of items rated by user i . Assume there are averagely \bar{d}_2 items rated by a user (i.e., $\bar{d}_2 = \text{mean}(|\bar{\Omega}_i|)$), then the collaborative ranking problem will have $O(d_1 \bar{d}_2^2)$ pairs and thus the size of Ω grows quadratically.

Unfortunately, all the existing algorithms have $O(|\Omega|r)$ complexity, so they cannot scale to large number of items. For example, the AltSVM (or referred to as Collrank) Algorithm in [81] will run out of memory when we subsample 500 rated items per user on Netflix dataset since its implementation¹ stores all the pairs in memory and therefore

¹Collrank code is available on <https://github.com/dhpark22/collranking>.

requires $O(|\Omega|)$ memory. So it cannot be used for the full Netflix dataset which has more than *20 billion pairs* and requires 300GB memory space. To the best of our knowledge, no collaborative ranking algorithms have been applied to the full Netflix data set. But in real life, we hope to make use of as much information as possible to make better recommendation. As shown in our experiments later, using full training data instead of sub-sampling (such as selecting a fixed number of rated items per user) achieves higher prediction and recommendation accuracy for the same test data.

To overcome this scalability issue, we propose two novel algorithms for solving problem (3.4), and both of them significantly reduce the time complexity over existing methods. If the input file is in the form of $|\Omega|$ pairwise comparisons, our proposed algorithm, Primal-CR, can reduce the time and space complexity from $O(|\Omega|r)$ to $O(|\Omega| + d_1\bar{d}_2r)$, where \bar{d}_2 is the average number of items compared by one user. If the input data is given as user-item ratings (e.g., Netflix, Yahoo-Music), the complexity is reduced from $O(d_1\bar{d}_2^2r)$ to $O(d_1\bar{d}_2r + d_1\bar{d}_2^2)$.

If the input file is given in ratings, we can further reduce the time complexity to $O(d_1\bar{d}_2r + d_1\bar{d}_2 \log \bar{d}_2)$ using exactly the same optimization algorithm but smarter ways to compute gradient and Hessian vector product. This time complexity is much smaller than the number of comparisons $|\Omega| = O(d_1\bar{d}_2^2)$, and we call this algorithm Primal-CR++.

We will first introduce Primal-CR in Section 3.4.2, and then present Primal-CR++ in Section 3.4.3.

3.4.2 Primal-CR: the proposed algorithm for pairwise input data

In the first setting, we consider the case where the pairwise comparisons $\{Y_{ijk} \mid (i, j, k) \in \Omega\}$ are given as input. To solve problem (3.4), we alternatively minimize U and V in the primal space (see Algorithm 2). First, we fix U and update V , and the subproblem for V while U is fixed can be written as follows:

$$V = \operatorname{argmin}_{V \in \mathbb{R}^{r \times d_2}} \left\{ \frac{\lambda}{2} \|V\|_F^2 + \sum_{(i,j,k) \in \Omega} L(Y_{ijk} \cdot u_i^T (v_j - v_k)) \right\} := f(V) \quad (3.6)$$

In [81], this subproblem is solved by stochastic dual coordinate descent, which requires $O(|\Omega|r)$ time and $O(|\Omega|)$ space complexity. Furthermore, the objective function decreases

Algorithm 2. Primal-CR / Primal-CR++: General Framework

Input: $\Omega, \{Y_{ijk} : (i, j, k) \in \Omega\}, \lambda \in \mathbb{R}^+$ ▷ for Primal-CR

Input: $M \in \mathbb{R}^{d_1 \times d_2}, \lambda \in \mathbb{R}^+$ ▷ for Primal-CR++

Output: $U \in \mathbb{R}^{r \times d_1}$ and $V \in \mathbb{R}^{r \times d_2}$

- 1: Randomly initialize U, V from Gaussian Distribution
- 2: **while** not converged **do**
- 3: **procedure** FIX U AND UPDATE V
- 4: **while** not converged **do**
- 5: Apply truncated Newton update (Algorithm 3)
- 6: **procedure** FIX V AND UPDATE U
- 7: **while** not converged **do**
- 8: Apply truncated Newton update (Algorithm 3)
- 9: **return** U, V ▷ recover score matrix X

for the dual problem sometimes does not imply the decrease of primal objective function value, which often results in slow convergence. We therefore propose to solve this subproblem for V using the primal truncated Newton method (Algorithm 3).

Newton method is a classical second-order optimization algorithm. For minimizing a vector-valued function $f(x)$, Newton method iteratively updates the solution by $x \leftarrow x - (\nabla^2 f(x))^{-1} \nabla f(x)$. However, the matrix inversion is usually hard to compute, so a truncated Newton method computes the update direction by solving the linear system $\nabla^2 f(x)a = \nabla f(x)$ up to a certain accuracy, usually using a linear conjugate gradient method. If we vectorized the problem for updating V in eq (3.6), the gradient is a (rd_2) -sized vector and the Hessian is an (rd_2) -by- (rd_2) matrix, so explicitly forming the Hessian is impossible. Below we discuss how to apply the truncated Newton method to solve our problem, and discuss efficient computations for each part.

Derivation of Gradient. When applying the truncated Newton method, the

Algorithm 3. Truncated Newton Update for V (same procedure can be used for updating U)

Input: Current solution U, V

Output: V

- 1: Compute $g = \text{vec}(\nabla f(V))$
 - 2: Let $H = \nabla^2 f(V)$ (do not explicitly compute H)
 - 3: **procedure** LINEAR CONJUGATE GRADIENT(g, F)
 - 4: Initialize $\delta_0 = 0$
 - 5: $r_0 = H\delta_0 - g, p_0 = -r_0$
 - 6: **for** $k = 0, 1, \dots, \text{maxiter}$ **do**
 - 7: Compute the Hessian-vector product $q = Hp_k$
 - 8: $\alpha_k = -r_k^T p_k / p_k^T q$
 - 9: $\delta_{k+1} = \delta_k + \alpha_k p_k$
 - 10: $r_{k+1} = r_k + \alpha_k q$
 - 11: **if** $\|r_{k+1}\|_2 < \|r_0\|_2 \cdot 10^{-2}$ **then**
 - 12: **break**
 - 13: $\beta_{k+1} = (r_{k+1}q) / p_k^T q$
 - 14: $p_{k+1} = -r_{k+1} + \beta_{k+1} p_k$
 - 15: **return** δ
 - 16: $V = V - s\delta$ (stepsize s found by line search)
 - 17: **return** U or V
-

gradient $\nabla f(V)$ is a $\mathbb{R}^{r \times d_2}$ matrix and can be computed explicitly:

$$\nabla f(V) = \sum_{i=1}^{d_1} \sum_{(j,k) \in \Omega_i} L'(Y_{ijk} \cdot u_i^T(v_j - v_k))(u_i e_j^T - u_i e_k^T) Y_{ijk} + \lambda V, \quad (3.7)$$

where $\nabla f(V) \in \mathbb{R}^{r \times d_2}$, $\Omega_i := \{(j, k) \mid (i, j, k) \in \Omega\}$ is the subset of pairs that associates with user i , and e_j is the indicator vector used to add the u_i vector to the j -th column of the output matrix. The first derivative for L_2 -hinge loss function (3.3) is

$$L'(a) = 2 \min(a - 1, 0) \quad (3.8)$$

For convenience, we define $g := \text{vec}(\nabla f(V))$ to be the vectorized form of gradient. One can easily see that computing g naively by going through all the pairwise comparisons (j, k) and adding up arrays is time-consuming and has $O(|\Omega|r)$ time complexity, which is the same with Collrank [81].

Fast computation for gradient Fortunately, we can reduce the time complexity to $O(|\Omega| + d_1 \bar{d}_2 r)$ by smartly rearranging the computations, so that the time is only linear to $|\Omega|$ and r , but not to $|\Omega|r$. The method is described below.

First, for each i , the first term of (3.7) can be represented by

$$\sum_{(j,k) \in \Omega_i} L'(Y_{ijk} \cdot u_i^T(v_j - v_k))(u_i e_j^T - u_i e_k^T) Y_{ijk} = \sum_{j \in \bar{d}_2(i)} t_j u_i e_j^T, \quad (3.9)$$

where $\bar{d}_2(i) := \{j \mid \exists k \text{ s.t. } (i, j, k) \in \Omega\}$ and t_j is some coefficient computed by summing over all the pairs in Ω_i . If we have t_j , the overall gradient can be computed by $O(\bar{d}_2(i)r)$ time for each i . To compute t_j , we first compute $u_i^T v_j$ for all $j \in \bar{d}_2(i)$ in $O(\bar{d}_2(i)r)$ time, and then go through all the (j, k) pairs while keep adding the coefficient related to this pair to t_j and t_k . Since there is no vector operations when we go through all pairs, this step only takes $O(\Omega_i)$ time. After getting all t_j , we can then conduct $\sum_{j \in \bar{d}_2(i)} t_j u_i e_j^T$ in $O(\bar{d}_2(i)r)$ time. Therefore, the overall complexity can be reduced to $O(|\Omega| + d_1 \bar{d}_2 r)$. The pseudo code is presented in Algorithm 4.

Derivation of Hessian-vector product Now we derive the Hessian $\nabla^2 f(V)$ for $f(V)$. We define $\nabla_j f(V) := \frac{\partial}{\partial v_j} f(V) \in \mathbb{R}^r$ and $\nabla_{j,k}^2 f(V) := \frac{\partial^2}{\partial v_j \partial v_k} f(V) \in \mathbb{R}^{r \times r}$ in the following

Algorithm 4. Primal-CR: efficient way to compute $\nabla f(V)$

Input: Ω , $\{Y_{ijk} : (i, j, k) \in \Omega\}$, $\lambda \in \mathbb{R}^+$, current variables U, V

Output: g, m $\triangleright g \in \mathbb{R}^{d_2 r}$ is the gradient for $f(V)$

- 1: Initialize $g = 0$ $\triangleright g \in \mathbb{R}^{r \times d_2}$
- 2: **for** $i = 1, 2, \dots, d_1$ **do**
- 3: **for all** $j \in \bar{d}_2(i)$ **do**
- 4: precompute $u_i^T v_j$ and store in a vector m_i
- 5: Initialize a zero array t of size d_2
- 6: **for all** $(j, k) \in \bar{d}_2(i)$ **do**
- 7: **if** $Y_{ijk}(m_i[j] - m_i[k]) < 1$ **then**
- 8: $s = 2(Y_{ijk}(m_i[j] - m_i[k]) - 1)$
- 9: $t[j] += Y_{ijk}s$
- 10: $t[k] -= Y_{ijk}s$ $\triangleright O(1)$ time per for loop iteration
- 11: **for all** $j \in \bar{d}_2(i)$ **do**
- 12: $g[:, j] += t[j] \cdot u_i$
- 13: $g = \text{vec}(g + \lambda V)$ \triangleright vectorize matrix $g \in \mathbb{R}^{r \times d_2}$
- 14: Form a sparse matrix $m = [m_1 \dots m_{d_1}]$ $\triangleright m$ can be reused later
- 15: **return** g, m

derivations. From the gradient derivation, we have

$$\nabla_j f(V) = \sum_{i: j \in \bar{d}_2(i)} \sum_{\substack{k \in \bar{d}_2(i) \\ k \neq j}} L'(Y_{ijk} \cdot u_i^T (v_j - v_k)) u_i Y_{ijk} + \lambda v_j.$$

Taking derivative again we can obtain

$$\nabla_{j,k}^2 f(V) = \begin{cases} \sum_{i:(j,k) \in \bar{d}_2(i)} L''(Y_{ijk} \cdot u_i^T (v_j - v_k)) (-u_i u_i^T) & \text{if } j \neq k \\ \sum_{i:j \in \bar{d}_2(i)} \sum_{k \in \bar{d}_2(i), k \neq j} L''(Y_{ijk} \cdot u_i^T (v_j - v_k)) u_i u_i^T + \lambda I_{r \times r} & \text{if } j = k \end{cases}$$

and the second derivative for

L_2 hinge loss function is given by:

$$L''(a) = \begin{cases} 2 & \text{if } a \leq 1 \\ 0 & \text{if } a > 1. \end{cases} \quad (3.10)$$

Note that if we write the full Hessian H as a $(d_2 r)$ by $(d_2 r)$ matrix, then $\nabla_{j,k}^2 f(V)$ is an $r \times r$ block in H , where there are totally d_2^2 of these blocks. In the CG update for solving $H^{-1}g$, we only need to compute $H \cdot a$ for some $a \in \mathbb{R}^{d_2 r}$. For convenience, we also partition this a into d_2 blocks, each subvector a_j has size r , so $a = [a_1; \dots; a_j]$. Similarly we can use subscript to denote the subarray $(H \cdot a)_j$ of the array $H \cdot a$, which becomes

$$(H \cdot a)_j = \sum_{k \neq j} \nabla_{j,k}^2 f(V) \cdot a_k + \nabla_{j,j}^2 f(V) \cdot a_j \quad (3.11)$$

$$= \lambda a_j + \sum_{i:j \in \bar{d}_2(i)} u_i \sum_{\substack{k \in \bar{d}_2(i) \\ k \neq j}} \quad (3.12)$$

$$+ L''(Y_{ijk} \cdot u_i^T (v_j - v_k)) (u_i^T a_j - u_i^T a_k). \quad (3.13)$$

Therefore, we have

$$\begin{aligned} H \cdot a &= \sum_j E_j (H \cdot a)_j \\ &= \lambda a + \sum_i \sum_{j \in \bar{d}_2(i)} E_j u_i \sum_{\substack{k \in \bar{d}_2(i) \\ k \neq j}} \end{aligned} \quad (3.14)$$

$$= L''(Y_{ijk} \cdot u_i^T (v_j - v_k)) (u_i^T a_j - u_i^T a_k) \quad (3.15)$$

where E_j is the projection matrix to the j -th block, indicating that we are only adding $(H \cdot a)_j$ to the j -th block of matrix, and setting 0 elsewhere.

Algorithm 5. Primal-CR: efficient way to compute Hessian vector product

Input: Ω , $\{Y_{ijk} : (i, j, k) \in \Omega\}$, $\lambda \in \mathbb{R}^+$, $a \in \mathbb{R}^{d_{2r}}$, m , U , V

Output: Ha

$\triangleright Ha \in \mathbb{R}^{d_{2r}}$ is needed in Linear CG

```

1:  $Ha = 0$   $\triangleright Ha \in \mathbb{R}^{d_{2r}}$ 
2: for  $i = 1, 2, \dots, d_1$  do
3:   for all  $j \in \bar{d}_2(i)$  do
4:     precompute  $u_i^T a_j$  and store it in array  $b$ 
5:   Initialize a zero array  $t$  of size  $d_2$ 
6:   for all  $(j, k) \in \bar{d}_2(i)$  do
7:     if  $Y_{ijk}(m_i[j] - m_i[k]) < 1.0$  then
8:        $s_{jk} = 2.0 \cdot (b[j] - b[k])$ 
9:        $t[j] += s_{jk}$ 
10:       $t[k] -= s_{jk}$   $\triangleright O(1)$  time per for loop iteration
11:   for all  $j \in \bar{d}_2(i)$  do
12:      $(Ha)[(p-1) \cdot r + 1 : p \cdot r] += t[j] \cdot u_i$ 
13: return  $Ha$ 

```

Fast computation for Hessian-vector product Similar to the case of gradient computation, using a naive way to compute $H \cdot a$ requires $O(|\Omega|r)$ time since we need to go through all the (i, j, k) tuples, and each of them requires $O(r)$ time. However, we can apply the similar trick in gradient computation to reduce the time complexity to $O(|\Omega| + d_1 \bar{d}_2 r)$ by pre-computing $u_i^T a_j$ and caching the coefficient using the array t . The detailed algorithm is given in Algorithm 5.

Note that in Algorithm 5, we can reuse the m (sparse array storing the current prediction) which has been pre-computed in the gradient computation (Algorithm 4), and that will cost only $O(d_1 \bar{d}_2)$ memory. Even without storing the m matrix, we can compute m in the loop of line 4 in Algorithm 5, which will not increase the overall computational

complexity.

Fix V and Update U After updating V by truncated Newton, we need to fix V and update U . The subproblem for U can be written as:

$$U = \operatorname{argmin}_{U \in \mathbb{R}^{r \times d_2}} \left\{ \frac{\lambda}{2} \|U\|_F^2 + \sum_{i=1}^{d_1} \sum_{(j,k) \in \bar{d}_2(i)} L(Y_{ijk} \cdot u_i^T (v_j - v_k)) \right\} \quad (3.16)$$

Since u_i , the i -th column of U , is independent from the rest of columns, equation 3.16 can be decomposed into d_1 independent problems for u_i :

$$u_i = \operatorname{argmin}_{u \in \mathbb{R}^r} \frac{\lambda}{2} \|u\|_2^2 + \sum_{(j,k) \in \bar{d}_2(i)} L(Y_{ijk} \cdot u^T (v_j - v_k)) := h(u) \quad (3.17)$$

Eq (3.17) is equivalent to an r -dimensional rankSVM problem. Since r is usually small, the problems are easy to solve. In fact, we can directly apply an efficient rankSVM algorithm proposed in [17] to solve each r -dimensional rankSVM problem. This algorithm requires $O(|\Omega_i| + r|\bar{d}_2(i)|)$ time for solving each subproblem with respect to u_i , so the overall complexity is $O(|\Omega| + rd_1\bar{d}_2)$ time per iteration.

Summary of time and space complexity When updating V , we first compute gradient by Algorithm 4, which takes $O(|\Omega| + d_1\bar{d}_2r)$ time, and each Hessian-vector product in 5 also takes the same time. The updates for U takes the same time complexity with updating V , so the overall time complexity is $O(|\Omega| + d_1\bar{d}_2r)$ per iteration. The whole algorithm only needs to store size $d_1 \times r$ and $d_2 \times r$ matrices for gradient and conjugate gradient method. The m matrix in Algorithm 4 is not needed, but in practice we find it can speedup the code by around 25%, and it only takes $d_1\bar{d}_2 \leq |\Omega|$ memory space (less than the input size). Therefore, our algorithm is very memory-efficient.

Before going to Primal-CR++, we discuss the time complexity of Primal-CR when the input data is the user-item rating matrix. Assume \bar{d}_2 is the averaged number of rated items per user, then there will be $|\Omega| = O(d_1\bar{d}_2^2)$ pairs, leading to $O(d_1\bar{d}_2^2 + d_1\bar{d}_2r)$ time complexity for Primal-CR. This is much better than the $O(d_1\bar{d}_2^2r)$ complexity for all the existing algorithms.

Algorithm 6. Primal-CR++: compute gradient part for $f(V)$

Input: $M \in \mathbb{R}^{d_1 \times d_2}$, $\lambda \in \mathbb{R}^+$, current U, V

Output: g

$\triangleright g \in \mathbb{R}^{d_2 r}$ is the gradient for $f(V)$

- 1: Initialize $g = 0$ $\triangleright g \in \mathbb{R}^{r \times d_2}$
 - 2: **for** $i = 1, 2, \dots, d_1$ **do**
 - 3: Let $\bar{d}_2 = |\bar{d}_2(i)|$ and $r_j = R_{i,j}$ for all j .
 - 4: Compute $m[j] = u_i^T v_j$ for all $j \in \bar{d}_2(i)$
 - 5: Sort $\bar{d}_2(i)$ according to the ascending order of m_i , so $m[\pi(1)] \leq \dots \leq m[\pi(\bar{d}_2)]$
 - 6: Initialize $s[1], \dots, s[L]$ and $c[1], \dots, c[L]$ with 0
 - 7: (Store in segment tree or Fenwick tree.)
 - 8: $p \leftarrow 1$
 - 9: **for all** $j = 1, \dots, \bar{d}_2$ **do**
 - 10: **while** $m[\pi(p)] \leq m_j + 1$ **do**
 - 11: $s[r_{\pi(p)}] += m[\pi(p)], c[r_{\pi(p)}] += 1$
 - 12: $p += 1$
 - 13: $S = \sum_{\ell \geq r_{\pi(p)}} s[\ell], C = \sum_{\ell \geq r_{\pi(p)}} c[\ell]$
 - 14: $t^+[\pi(j)] = 2(C \cdot (m[\pi(j)] + 1) - S)$
 - 15: Do another scan j from \bar{d}_2 to 1 to compute $t^-[\pi(j)]$ for all j
 - 16: $g[:, j] += (t^+[j] + t^-[j]) \cdot u_i$ for all j
 - 17: $g = \text{vec}(g + \lambda V)$ \triangleright vectorize matrix $g \in \mathbb{R}^{r \times d_2}$
 - 18: **return** g
-

3.4.3 Primal-CR++: the proposed algorithm for rating data

Now we discuss a more realistic scenario, where the input data is a rating matrix $\{R_{ij} \mid (i, j) \in \bar{\Omega}\}$ and $\bar{\Omega}$ is the observed set of user-item ratings. We assume there are only L levels of ratings, so $R_{ij} \in \{1, 2, \dots, L\}$. Also, we use $\bar{d}_2(i) := \{j \mid (i, j) \in \bar{\Omega}\}$ to denote the rated items for user i .

Given this data, the goal is to solve the collaborative ranking problem (3.4) with all the pairwise comparisons in the rating dataset as defined in (3.5). There are totally $O(d_1 \bar{d}_2^2)$ pairs, and the question is: Can we have an algorithm with near-linear time with respect to number of observed ratings $|\bar{\Omega}| = d_1 \bar{d}_2$? We answer this question in the affirmative by proposing Primal-CR++, a near-linear time algorithm for solving problem (3.4) with L2-hinge loss.

The algorithm of Primal-CR++ is exactly the same with Primal-CR, but we use a smarter algorithm to compute gradient and Hessian vector product in near-linear time, by exploiting the structure of the input data.

We first discuss how to speed up the gradient computation of (3.7), where the main computation is to compute (3.9) for each i . When the loss function is L2-hinge loss, we can explicitly write down the coefficients t_j in (3.9) by

$$t_j = \sum_{k \in \bar{d}_2(i)} 2(m_j - m_k - Y_{ijk}) I[Y_{ijk}(m_j - m_k) \leq 1], \quad (3.18)$$

where $m_j := u_i^T v_j$ and $I[\cdot]$ is an indicator function such that $I[a \leq b] = 1$ if $a \leq b$, and $I[a \leq b] = 0$ otherwise. By splitting the cases of $Y_{ijk} = 1$ and $Y_{ijk} = -1$, we get

$$\begin{aligned} t_j &= t_j^+ + t_j^- \\ &= \sum_{\substack{k \in \bar{d}_2(i) \\ m_k \leq m_j + 1, Y_{ijk} = -1}} 2(m_j - m_k + 1) + \sum_{\substack{k \in \bar{d}_2(i) \\ m_k \geq m_j - 1, Y_{ijk} = 1}} 2(m_j - m_k - 1). \end{aligned} \quad (3.19)$$

Assume the indexes in $\bar{d}_2(i)$ are sorted by the ascending order of m_j . Then we can scan from left to right, and maintain the current accumulated sum s_1, \dots, s_L and the current index counts c_1, \dots, c_L for each rating level. If the current pointer is p , then

$$s_\ell[p] = \sum_{j: m_j \leq p, R_{ij} = \ell} m_j \quad \text{and} \quad c_\ell[p] = |\{j : m_j \leq p, R_{ij} = \ell\}|.$$

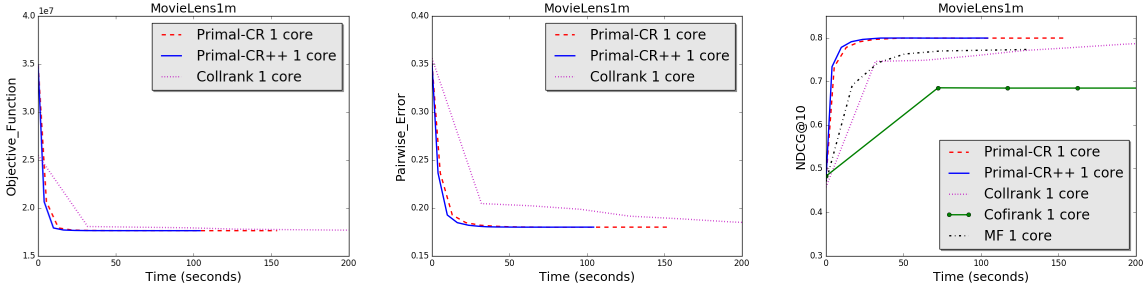


Figure 3.1. Comparing Primal-CR, Primal-CR++ and Collrank, MovieLens1m data, 200 ratings/user, rank 100, lambda = 5000

Since we scan from left to right, these numbers can be maintained in constant time at each step. Now assume we scan over the numbers $m_1 + 1, m_2 + 1, \dots$, then at each point we can compute

$$t_j^+ = \sum_{\ell=R_{i,j}+1}^L 2\{(m_j + 1)c_\ell[m_j + 1] - s_\ell[m_j + 1]\},$$

which can be computed in $O(L)$ time.

Although we observe that $O(L)$ time is already small in practice (since L usually smaller than 10), in the following we show there is a way to remove the dependency on L by using a simple Fenwick tree [27], F+tree [125] or segment tree. If we store the set $\{s_1, \dots, s_L\}$ in Fenwick tree, then each query of $\sum_{\ell \geq r} s_\ell$ can be done in $O(\log L)$ time, and since each step we only need to change one element into the set, the updating time is also $O(\log L)$. Note that t_j^- can be computed in the same way by scanning from largest m_j to the smallest one.

To sum up, the algorithm first computes all m_j in $O(\bar{d}_2 r)$ time, then sort these numbers using $O(\bar{d}_2 \log \bar{d}_2)$ time, and then compute t_j for all j using two linear scans in $O(\bar{d}_2 \log L)$ time. Here $\log L$ is dominated by $\log \bar{d}_2$ since L can be the number of unique rating levels in the current set $\bar{d}_2(i)$. Therefore, after computing this for all users $i = 1, \dots, d_1$, the time complexity for computing gradient is

$$O(d_1 \bar{d}_2 \log \bar{d}_2 + d_1 \bar{d}_2 r) = O(|\bar{\Omega}|(\log \bar{d}_2 + r)).$$

A similar procedure can also be used for computing the Hessian-vector product, and the computation of updating U with fixed V is simpler since the problem becomes

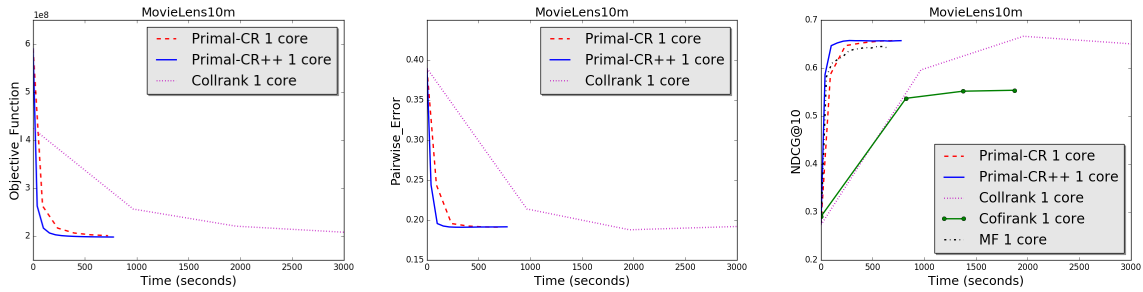


Figure 3.2. Comparing Primal-CR, Primal-CR++ and Collrank, MovieLens10m data, 500 ratings/user, rank 100, lambda = 7000

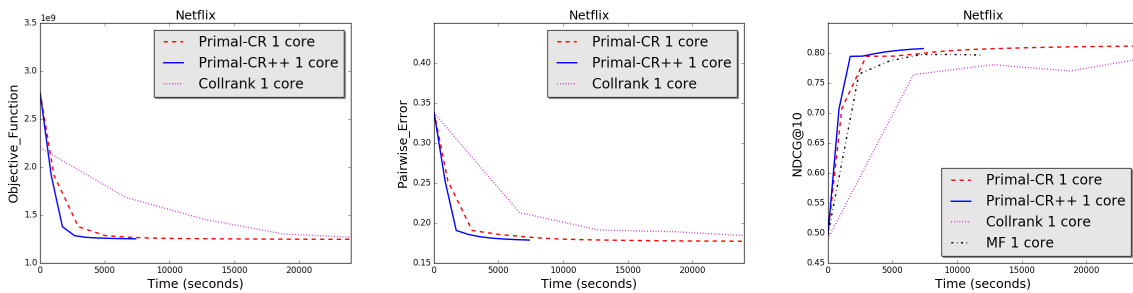


Figure 3.3. Comparing Primal-CR, Primal-CR++ and Collrank, Netflix data, 200 ratings/user, rank 100, lambda = 10000

decomposable to d_1 independent problems, see eq (3.17). Due to the page limit we omit the details here; interesting readers can check our code on github.

Compared with the classical matrix factorization, where both ALS and SGD requires $O(|\bar{\Omega}|r)$ time per iteration [61], our algorithm has almost the same complexity, since $\log \bar{d}_2$ is usually smaller than r (typically $r = 100$). Also, since all the temporary memory when computing user i can be released immediately, the only memory cost is still the same with Primal-CR++, which is $O(d_1r + d_2r)$.

3.4.4 Parallelization

Updating U while fixing V can be parallelized easily because each column of U is independent and we can actually solve d_1 independent subproblems at the same time. For the other side, updating V while fixing U can also be parallelized by parallelizing “computing g ” part and “computing Ha ” part respectively. We implemented the algorithm using parallel computing techniques in Julia by computing g and Ha distributedly and summing

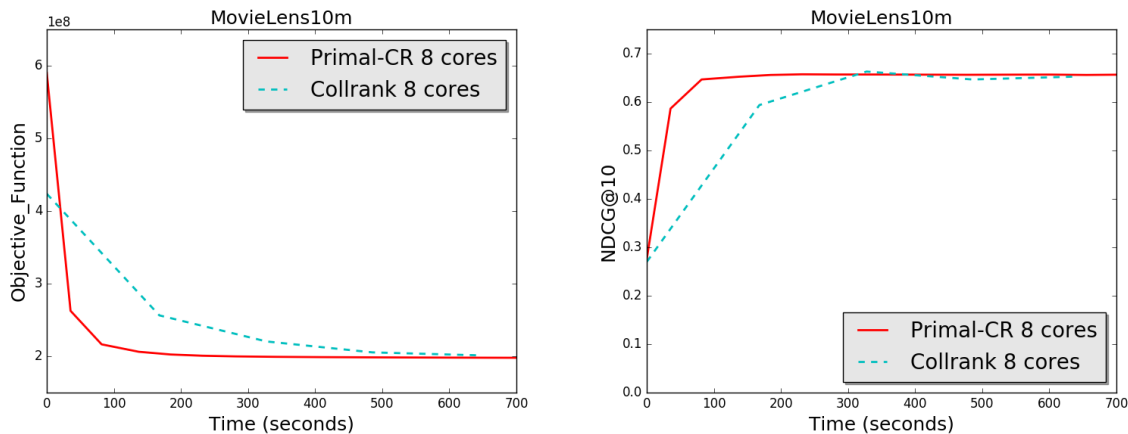


Figure 3.4. Comparing parallel version of Primal-CR and Collrank, MovieLens10m data, 500 ratings/user, rank 100, lambda = 7000

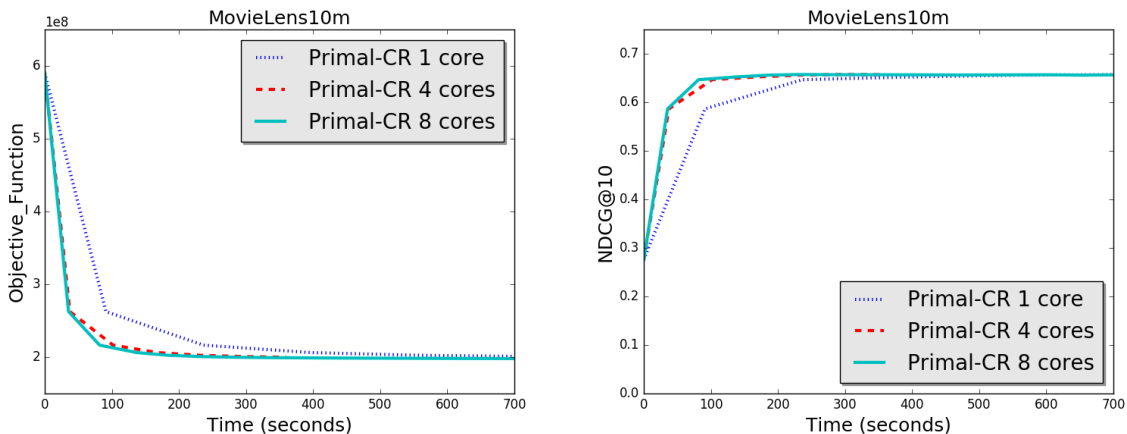


Figure 3.5. Speedup of Primal-CR, MovieLens10m data, 500 ratings/user, rank 100, lambda = 7000

them up in the end. We show in section 3.5.2 that our parallel version of the proposed new algorithm works better than the paralleled version of Collrank algorithm [81].

3.5 Experiments

In this section, we test the performance of our proposed algorithms Primal-CR and Primal-CR++ on real world datasets, and compare with existing methods. All experiments are conducted on the UC Davis Illidan server with an Intel Xeon E5-2640 2.40GHz CPU and 64G RAM. We compare the following methods:

- Primal-CR and Primal-CR++: our proposed methods implemented in Julia. ²
- Collrank: the collaborative ranking algorithm proposed in [81]. We use the C++ code released by the authors, and they parallelized their algorithm using OpenMP.
- Cofrank: the classical collaborative ranking algorithm proposed in [114]. We use the C++ code released by the authors.
- MF: the classical matrix factorization model in (3.1) solved by SGD [61].

We used three data sets (MovieLens1m, Movielens10m, Netflix data) to compare these algorithms. The dataset statistics are summarized in Table 3.1. The regularization parameter λ used for each datasets are chosen by a random sampled validation set. For the pair-wise based algorithms, we covert the ratings into pair-wise comparisons, by saying that item j is preferred over item k by user i if user i gives a higher rating to item j over item k , and there will be no pair between two items if they have the same rating.

We compare the algorithms in the following three different ways:

- Objective function: since Collrank, Primal-CR, Primal-CR++ have the same objective function, we can compare the convergence speed in terms of the objective function (3.4) with squared hinge loss.
- Predicted pairwise error: the proportion of pairwise preference comparisons that we predicted correctly out of all the pairwise comparisons in the testing data:

$$\text{pairwise error} = \frac{1}{|\mathcal{T}|} \sum_{\substack{(i,j,k) \in \mathcal{T} \\ Y_{ijk}=1}} \mathbb{1}(X_{ij} > X_{ik}), \quad (3.20)$$

where \mathcal{T} represents the test data set and $|\mathcal{T}|$ denotes the size of test data set.

- NDCG@ k : a standard performance measure of ranking, defined as:

$$\text{NDCG}@k = \frac{1}{d_1} \sum_{i=1}^{d_1} \frac{\text{DCG}@k(i, \pi_i)}{\text{DCG}@k(i, \pi_i^*)}, \quad (3.21)$$

where i represents i -th user and

$$\text{DCG}@k(i, \pi_i) = \sum_{l=1}^k \frac{2^{M_i \pi_i(l)} - 1}{\log_2(l+1)}. \quad (3.22)$$

²Our code is available on <https://github.com/wuliwei9278/ml-1m>.

In the DCG definition, $\pi_i(l)$ represents the index of the l -th ranked item for user i in test data based on the score matrix $X = U^T V$ generated, M is the rating matrix and M_{ij} is the rating given to item j by user i . π_i^* is the ordering provided by the underlying ground truth of the rating.

3.5.1 Compare single thread versions using the same subsamples

Since Collrank cannot scale to the full dataset of Movielens10m and Netflix, we sub-sample data using the same approach in their paper [81] and compare all the methods using the smaller training sets. More specifically, for each data set, we subsampled N ratings for training data and used the rest of ratings as test data. For this subsampled data, we discard users with less than $N + 10$ ratings, since we need at least 10 ratings for test data to compute the NDCG@10.

As shown in Figure 3.1, 3.2, 3.3, both Primal-CR and Primal-CR++ perform considerably better than the existing Collrank algorithm. As data size increases, the performance gap becomes larger. As one can see, for Netflix data where $N = 200$, the speedup is more than 10 times compared to Collrank.

For Cofrank, we observe that it is even slower than Collrank, which confirms the experiments conducted in [81]. Furthermore, Cofrank cannot scale to larger datasets, so we omit the results in Figure 3.2 and 3.3.

We also include the classical matrix factorization algorithm in the NDCG comparisons. As shown in our complexity analysis, our proposed algorithms are competitive with MF in terms of speed, and MF is much faster than other collaborative ranking algorithms. Also, we observe that MF converges to a slightly worse solution in MovieLens10m and Netflix datasets, and converges to a much worse solution in MovieLens1m. The reason is that MF minimizes a simple mean square error, while our algorithms are minimizing ranking loss. Based on the experimental results, our algorithm Primal-CR++ should be able to replace MF in many real world recommender systems.

Table 3.1. Datasets used for experiments

	MovieLens1m	MovieLens10m	Netflix
Users	6,040	71,567	2,649,430
Items	3,952	65,134	17,771
Ratings	1,000,209	10,000,054	99,072,112
λ	5,000	7,000	10,000

Table 3.2. Scalability of Primal-CR and Collrank on MovieLens10m

# cores	1	4	8
Speedup for Primal-CR	1x	2.46x	3.12x
Speedup for Collrank	1x	2.95x	3.47x

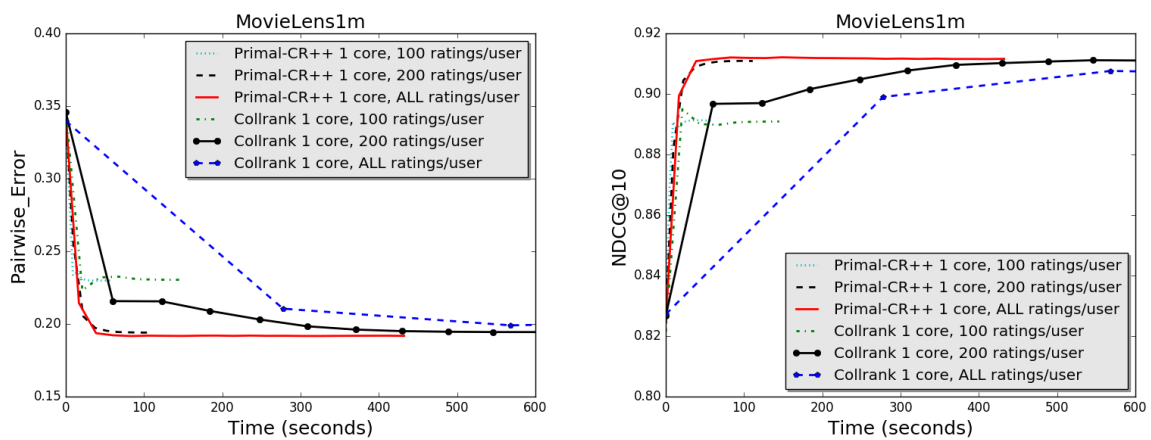


Figure 3.6. Varies number of ratings per user in training data, MovieLens1m data, rank 100

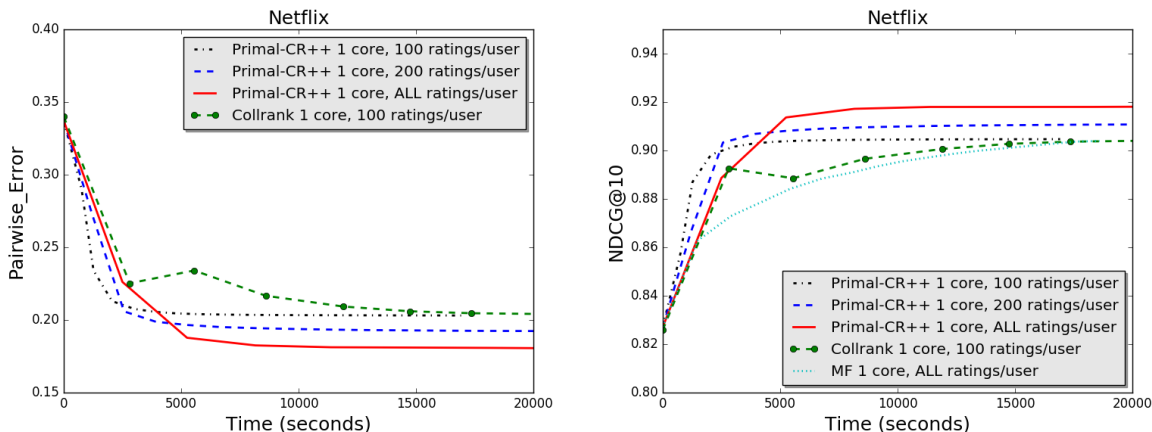


Figure 3.7. Varies number of ratings per user in training data, Netflix data, rank 100

3.5.2 Compare parallel versions

Since Collrank can be implemented in a parallel fashion, we also implemented the parallel version of our algorithm in Julia. We want to show our algorithm scales up well and is still much faster than Collrank in the multi-core shared memory setting. As shown in Figure 3.5, Primal-CR is still much faster than Collrank when 8 cores are used. Comparing our Primal-CR algorithm in 1 core, 4 cores and 8 cores on the same machine in Figure 3.5, the speedup is desirable. The speedup of Primal-CR and Collrank is summarized in the Table 3.2. One can see from the table that the speedup of our Primal-CR algorithm is comparable to Collrank.

3.5.3 Performance using Full Training Data

Due to the $O(|\Omega|k)$ complexity, existing algorithms cannot deal with large number of pairs, so they always sub-sample a limited number of pairs per user when solving MovieLens10m or Netflix data. For example, for Collrank, the authors fixed number of ratings per user in training as N and only reported N up to 100 for Netflix data. When we tried to apply their code for $N = 200$, the algorithm gets very slow and reports memory error for $N = 500$.

Using our algorithm, we have the ability to solve the full Netflix problem, so a natural question to ask is: Does using more training data help us predict and recommend better? The answer is yes! We conduct the following experiments to verify this: For all the users with more than 20 ratings, we randomly choose 10 ratings as test data and out of the

rest ratings we randomly choose up to C ratings per user as training data. One can see in Figure 3.6, for the same test data, more training data leads to better prediction performance in terms of pairwise error and NDCG. Using all available ratings ($C = d_2$) gives lowest pairwise error and highest NDCG@10, using up to 200 ratings per user ($C = 200$) gives second lowest pairwise error and second highest NDCG@10, and using up to 100 ratings per user ($C = 100$) has the highest pairwise error and lowest NDCG@10. Similar phenomenon is observed for Netflix data in Figure 3.7. Collrank code does not work for $C = 200$ and $C = d_2$ and even for $C = 100$, it takes more than 20,000 secs to converge while our Primal-CR++ takes less than 5,000 secs for the full Netflix data. The speedup of our algorithm will be even more for a larger C or larger data size d_1 and d_2 . We tried to create input file without subsampling for Collrank, we created 344GB input data file and Collrank reported memory error message "Segmentation Fault". We also tried $C = 200$, still got the same error message. It is possible to implement Collrank algorithm by directly working on the rating data, but the time complexity remains the same, so it is clear that our proposed Primal-CR and Primal-CR++ algorithms are much faster.

To the best of our knowledge, our algorithm is the first ranking-based algorithm that can scale to full Netflix data set using a single core, and without sub-sampling. Our proposed algorithm makes the Collaborative Ranking Model in (3.4) a clear better choice for large-scale recommendation system over standard Matrix Factorization techniques, since we have the same scalability but achieve much better accuracy. Also, our experiments suggest that in practice, when we are given a set of training data, we should try to use all the training data instead of doing sub-sampling as existing algorithms do, and only Primal-CR and Primal-CR++ can scale up to all the ratings.

3.6 Conclusions

We considered the collaborative ranking problem setting in which a low-rank matrix is fitted to the data in the form of pairwise comparisons or numerical ratings. We proposed our new optimization algorithms Primal-CR and Primal-CR++ where the time complexity is much better than all the existing approaches. We showed that our algorithms are much faster

than state-of-the-art collaborative ranking algorithms on real data sets (MovieLens1m, Movielens10m and Netflix) using same subsampling scheme, and moreover our algorithm is the only one that can scale to the full Movielens10m and Netflix data. We observed that our algorithm has the same efficiency with matrix factorization, while achieving better NDCG since we minimize ranking loss. As a result, we expect our algorithm to be able to replace matrix factorization in many real applications.

Chapter 4

SQL-Rank: A Listwise Approach to Collaborative Ranking

4.1 Introduction

We study a novel approach to collaborative ranking—the personalized ranking of items for users based on their observed preferences—through the use of listwise losses, which are dependent only on the observed rankings of items by users. We propose the SQL-Rank algorithm, which can handle ties and missingness, incorporate both explicit ratings and more implicit feedback, provides personalized rankings, and is based on the relative rankings of items. To better understand the proposed contributions, let us begin with a brief history of the topic.

4.1.1 A brief history of collaborative ranking

Recommendation systems, found in many modern web applications, movie streaming services, and social media, rank new items for users and are judged based on user engagement (implicit feedback) and ratings (explicit feedback) of the recommended items. A high-quality recommendation system must understand the popularity of an item and infer a user’s specific preferences with limited data. Collaborative filtering, introduced in [44], refers to the use of an entire community’s preferences to better predict the preferences of an individual (see [95] for an overview). In systems where users provide ratings of items, collaborative filtering can be approached as a point-wise prediction task, in which we

attempt to predict the unobserved ratings [80]. Low rank methods, in which the rating distribution is parametrized by a low rank matrix (meaning that there are a few latent factors) provides a powerful framework for estimating ratings [59, 73]. There are several issues with this approach. One issue is that the feedback may not be representative of the unobserved entries due to a sampling bias, an effect that is prevalent when the items are only ‘liked’ or the feedback is implicit because it is inferred from user engagement. Augmenting techniques like weighting were introduced to the matrix factorization objective to overcome this problem [48, 49]. Many other techniques are also introduced [55, 112, 120]. Another methodology worth noting is the CofiRank algorithm of [113] which minimizes a convex surrogate of the normalized discounted cumulative gain (NDCG). The pointwise framework has other flaws, chief among them is that in recommendation systems we are not interested in predicting ratings or engagement, but rather we must rank the items.

Ranking is an inherently relative exercise. Because users have different standards for ratings, it is often desirable for ranking algorithms to rely only on relative rankings and not absolute ratings. A ranking loss is one that only considers a user’s relative preferences between items, and ignores the absolute value of the ratings entirely, thus deviating from the pointwise framework. Ranking losses can be characterized as pairwise and listwise. A pairwise method decomposes the objective into pairs of items j, k for a user i , and effectively asks ‘did we successfully predict the comparison between j and k for user i ?’. The comparison is a binary response—user i liked j more than or less than k —with possible missing values in the event of ties or unobserved preferences. Because the pairwise model has cast the problem in the classification framework, then tools like support vector machines were used to learn rankings; [54] introduces rankSVM and efficient solvers can be found in [17]. Much of the existing literature focuses on learning a single ranking for all users, which we will call simple ranking [2, 28, 77]. This work will focus on the personalized ranking setting, in which the ranking is dependent on the user.

Pairwise methods for personalized ranking have seen great advances in recent years, with the AltSVM algorithm of [81], Bayesian personalized ranking (BPR) of [91], and the near linear-time algorithm of [115]. Nevertheless, pairwise algorithms implicitly assume

that the item comparisons are independent, because the objective can be decomposed where each comparison has equal weight. Listwise losses instead assign a loss, via a generative model, to the entire observed ranking, which can be thought of as a permutation of the m items, instead of each comparison independently. The listwise permutation model, introduced in [15], can be thought of as a weighted urn model, where items correspond to balls in an urn and they are sequentially plucked from the urn with probability proportional to $\phi(X_{ij})$ where X_{ij} is the latent score for user i and item j and ϕ is some non-negative function. They proposed to learn rankings by optimizing a cross entropy between the probability of k items being at the top of the ranking and the observed ranking, which they combine with a neural network, resulting in the ListNet algorithm. [100] applies this idea to collaborative ranking, but uses only the top-1 probability because of the computational complexity of using top-k in this setting. This was extended in [50] to incorporate neighborhood information. [121] instead proposes a maximum likelihood framework that uses the permutation probability directly, which enjoyed some empirical success.

Very little is understood about the theoretical performance of listwise methods. [15] demonstrates that the listwise loss has some basic desirable properties such as monotonicity, i.e. increasing the score of an item will tend to make it more highly ranked. [65] studies the generalizability of several listwise losses, using the local Rademacher complexity, and found that the excess risk could be bounded by a $1/\sqrt{n}$ term (recall, n is the number of users). Two main issues with this work are that no dependence on the number of items is given—it seems these results do not hold when m is increasing—and the scores are not personalized to specific users, meaning that they assume that each user is an independent and identically distributed observation. A simple open problem is: can we consistently learn preferences from a single user’s data if we are given item features and we assume a simple parametric model? ($n = 1, m \rightarrow \infty$.)

4.1.2 Contributions of this work

We can summarize the shortcomings of the existing work: current listwise methods for collaborative ranking rely on the top-1 loss, algorithms involving the full permutation probability are computationally expensive, little is known about the theoretical performance

of listwise methods, and few frameworks are flexible enough to handle explicit and implicit data with ties and missingness. This chapter addresses each of these in turn by proposing and analyzing the SQL-rank algorithm.

- We propose the SQL-Rank method, which is motivated by the permutation probability, and has advantages over the previous listwise method using cross entropy loss.
- We provide an $O(\text{iter} \cdot (|\Omega|r))$ linear algorithm based on stochastic gradient descent, where Ω is the set of observed ratings and r is the rank.
- The methodology can incorporate both implicit and explicit feedback, and can gracefully handle ties and missing data.
- We provide a theoretical framework for analyzing listwise methods, and apply this to the simple ranking and personalized ranking settings, highlighting the dependence on the number of users and items.

4.2 Methodology

4.2.1 Permutation probability

The permutation probability, [15], is a generative model for the ranking parametrized by latent scores. First assume there exists a ranking function that assigns scores to all the items. Let's say we have m items, then the scores assigned can be represented as a vector $s = (s_1, s_2, \dots, s_m)$. Denote a particular permutation (or ordering) of the m items as π , which is a random variable and takes values from the set of all possible permutations S_m (the symmetric group on m elements). π_1 denotes the index of highest ranked item and π_m is the lowest ranked. The probability of obtaining π is defined to be

$$P_s(\pi) := \prod_{j=1}^m \frac{\phi(s_{\pi_j})}{\sum_{l=j}^m \phi(s_{\pi_l})}, \quad (4.1)$$

where $\phi(\cdot)$ is an increasing and strictly positive function. An interpretation of this model is that each item is drawn without replacement with probability proportional to $\phi(s_i)$ for item i in each step. One can easily show that $P_s(\pi)$ is a valid probability distribution, i.e. $\sum_{\pi \in S_m} P_s(\pi) = 1, P_s(\pi) > 0, \forall \pi$. Furthermore, this definition of permutation probability enjoys several favorable properties (see [15]). For any permutation π if you swap two

elements ranked at $i < j$ generating the permutation π' ($\pi'_i = \pi_j$, $\pi'_j = \pi_i$, $\pi_k = \pi'_k$, $k \neq i, j$), if $s_{\pi_i} > s_{\pi_j}$ then $P_s(\pi) > P_s(\pi')$. Also, if permutation π satisfies $s_{\pi_i} > s_{\pi_{i+1}}$, $\forall i$, then we have $\pi = \arg \max_{\pi' \in S_m} P_s(\pi')$. Both of these properties can be summarized: larger scores will tend to be ranked more highly than lower scores. These properties are required for the negative log-likelihood to be considered sound for ranking [121].

In recommendation systems, the top ranked items can be more impactful for the performance. In order to focus on the top k ranked items, we can compute the partial-ranking marginal probability,

$$P_s^{(k, \bar{m})}(\pi) = \prod_{j=1}^{\min\{k, \bar{m}\}} \frac{\phi(s_{\pi_j})}{\sum_{l=j}^{\bar{m}} \phi(s_{\pi_l})}. \quad (4.2)$$

It is a common occurrence that only a proportion of the m items are ranked, and in that case we will allow $\bar{m} \leq m$ to be the number of observed rankings (we assume that $\pi_1, \dots, \pi_{\bar{m}}$ are the complete list of ranked items). When $k = 1$, the first summation vanishes and top-1 probability can be calculated straightforwardly, which is why $k = 1$ is widely used in previous listwise approaches for collaborative ranking. Counter-intuitively, we demonstrate that using a larger k tends to improve the ranking performance.

We see that computing the likelihood loss is linear in the number of ranked items, which is in contrast to the cross-entropy loss used in [15], which takes exponential time in k . The cross-entropy loss is also not sound, i.e. it can rank worse scoring permutations more highly, but the negative log-likelihood is sound. We will discuss how we can deal with ties in the following subsection, namely, when the ranking is derived from ratings and multiple items receive the same rating, then there is ambiguity as to the order of the tied items. This is a common occurrence when the data is implicit, namely the output is whether the user engaged with the item or not, yet did not provide explicit feedback. Because the output is binary, the cross-entropy loss (which is based on top- k probability with k very small) will perform very poorly because there will be many ties for the top ranked items. To this end, we propose a collaborative ranking algorithm using the listwise likelihood that can accommodate ties and missingness, which we call Stochastic Queuing Listwise Ranking, or SQL-Rank.

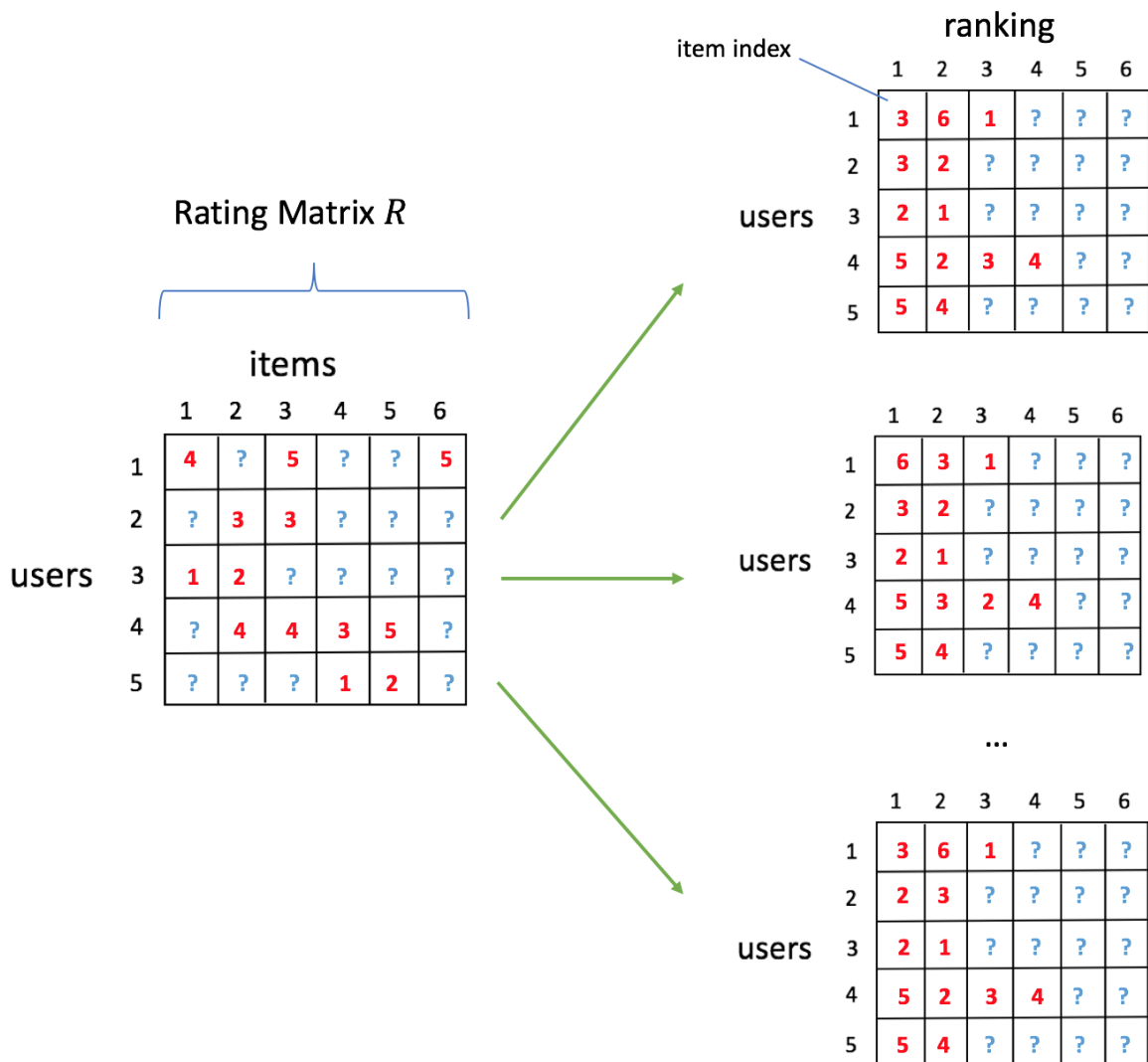


Figure 4.1. Demonstration of Stochastic Queuing Process—the rating matrix R (left) generates multiple possible rankings Π 's (right), $\Pi \in \mathcal{S}(R, \Omega)$ by breaking ties randomly.

4.2.2 Deriving objective function for SQL-Rank

The goal of collaborative ranking is to predict a personalized score X_{ij} that reflects the preference level of user i towards item j , where $1 \leq i \leq n$ and $1 \leq j \leq m$. It is reasonable to assume the matrix $X \in \mathbb{R}^{n \times m}$ to be low rank because there are only a small number of latent factors contributing to users' preferences. The input data is given in the form of “user i gives item j a relevance score R_{ij} ”. Note that for simplicity we assume all the users

have the same number \bar{m} of ratings, but this can be easily generalized to the non-uniform case by replacing \bar{m} with m_i (number of ratings for user i).

With our scores X and our ratings R , we can specify our collaborative ranking model using the permutation probability (4.2). Let Π_i be a ranking permutation of items for user i (extracted from R), we can stack Π_1, \dots, Π_n , row by row, to get the permutation matrix $\Pi \in \mathbb{R}^{n \times m}$. Assuming users are independent with each other, the probability of observing a particular Π given the scoring matrix X can be written as

$$P_X^{(k, \bar{m})}(\Pi) = \prod_{i=1}^n P_{X_i}^{(k, \bar{m})}(\Pi_i). \quad (4.3)$$

We will assume that $\log \phi(x) = 1/(1 + \exp(-x))$ is the sigmoid function. This has the advantage of bounding the resulting weights, $\phi(X_{ij})$, and maintaining their positivity without adding additional constraints.

Typical rating data will contain many ties within each row. In such cases, the permutation Π is no longer unique and there is a set of permutations that coincides with rating because with any candidate Π we can arbitrarily shuffle the ordering of items with the same relevance scores to generate a new candidate matrix Π' which is still valid (see Figure 4.1). We denote the set of valid permutations as $\mathcal{S}(R, \Omega)$, where Ω is the set of all pairs (i, j) such that $R_{i,j}$ is observed. We call this shuffling process the *Stochastic Queuing Process*, since one can imagine that by permuting ties we are stochastically queuing new Π 's for future use in the algorithm.

The probability of observing R therefore should be defined as $P_X^{(k, \bar{m})}(R) = \sum_{\Pi \in \mathcal{S}(R, \Omega)} P_X(\Pi)$. To learn the scoring matrix X , we can naturally solve the following maximum likelihood estimator with low-rank constraint:

$$\min_{X \in \mathcal{X}} -\log \sum_{\Pi \in \mathcal{S}(R, \Omega)} P_X^{(k, \bar{m})}(\Pi), \quad (4.4)$$

where \mathcal{X} is the structural constraint of the scoring matrix. To enforce low-rankness, we use the nuclear norm regularization $\mathcal{X} = \{X : \|X\|_* \leq r\}$.

Eq (4.4) is hard to optimize since there is a summation inside the log. But by Jensen's inequality and convexity of $-\log$ function, we can move the summation outside log and

obtain an upper bound of the original negative log-likelihood, leading to the following optimization problem:

$$\min_{X \in \mathcal{X}} - \sum_{\Pi \in \mathcal{S}(R, \Omega)} \log P_X^{(k, \bar{m})}(\Pi) \quad (4.5)$$

This upper bound is much easier to optimize and can be solved using Stochastic Gradient Descent (SGD).

Next we discuss how to apply our model for explicit and implicit feedback settings. In the explicit feedback setting, it is assumed that the matrix R is partially observed and the observed entries are explicit ratings in a range (e.g., 1 to 5). We will show in the experiments that $k = \bar{m}$ (using the full list) leads to the best results. [50] also observed that increasing k is useful for their cross-entropy loss, but they were not able to increase k since their model has time complexity exponential to k .

In the implicit feedback setting each element of R_{ij} is either 1 or 0, where 1 means positive actions (e.g., click or like) and 0 means no action is observed. Directly solving (4.5) will be expensive since $\bar{m} = m$ and the computation will involve all the mn elements at each iteration. Moreover, the 0's in the matrix could mean either a lower relevance score or missing, thus should contribute less to the objective function. Therefore, we adopt the idea of negative sampling [71] in our list-wise formulation. For each user (row of R), assume there are \tilde{m} 1's, we then sample $\rho\tilde{m}$ unobserved entries uniformly from the same row and append to the back of the list. This then becomes the problem with $\bar{m} = (1 + \rho)\tilde{m}$ and then we use the same algorithm in explicit feedback setting to conduct updates. We then repeat the sampling process at the end of each iteration, so the update will be based on different set of 0's at each time.

4.2.3 Non-convex implementation

Despite the advantage of the objective function in equation (4.5) being convex, it is still not feasible for large-scale problems since the scoring matrix $X \in \mathbb{R}^{n \times m}$ leads to high computational and memory cost. We follow a common trick to transform (4.5) to the non-convex form by replacing $X = U^T V$: with $U \in \mathbb{R}^{r \times n}$, $V \in \mathbb{R}^{r \times m}$ so that the objective

Algorithm 7. SQL-Rank: General Framework

Input: $\Omega, \{R_{ij} : (i, j) \in \Omega\}, \lambda \in \mathbb{R}^+, ss, rate, \rho$

Output: $U \in \mathbb{R}^{r \times n}$ and $V \in \mathbb{R}^{r \times m}$

Randomly initialize U, V from Gaussian Distribution

repeat

Generate a new permutation matrix Π ▷ see alg 8

Apply gradient update to U while fixing V

Apply gradient update to V while fixing U ▷ see alg 13

until performance for validation set is good

return U, V ▷ recover score matrix X

Algorithm 8. Stochastic Queuing Process

Input: $\Omega, \{R_{ij} : (i, j) \in \Omega\}, \rho$

Output: $\Pi \in \mathbb{R}^{n \times m}$

for $i = 1$ **to** n **do**

Sort items based on observed relevance levels R_i

Form Π_i based on indices of items in the sorted list

Shuffle Π_i for items within the same relevance level

if Dataset is implicit feedback **then**

Uniformly sample $\rho \tilde{m}$ items from unobserved items

Append sampled indices to the back of Π_i

Stack Π_i as rows to form matrix Π

Return Π ▷ Used later to compute gradient

is,

$$\sum_{\Pi \in \mathcal{S}(R, \Omega)} - \underbrace{\sum_{i=1}^n \sum_{j=1}^{\tilde{m}} \log \frac{\phi(u_i^T v_{\Pi_{ij}})}{\sum_{l=j}^{\tilde{m}} \phi(u_i^T v_{\Pi_{il}})}}_{f(U, V)} + \frac{\lambda}{2} (\|U\|_F^2 + \|V\|_F^2),$$

where u_i, v_j are columns of U, V respectively. We apply stochastic gradient descent to solve this problem. At each step, we choose a permutation matrix $\Pi \in \mathcal{S}(R, \Omega)$ using

the stochastic queuing process (Algorithm 8) and then update U, V by $\nabla f(U, V)$. For example, the gradient with respect to V is ($g = \log \phi$ is the sigmoid function),

$$\frac{\partial f}{\partial v_j} = \sum_{i \in \Omega_j} \sum_{t=1}^{\bar{m}} \left\{ -g'(u_i^T v_t) u_i + \frac{\mathbb{1}(\text{rank}_i(j) \geq t) \phi(u_i^T v_j)}{\sum_{l=t}^{\bar{m}} \phi(u_i^T v_{\Pi_{il}})} g'(u_i^T v_j) u_i \right\}$$

where Ω_j denotes the set of users that have rated the item j and $\text{rank}_i(j)$ is a function gives the rank of the item j for that user i . Because g is the sigmoid function, $g' = g \cdot (1 - g)$. The gradient with respect to U can be derived similarly.

As one can see, a naive way to compute the gradient of f requires $O(n\bar{m}^2r)$ time, which is very slow even for one iteration. However, we show in Algorithm 12 (in the appendix) that there is a smart way to re-arranging the computation so that $\nabla_V f(U, V)$ can be computed in $O(n\bar{m}r)$ time, which makes our SQL-Rank a linear-time algorithm (with the same per-iteration complexity as classical matrix factorization).

4.3 Experiments

In this section, we compare our proposed algorithm (SQL-Rank) with other state-of-the-art algorithms on real world datasets. Note that our algorithm works for both implicit feedback and explicit feedback settings. In the implicit feedback setting, all the ratings are 0 or 1; in the explicit feedback setting, explicit ratings (e.g., 1 to 5) are given but only to a subset of user-item pairs. Since many real world recommendation systems follow the implicit feedback setting (e.g., purchases, clicks, or checkins), we will first compare SQL-Rank on implicit feedback datasets and show it outperforms state-of-the-art algorithms. Then we will verify that our algorithm also performs well on explicit feedback problems. All experiments are conducted on a server with an Intel Xeon E5-2640 2.40GHz CPU and 64G RAM.

4.3.1 Implicit Feedback

In the implicit feedback setting we compare the following methods:

- SQL-Rank: our proposed algorithm implemented in Julia ¹.

¹<https://github.com/wuliwei9278/SQL-Rank>

- Weighted-MF: the weighted matrix factorization algorithm by putting different weights on 0 and 1’s [48, 49].
- BPR: the Bayesian personalized ranking method motivated by MLE [91]. For both Weighted-MF and BPR, we use the C++ code by Quora ².

Note that other collaborative ranking methods such as Pirmal-CR++ [115] and List-MF [100] do not work for implicit feedback data, and we will compare with them later in the explicit feedback experiments. For the performance metric, we use precision@ k for $k = 1, 5, 10$ defined by

$$\text{precision@}k = \frac{\sum_{i=1}^n |\{1 \leq l \leq k : R_{i\Pi_{il}} = 1\}|}{n \cdot k}, \quad (4.6)$$

where R is the rating matrix and Π_{il} gives the index of the l -th ranked item for user i among all the items not rated by user i in the training set.

We use rank $r = 100$ and tune regularization parameters for all three algorithms using a random sampled validation set. For Weighted-MF, we also tune the confidence weights on unobserved data. For BPR and SQL-Rank, we fix the ratio of subsampled unobserved 0’s versus observed 1’s to be 3 : 1, which gives the best performance for both BPR and SQL-rank in practice.

We experiment on the following four datasets. Note that the original data of Movielens1m, Amazon and Yahoo-music are ratings from 1 to 5, so we follow the procedure in [91, 124] to preprocess the data. We transform ratings of 4, 5 into 1’s and the rest entries (with rating 1, 2, 3 and unknown) as 0’s. Also, we remove users with very few 1’s in the corresponding row to make sure there are enough 1’s for both training and testing. For Amazon, Yahoo-music and Foursquare, we discard users with less than 20 ratings and randomly select 10 1’s as training and use the rest as testing. Movielens1m has more ratings than others, so we keep users with more than 60 ratings, and randomly sample 50 of them as training.

- Movielens1m: a popular movie recommendation data with 6,040 users and 3,952 items.

²<https://github.com/quora/qmf>

Table 4.1. Comparing implicit feedback methods on various datasets.

DATASET	METHOD	P@1	P@5	P@10
MOVIELENS1M	SQL-RANK	0.73685	0.67167	0.61833
	WEIGHTED-MF	0.54686	0.49423	0.46123
	BPR	0.69951	0.65608	0.62494
AMAZON	SQL-RANK	0.04255	0.02978	0.02158
	WEIGHTED-MF	0.03647	0.02492	0.01914
	BPR	0.04863	0.01762	0.01306
YAHOO MUSIC	SQL-RANK	0.45512	0.36137	0.30689
	WEIGHTED-MF	0.39075	0.31024	0.27008
	BPR	0.37624	0.32184	0.28105
FOURSQUARE	SQL-RANK	0.05825	0.01941	0.01699
	WEIGHTED-MF	0.02184	0.01553	0.01407
	BPR	0.03398	0.01796	0.01359

- Amazon: the Amazon purchase rating data for musical instruments ³ with 339,232 users and 83,047 items.
- Yahoo-music: the Yahoo music rating data set ⁴ which contains 15,400 users and 1,000 items.
- Foursquare: a location check-in data⁵. The data set contains 3,112 users and 3,298 venues with 27,149 check-ins. The data set is already in the form of “0/1” so we do not need to do any transformation.

The experimental results are shown in Table 4.1. We find that SQL-Rank outperforms both Weighted-MF and BPR in most cases.

³<http://jmcauley.ucsd.edu/data/amazon/>

⁴<https://webscope.sandbox.yahoo.com/catalog.php?datatype=r&did=3>

⁵<https://sites.google.com/site/yangdingqi/home/foursquare-dataset>

4.3.2 Explicit Feedback

Next we compare the following methods in the explicit feedback setting:

- SQL-Rank: our proposed algorithm implemented in Julia. Note that in the explicit feedback setting our algorithm only considers pairs with explicit ratings.
- List-MF: the listwise algorithm using the cross entropy loss between observed rating and top 1 probability [100]. We use the C++ implementation on github⁶.
- MF: the classical matrix factorization algorithm in [59] utilizing a pointwise loss solved by SGD. We implemented SGD in Julia.
- Primal-CR++: the recently proposed pairwise algorithm in [115]. We use the Julia implementation released by the authors⁷.

Experiments are conducted on Movielens1m and Yahoo-music datasets. We perform the same procedure as in implicit feedback setting except that we do not need to mask the ratings into “0/1”.

We measure the performance in the following two ways:

- NDCG@ k : defined as:

$$\text{NDCG}@k = \frac{1}{n} \sum_{i=1}^n \frac{\text{DCG}@k(i, \Pi_i)}{\text{DCG}@k(i, \Pi_i^*)},$$

where i represents i -th user and

$$\text{DCG}@k(i, \Pi_i) = \sum_{l=1}^k \frac{2^{R_{i\Pi_{il}}} - 1}{\log_2(l + 1)}.$$

In the DCG definition, Π_{il} represents the index of the l -th ranked item for user i in test data based on the learned score matrix X . R is the rating matrix and R_{ij} is the rating given to item j by user i . Π_i^* is the ordering provided by the ground truth rating.

- Precision@ k : defined as a fraction of relevant items among the top k recommended items:

$$\text{precision}@k = \frac{\sum_{i=1}^n |\{1 \leq l \leq k : 4 \leq R_{i\Pi_{il}} \leq 5\}|}{n \cdot k},$$

⁶<https://github.com/gpoesia/listrankmf>

⁷<https://github.com/wuliwei9278/ml-1m>

Table 4.2. Comparing explicit feedback methods on various datasets.

DATASET	METHOD	NDCG@10	P@1	P@5	P@10
MOVIELENS1M	SQL-RANK	0.75076	0.50736	0.43692	0.40248
	LIST-MF	0.73307	0.45226	0.40482	0.38958
	PRIMAL-CR++	0.76826	0.49365	0.43098	0.39779
	MF	0.74661	0.00050	0.00096	0.00134
YAHOO MUSIC	SQL-RANK	0.66150	0.14983	0.12144	0.10192
	LIST-MF	0.67490	0.12646	0.11301	0.09865
	PRIMAL-CR++	0.66420	0.14291	0.10787	0.09104
	MF	0.69916	0.04944	0.03105	0.04787

here we consider items with ratings assigned as 4 or 5 as relevant. R_{ij} follows the same definitions above but unlike before Π_{il} gives the index of the l -th ranked item for user i among all the items that are not rated by user i in the training set (including both rated test items and unobserved items).

As shown in Table 4.2, our proposed listwise algorithm SQL-Rank outperforms previous listwise method List-MF in both NDCG@10 and precision@1, 5, 10. It verifies the claim that log-likelihood loss outperforms the cross entropy loss if we use it correctly. When listwise algorithm SQL-Rank is compared with pairwise algorithm Primal-CR++, the performances between SQL-Rank and Primal-CR++ are quite similar, slightly lower for NDCG@10 but higher for precision@1, 5, 10. Pointwise method MF is doing okay in NDCG but really bad in terms of precision. Despite having comparable NDCG, the predicted top k items given by MF are quite different from those given by other algorithms utilizing a ranking loss. The ordered lists based on SQL-Rank, Primal-CR++ and List-MF, on the other hand, share a lot of similarity and only have minor difference in ranking of some items. It is an interesting phenomenon that we think is worth exploring further in the future.

4.3.3 Training speed

To illustrate the training speed of our algorithm, we plot precision@1 versus training time for the Movielens1m dataset and the Foursquare dataset. Figure 4.2 and Figure 8.2 (in the appendix) show that our algorithm SQL-Rank is faster than BPR and Weighted-MF. Note that our algorithm is implemented in Julia while BPR and Weighted-MF are highly-optimized C++ codes (usually at least 2 times faster than Julia) released by Quora. This speed difference makes sense as our algorithm takes $O(n\bar{m}r)$ time, which is linearly to the observed ratings. In comparison, pair-wise model such as BPR has $O(n\bar{m}^2)$ pairs, so will take $O(n\bar{m}^2r)$ time for each epoch.

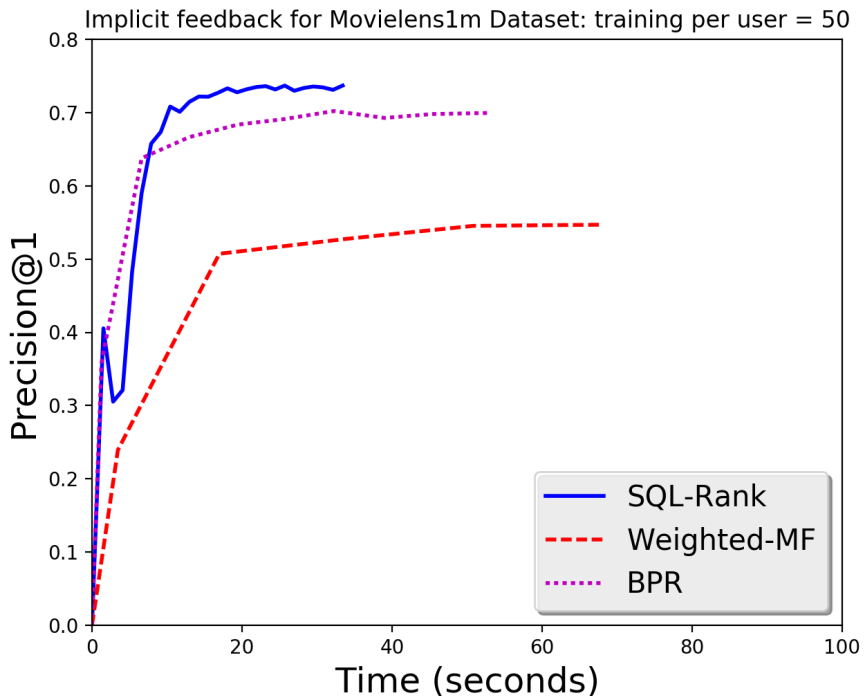


Figure 4.2. Training time of implicit feedback methods.

4.3.4 Effectiveness of Stochastic Queuing (SQ)

One important innovation in our SQL-Rank algorithm is the Stochastic Queuing (SQ) Process for handling ties. To illustrate the effectiveness of the SQ process, we compare

Table 4.3. Effectiveness of Stochastic Queuing Process.

METHOD	P@1	P@5	P@10
WITH SQ	0.73685	0.67167	0.61833
WITHOUT SQ	0.62763	0.58420	0.55036

our algorithm with and without SQ. Recall that without SQ means we fix a certain permutation matrix Π and optimize with respect to it throughout all iterations without generating new Π , while SQ allows us to update using a new permutation at each time. As shown Table 4.3 and Figure 8.3 (in the appendix), the performance gain from SQ in terms of precision is substantial (more than 10%) on Movielens1m dataset. It verifies the claim that our way of handling ties and missing data is very effective and improves the ranking results by a lot.

Table 4.4. Comparing different k on Movielens1m data set using 50 training data per user.

k	NDCG@10	P@1	P@5	P@10
5	0.64807	0.39156	0.33591	0.29855
10	0.67746	0.43118	0.34220	0.33339
25	0.74589	0.47003	0.42874	0.39796
50 (FULL LIST)	0.75076	0.50736	0.43692	0.40248

4.3.5 Effectiveness of using the Full List

Another benefit of our algorithm is that we are able to minimize top k probability with much larger k and without much overhead. Previous approaches [50] already pointed out increasing k leads to better ranking results, but their complexity is exponential to k so they were not able to have $k > 1$. To show the effectiveness of using permutation probability for full lists rather than using the top k probability for top k partial lists in the likelihood loss, we fix everything else to be the same and only vary k in Equation (4.5).

We obtain the results in Table 4.4 and Figure 8.4 (in the appendix). It shows that the larger k we use, the better the results we can get. Therefore, in the final model, we set k to be the maximum number (length of the observed list.)

4.4 Conclusions

In this chapter, we propose a listwise approach for collaborative ranking and provide an efficient algorithm to solve it. Our methodology can incorporate both implicit and explicit feedback, and can gracefully handle ties and missing data. In experiments, we demonstrate our algorithm outperforms existing state-of-the-art methods in terms of top k recommendation precision. We also provide a theoretical framework for analyzing listwise methods highlighting the dependence on the number of users and items.

Chapter 5

Stochastic Shared Embeddings: Data-driven Regularization of Embedding Layers

5.1 Introduction

Recently, embedding representations have been widely used in almost all AI-related fields, from feature maps [63] in computer vision, to word embeddings [71, 83] in natural language processing, to user/item embeddings [49, 73] in recommender systems. Usually, the embeddings are high-dimensional vectors. Take language models for example, in GPT [87] and Bert-Base model [25], 768-dimensional vectors are used to represent words. Bert-Large model utilizes 1024-dimensional vectors and GPT-2 [88] may have used even higher dimensions in their unreleased large models. In recommender systems, things are slightly different: the dimension of user/item embeddings are usually set to be reasonably small, 50 or 100. But the number of users and items is on a much bigger scale. Contrast this with the fact that the size of word vocabulary that normally ranges from 50,000 to 150,000, the number of users and items can be millions or even billions in large-scale real-world commercial recommender systems [6].

Given the massive number of parameters in modern neural networks with embedding layers, mitigating over-parameterization can play a big role in preventing over-fitting in deep learning. We propose a regularization method, Stochastic Shared Embeddings (SSE),

that uses prior information about similarities between embeddings, such as semantically and grammatically related words in natural languages or real-world users who share social relationships. Critically, SSE progresses by stochastically transitioning between embeddings as opposed to a more brute-force regularization such as graph-based Laplacian regularization and ridge regularization. Thus, SSE integrates seamlessly with existing stochastic optimization methods and the resulting regularization is data-driven.

We will begin the paper with the mathematical formulation of the problem, propose SSE, and provide the motivations behind SSE. We provide a theoretical analysis of SSE that can be compared with excess risk bounds based on empirical Rademacher complexity. We then conducted experiments for a total of 6 tasks from simple neural networks with one hidden layer in recommender systems, to the transformer and BERT in natural languages and find that when used along with widely-used regularization methods such as weight decay and dropout, our proposed methods can further reduce over-fitting, which often leads to more favorable generalization results.

5.2 Related Work

Regularization techniques are used to control model complexity and avoid over-fitting. ℓ_2 regularization [46] is the most widely used approach and has been used in many matrix factorization models in recommender systems; ℓ_1 regularization [109] is used when a sparse model is preferred. For deep neural networks, it has been shown that ℓ_p regularizations are often too weak, while dropout [45, 106] is more effective in practice. There are many other regularization techniques, including parameter sharing [30], max-norm regularization [104], gradient clipping [82], etc.

Our proposed SSE-graph is very different from graph Laplacian regularization [13], in which the distances of any two embeddings connected over the graph are directly penalized. Hard parameter sharing uses one embedding to replace all distinct embeddings in the same group, which inevitably introduces a significant bias. Soft parameter sharing [75] is similar to the graph Laplacian, penalizing the ℓ_2 distances between any two embeddings. These methods have no dependence on the loss, while the proposed SSE-graph method is

Algorithm 9. SSE-Graph for Neural Networks with Embeddings

- 1: **Input:** input x_i , label y_i , backpropagate T steps, mini-batch size m , knowledge graphs on embeddings $\{E_1, \dots, E_M\}$
 - 2: Define $p_l(\cdot, \cdot | \Phi)$ based on knowledge graphs on embeddings, $l = 1, \dots, M$
 - 3: **for** $t = 1$ **to** T **do**
 - 4: Sample one mini-batch $\{x_1, \dots, x_m\}$
 - 5: **for** $i = 1$ **to** m **do**
 - 6: Identify the set of embeddings $\mathcal{S}_i = \{E_1[j_1^i], \dots, E_M[j_M^i]\}$ for input x_i and label y_i
 - 7: **for** each embedding $E_l[j_l^i] \in \mathcal{S}_i$ **do**
 - 8: Replace $E_l[j_l^i]$ with $E_l[k_l]$, where $k_l \sim p_l(j_l^i, \cdot | \Phi)$
 - 9: Forward and backward pass with the new embeddings
 - 10: Return embeddings $\{E_1, \dots, E_M\}$, and neural network parameters Θ
-

data-driven in that the loss influences the effect of regularization. Unlike graph Laplacian regularization, hard and soft parameter sharing, our method is stochastic by nature. This allows our model to enjoy similar advantages as dropout [106].

Interestingly, in the original BERT model’s pre-training stage [25], a variant of SSE-SE is already implicitly used for token embeddings but for a different reason. In [25], the authors masked 15% of words and 10% of the time replaced the [mask] token with a random token. In the next section, we discuss how SSE-SE differs from this heuristic. Another closely related technique to ours is the label smoothing [107], which is widely used in the computer vision community. We find that in the classification setting if we apply SSE-SE to one-hot encodings associated with output y_i only, our SSE-SE is closely related to the label smoothing, which can be treated as a special case of our proposed method.

5.3 Stochastic Shared Embeddings

Throughout this chapter, the network input x_i and label y_i will be encoded into indices j_1^i, \dots, j_M^i which are elements of $\mathcal{I}_1 \times \dots \times \mathcal{I}_M$, the index sets of embedding tables. A typical

choice is that the indices are the encoding of a dictionary for words in natural language applications, or user and item tables in recommendation systems. Each index, j_l , within the l th table, is associated with an embedding $E_l[j_l]$ which is a trainable vector in \mathbb{R}^{d_l} . The embeddings associated with label y_i are usually non-trainable one-hot vectors corresponding to label look-up tables while embeddings associated with input x_i are trainable embedding vectors for embedding look-up tables. In natural language applications, we appropriately modify this framework to accommodate sequences such as sentences.

The loss function can be written as the functions of embeddings:

$$R_n(\Theta) = \sum_i \ell(x_i, y_i | \Theta) = \sum_i \ell(E_1[j_1^i], \dots, E_M[j_M^i] | \Theta), \quad (5.1)$$

where y_i is the label and Θ encompasses all trainable parameters including the embeddings, $\{E_l[j_l] : j_l \in \mathcal{I}_l\}$. The loss function ℓ is a mapping from embedding spaces to the reals. For text input, each $E_l[j_l^i]$ is a word embedding vector in the input sentence or document. For recommender systems, usually there are two embedding look-up tables: one for users and one for items [41]. So the objective function, such as mean squared loss or some ranking losses, will comprise both user and item embeddings for each input. We can more succinctly write the matrix of all embeddings for the i th sample as $\mathbf{E}[\mathbf{j}^i] = (E_1[j_1^i], \dots, E_M[j_M^i])$ where $\mathbf{j}^i = (j_1^i, \dots, j_M^i) \in \mathcal{I}$. By an abuse of notation we write the loss as a function of the embedding matrix, $\ell(\mathbf{E}[\mathbf{j}^i] | \Theta)$.

Suppose that we have access to knowledge graphs [66, 72] over embeddings, and we have a prior belief that two embeddings will share information and replacing one with the other should not incur a significant change in the loss distribution. For example, if two movies are both comedies and they are starred by the same actors, it is very likely that for the same user, replacing one comedy movie with the other comedy movie will result in little change in the loss distribution. In stochastic optimization, we can replace the loss gradient for one movie's embedding with the other similar movie's embedding, and this will not significantly bias the gradient if the prior belief is accurate. On the other hand, if this exchange is stochastic, then it will act to smooth the gradient steps in the long run, thus regularizing the gradient updates.

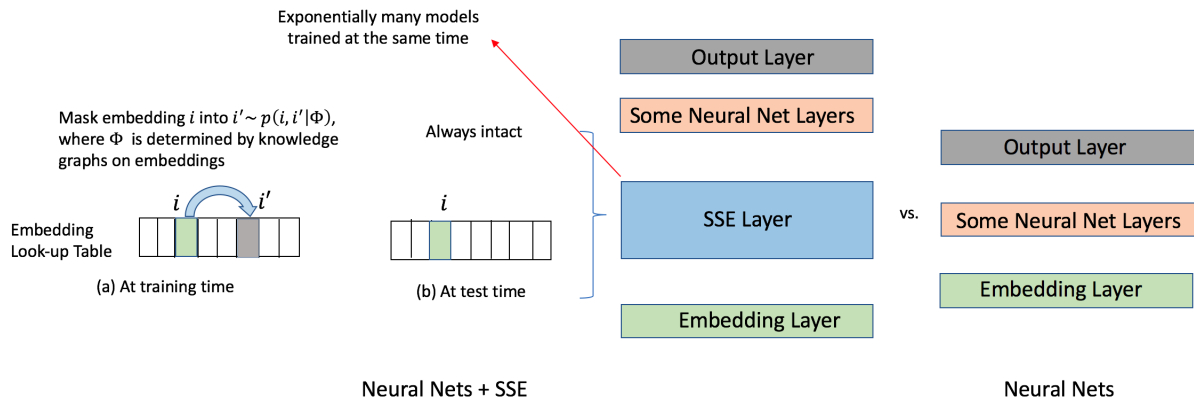


Figure 5.1. SSE-Graph described in Algorithm 9 and Figure 5.2 can be viewed as adding exponentially many distinct reordering layers above the embedding layer. A modified backpropagation procedure in Algorithm 9 is used to train exponentially many such neural networks at the same time.

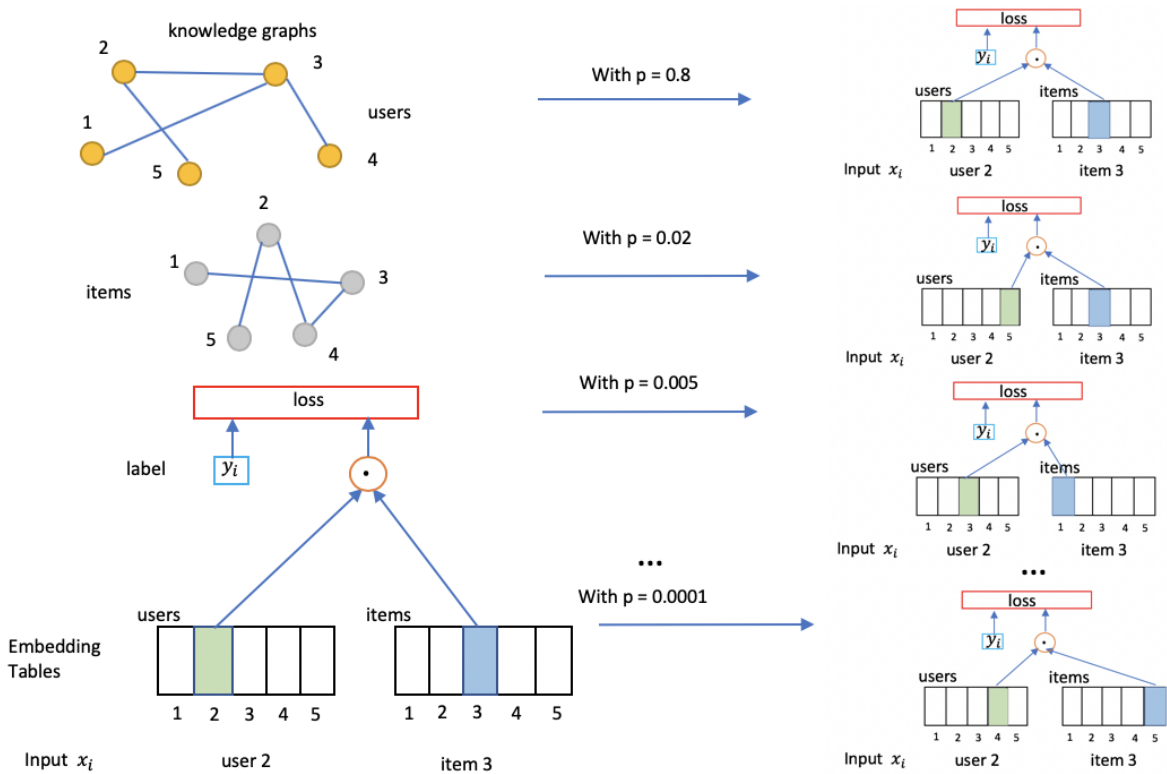


Figure 5.2. Illustration of how SSE-Graph algorithm in Figure 5.1 works for a simple neural network.

5.3.1 General SSE with Knowledge Graphs: SSE-Graph

Instead of optimizing objective function $R_n(\Theta)$ in (5.1), SSE-Graph described in Algorithm 9, Figure 5.1, and Figure 5.2 is approximately optimizing the objective function below:

$$S_n(\Theta) = \sum_i \sum_{\mathbf{k} \in \mathcal{I}} p(\mathbf{j}^i, \mathbf{k} | \Phi) \ell(\mathbf{E}[\mathbf{k}] | \Theta), \quad (5.2)$$

where $p(\mathbf{j}, \mathbf{k} | \Phi)$ is the transition probability (with parameters Φ) of exchanging the encoding vector $\mathbf{j} \in \mathcal{I}$ with a new encoding vector $\mathbf{k} \in \mathcal{I}$ in the Cartesian product index set of all embedding tables. When there is a single embedding table ($M = 1$) then there are no hard restrictions on the transition probabilities, $p(\cdot, \cdot)$, but when there are multiple tables ($M > 1$) then we will enforce that $p(\cdot, \cdot)$ takes a tensor product form (see (5.4)). When we are assuming that there is only a single embedding table ($M = 1$) we will not bold $j, E[j]$ and suppress their indices.

In the single embedding table case, $M = 1$, there are many ways to define transition probability from j to k . One simple and effective way is to use a random walk (with random restart and self-loops) on a knowledge graph \mathcal{G} , i.e. when embedding j is connected with k but not with l , we can set the ratio of $p(j, k | \Phi)$ and $p(j, l | \Phi)$ to be a constant greater than 1. In more formal notation, we have

$$j \sim k, j \not\sim l \longrightarrow p(j, k | \Phi) / p(j, l | \Phi) = \rho, \quad (5.3)$$

where $\rho > 1$ and is a tuning parameter. It is motivated by the fact that embeddings connected with each other in knowledge graphs should bear more resemblance and thus be more likely replaced by each other. Also, we let $p(j, j | \Phi) = 1 - p_0$, where p_0 is called the *SSE probability* and embedding retainment probability is $1 - p_0$. We treat both p_0 and ρ as tuning hyper-parameters in experiments. With (5.3) and $\sum_k p(j, k | \Phi) = 1$, we can derive transition probabilities between any two embeddings to fill out the transition probability table.

When there are multiple embedding tables, $M > 1$, then we will force that the transition from \mathbf{j} to \mathbf{k} can be thought of as independent transitions from j_l to k_l within embedding table l (and index set \mathcal{I}_l). Each table may have its own knowledge graph, resulting in its

own transition probabilities $p_l(\cdot, \cdot)$. The more general form of the SSE-graph objective is given below:

$$S_n(\Theta) = \sum_i \sum_{k_1, \dots, k_M} p_1(j_1^i, k_1 | \Phi) \cdots p_M(j_M^i, k_M | \Phi) \ell(E_1[k_1], \dots, E_M[k_M] | \Theta), \quad (5.4)$$

Intuitively, this SSE objective could reduce the variance of the estimator.

Optimizing (5.4) with SGD or its variants (Adagrad [26], Adam [57]) is simple. We just need to randomly switch each original embedding tensor $\mathbf{E}[\mathbf{j}^i]$ with another embedding tensor $\mathbf{E}[\mathbf{k}]$ randomly sampled according to the transition probability (see Algorithm 9). This is equivalent to have a randomized embedding look-up layer as shown in Figure 5.1.

We can also accommodate sequences of embeddings, which commonly occur in natural language application, by considering $(j_{l,1}^i, k_{l,1}), \dots, (j_{l,n_l}^i, k_{l,n_l}^i)$ instead of (j_l^i, k_l) for l -th embedding table in (5.4), where $1 \leq l \leq M$ and n_l^i is the number of embeddings in table l that are associated with (x_i, y_i) . When there is more than one embedding look-up table, we sometimes prefer to use different p_0 and ρ for different look-up tables in (5.3) and the SSE probability constraint. For example, in recommender systems, we would use p_u, ρ_u for user embedding table and p_i, ρ_i for item embedding table.

We find that SSE with knowledge graphs, i.e., SSE-Graph, can force similar embeddings to cluster when compared to the original neural network without SSE-Graph. In Figure 5.3, one can easily see that more embeddings tend to cluster into 2 black holes after applying SSE-Graph when embeddings are projected into 3D spaces using PCA. Interestingly, a similar phenomenon occurs when assuming the knowledge graph is a complete graph, which we would introduce as SSE-SE below.

5.3.2 Simplified SSE with Complete Graph: SSE-SE

One clear limitation of applying the SSE-Graph is that not every dataset comes with good-quality knowledge graphs on embeddings. For those cases, we could assume there is a complete graph over all embeddings so there is a small transition probability between every pair of different embeddings:

$$p(j, k | \Phi) = \frac{p_0}{N-1}, \quad \forall 1 \leq k \neq j \leq N, \quad (5.5)$$

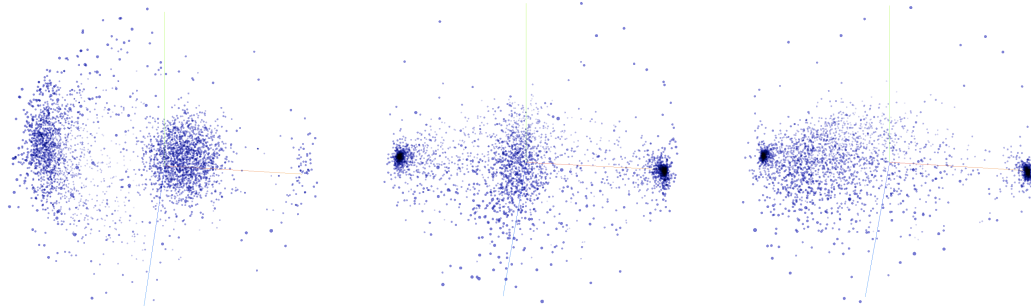


Figure 5.3. Projecting 50-dimensional embeddings obtained by training a simple neural network without SSE (Left), and with SSE-Graph (Center) , SSE-SE (Right) into 3D space using PCA.

where N is the size of the embedding table. The SGD procedure in Algorithm 9 can still be applied and we call this algorithm SSE-SE (Stochastic Shared Embeddings - Simple and Easy). It is worth noting that SSE-Graph and SSE-SE are applied to embeddings associated with not only input x_i but also those with output y_i . Unless there are considerably many more embeddings than data points and model is significantly overfitting, normally $p_0 = 0.01$ gives reasonably good results.

Interestingly, we found that the SSE-SE framework is related to several techniques used in practice. For example, BERT pre-training unintentionally applied a method similar to SSE-SE to input x_i by replacing the masked word with a random word. This would implicitly introduce an SSE layer for input x_i in Figure 5.1, because now embeddings associated with input x_i be stochastically mapped according to (5.5). The main difference between this and SSE-SE is that it merely augments the input once, while SSE introduces randomization at every iteration, and we can also accommodate label embeddings. In experimental Section 5.4.4, we will show that SSE-SE would improve original BERT pre-training procedure as well as fine-tuning procedure.

5.3.3 Theoretical Guarantees

We explain why SSE can reduce the variance of estimators and thus leads to better generalization performance. For simplicity, we consider the SSE-graph objective (5.2) where there is no transition associated with the label y_i , and only the embeddings associated with the input x_i undergo a transition. When this is the case, we can think of the loss as a

function of the x_i embedding and the label, $\ell(\mathbf{E}[\mathbf{j}^i], y_i; \Theta)$. We take this approach because it is more straightforward to compare our resulting theory to existing excess risk bounds.

The SSE objective in the case of only input transitions can be written as,

$$S_n(\Theta) = \sum_i \sum_{\mathbf{k}} p(\mathbf{j}^i, \mathbf{k}) \cdot \ell(\mathbf{E}[\mathbf{k}], y_i | \Theta), \quad (5.6)$$

and there may be some constraint on Θ . Let $\hat{\Theta}$ denote the minimizer of S_n subject to this constraint. We will show in the subsequent theory that minimizing S_n will get us close to a minimizer of $S(\Theta) = \mathbb{E}S_n(\Theta)$, and that under some conditions this will get us close to the Bayes risk. We will use the standard definitions of empirical and true risk, $R_n(\Theta) = \sum_i \ell(x_i, y_i | \Theta)$ and $R(\Theta) = \mathbb{E}R_n(\Theta)$.

Our results depend on the following decomposition of the risk. By optimality of $\hat{\Theta}$,

$$R(\hat{\Theta}) = S_n(\hat{\Theta}) + [R(\hat{\Theta}) - S(\hat{\Theta})] + [S(\hat{\Theta}) - S_n(\hat{\Theta})] \leq S_n(\Theta^*) + B(\hat{\Theta}) + \mathcal{E}(\hat{\Theta}) \quad (5.7)$$

where $B(\Theta) = |R(\Theta) - S(\Theta)|$, and $E(\Theta) = |S(\Theta) - S_n(\Theta)|$. We can think of $B(\Theta)$ as representing the bias due to SSE, and $E(\Theta)$ as an SSE form of excess risk. Then by another application of similar bounds,

$$R(\hat{\Theta}) \leq R(\Theta^*) + B(\hat{\Theta}) + B(\Theta^*) + E(\hat{\Theta}) + E(\Theta^*). \quad (5.8)$$

The high level idea behind the following results is that when the SSE protocol reflects the underlying distribution of the data, then the bias term $B(\Theta)$ is small, and if the SSE transitions are well mixing then the SSE excess risk $E(\Theta)$ will be of smaller order than the standard Rademacher complexity. This will result in a small excess risk.

Theorem 1. *Consider SSE-graph with only input transitions.*

Let $L(\mathbf{E}[\mathbf{j}^i]) = \mathbb{E}_{Y|X=x^i} \ell(\mathbf{E}[\mathbf{j}^i], Y | \Theta)$ be the expected loss conditional on input x^i and $e(\mathbf{E}[\mathbf{j}^i], y; \Theta) = \ell(\mathbf{E}[\mathbf{j}^i], y | \Theta) - L(\mathbf{E}[\mathbf{j}^i] | \Theta)$ be the residual loss. Define the conditional and residual SSE empirical Rademacher complexities to be

$$\rho_{L,n} = \mathbb{E}_\sigma \sup_{\Theta} \left| \sum_i \sigma_i \sum_{\mathbf{k}} p(\mathbf{j}^i, \mathbf{k}) \cdot L(\mathbf{E}[\mathbf{k}] | \Theta) \right|, \quad (5.9)$$

$$\rho_{e,n} = \mathbb{E}_\sigma \sup_{\Theta} \left| \sum_i \sigma_i \sum_{\mathbf{k}} p(\mathbf{j}^i, \mathbf{k}) \cdot e(\mathbf{E}[\mathbf{k}], y_i; \Theta) \right|, \quad (5.10)$$

respectively where σ is a Rademacher ± 1 random vectors in \mathbb{R}^n . Then we can decompose the SSE empirical risk into

$$\mathbb{E} \sup_{\Theta} |S_n(\Theta) - S(\Theta)| \leq 2\mathbb{E}[\rho_{L,n} + \rho_{e,n}]. \quad (5.11)$$

Remark 1. The transition probabilities in (5.9), (5.10) act to smooth the empirical Rademacher complexity. To see this, notice that we can write the inner term of (5.9) as $(P\sigma)^\top L$, where we have vectorized $\sigma_i, L(x_i; \Theta)$ and formed the transition matrix P . Transition matrices are contractive and will induce dependencies between the Rademacher random variables, thereby stochastically reducing the supremum. In the case of no label noise, namely that $Y|X$ is a point mass, $e(x, y; \Theta) = 0$, and $\rho_{e,n} = 0$. The use of L as opposed to the losses, ℓ , will also make $\rho_{L,n}$ of smaller order than the standard empirical Rademacher complexity. We demonstrate this with a partial simulation of $\rho_{L,n}$ on the *Movielens1m* dataset in Figure 8.5 of the Appendix.

Theorem 2. Let the SSE-bias be defined as

$$\mathcal{B} = \sup_{\Theta} \left| \mathbb{E} \left[\sum_i \sum_{\mathbf{k}} p(\mathbf{j}^i, \mathbf{k}) \cdot (\ell(\mathbf{E}[\mathbf{k}], y_i | \Theta) - \ell(\mathbf{E}[\mathbf{j}^i], y_i | \Theta)) \right] \right|.$$

Suppose that $0 \leq \ell(\cdot, \cdot; \Theta) \leq b$ for some $b > 0$, then

$$\mathbb{P} \left\{ R(\hat{\Theta}) > R(\Theta^*) + 2\mathcal{B} + 4\mathbb{E}[\rho_{L,n} + \rho_{e,n}] + \sqrt{nu} \right\} \leq e^{-\frac{u^2}{2b^2}}.$$

Remark 2. The price for ‘smoothing’ the Rademacher complexity in Theorem 1 is that SSE may introduce a bias. This will be particularly prominent when the SSE transitions have little to do with the underlying distribution of Y, X . On the other extreme, suppose that $p(\mathbf{j}, \mathbf{k})$ is non-zero over a neighborhood \mathcal{N}_j of \mathbf{j} , and that for data x', y' with encoding $\mathbf{k} \in \mathcal{N}_j$, x', y' is identically distributed with x_i, y_i , then $\mathcal{B} = 0$. In all likelihood, the SSE transition probabilities will not be supported over neighborhoods of iid random pairs, but with a well chosen SSE protocol the neighborhoods contain approximately iid pairs and \mathcal{B} is small.

Table 5.1. Compare SSE-Graph and SSE-SE against ALS-MF with Graph Laplacian Regularization. The p_u and p_i are the SSE probabilities for user and item embedding tables respectively, as in (5.5). Definitions of ρ_u and ρ_i can be found in (5.3). Movielens10m does not have user graphs.

Model	Movielens1m					Movielens10m				
	RMSE	ρ_u	ρ_i	p_u	p_i	RMSE	ρ_u	ρ_i	p_u	p_i
SGD-MF	1.0984	-	-	-	-	1.9490	-	-	-	-
Graph Laplacian + ALS-MF	1.0464	-	-	-	-	1.9755	-	-	-	-
SSE-Graph + SGD-MF	1.0145	500	200	0.005	0.005	1.9019	1	500	0.01	0.01
SSE-SE + SGD-MF	1.0150	1	1	0.005	0.005	1.9085	1	1	0.01	0.01

Table 5.2. SSE-SE outperforms Dropout for Neural Networks with One Hidden Layer such as Matrix Factorization Algorithm regardless of dimensionality we use. p_s is the SSE probability for both user and item embedding tables and p_d is the dropout probability.

Model	Douban			Movielens10m			Netflix		
	RMSE	p_d	p_s	RMSE	p_d	p_s	RMSE	p_d	p_s
MF	0.7339	-	-	0.8851	-	-	0.8941	-	-
Dropout + MF	0.7296	0.1	-	0.8813	0.1	-	0.8897	0.1	-
SSE-SE + MF	0.7201	-	0.008	0.8715	-	0.008	0.8842	-	0.008
SSE-SE + Dropout + MF	0.7185	0.1	0.005	0.8678	0.1	0.005	0.8790	0.1	0.005

Table 5.3. SSE-SE outperforms dropout for Neural Networks with One Hidden Layer such as Bayesian Personalized Ranking Algorithm regardless of dimensionality we use. We report the metric precision for top k recommendations as $P@k$.

Model	Movielens1m			Yahoo Music			Foursquare		
	$P@1$	$P@5$	$P@10$	$P@1$	$P@5$	$P@10$	$P@1$	$P@5$	$P@10$
SQL-Rank (2018)	0.7369	0.6717	0.6183	0.4551	0.3614	0.3069	0.0583	0.0194	0.0170
BPR	0.6977	0.6568	0.6257	0.3971	0.3295	0.2806	0.0437	0.0189	0.0143
Dropout + BPR	0.7031	0.6548	0.6273	0.4080	0.3315	0.2847	0.0437	0.0184	0.0146
SSE-SE + BPR	0.7254	0.6813	0.6469	0.4297	0.3498	0.3005	0.0609	0.0262	0.0155

5.4 Experiments

We have conducted extensive experiments on 6 tasks, including 3 recommendation tasks (explicit feedback, implicit feedback and sequential recommendation) and 3 NLP tasks (neural machine translation, BERT pre-training, and BERT fine-tuning for sentiment classification) and found that our proposed SSE can effectively improve generalization performances on a wide variety of tasks. Note that the details about datasets and parameter settings can be found in the appendix.

5.4.1 Neural Networks with One Hidden Layer (Matrix Factorization and BPR)

Matrix Factorization Algorithm (MF) [73] and Bayesian Personalized Ranking Algorithm (BPR) [91] can be viewed as neural networks with one hidden layer (latent features) and are quite popular in recommendation tasks. MF uses the squared loss designed for explicit feedback data while BPR uses the pairwise ranking loss designed for implicit feedback data.

First, we conduct experiments on two explicit feedback datasets: Movielens1m and Movielens10m. For these datasets, we can construct graphs based on actors/actresses starring the movies. We compare SSE-graph and the popular Graph Laplacian Regularization (GLR) method [89] in Table 5.1. The results show that SSE-graph consistently outperforms GLR. This indicates that our SSE-Graph has greater potentials over graph Laplacian regularization as we do not explicitly penalize the distances across embeddings, but rather we implicitly penalize the effects of similar embeddings on the loss. Furthermore, we show that even without existing knowledge graphs of embeddings, our SSE-SE performs only slightly worse than SSE-Graph but still much better than GLR and MF.

In general, SSE-SE is a good alternative when graph information is not available. We then show that our proposed SSE-SE can be used together with standard regularization techniques such as dropout and weight decay to improve recommendation results regardless of the loss functions and dimensionality of embeddings. This is evident in Table 5.2 and Table 5.3. With the help of SSE-SE, BPR can perform better than the state-of-art listwise approach SQL-Rank [116] in most cases. We include the optimal SSE parameters in the

Table 5.4. SSE-SE has two tuning parameters: probability p_x to replace embeddings associated with input x_i and probability p_y to replace embeddings associated with output y_i . We use the dropout probability of 0.1, weight decay of $1e^{-5}$, and learning rate of $1e^{-3}$ for all experiments.

Model	Movielens1m		Dimension	# of Blocks	SSE-SE Parameters	
	NDCG@10	Hit Ratio@10	d	b	p_x	p_y
SASRec	0.5941	0.8182	100	2	-	-
SASRec	0.5996	0.8272	100	6	-	-
SSE-SE + SASRec	0.6092	0.8250	100	2	0.1	0
SSE-SE + SASRec	0.6085	0.8293	100	2	0	0.1
SSE-SE + SASRec	0.6200	0.8315	100	2	0.1	0.1
SSE-SE + SASRec	0.6265	0.8364	100	6	0.1	0.1

table for references and leave out other experiment details to the appendix. In the rest of the paper, we would mostly focus on SSE-SE as we do not have high-quality graphs of embeddings on most datasets.

5.4.2 Transformer Encoder Model for Sequential Recommendation

SASRec [56] is the state-of-the-arts algorithm for sequential recommendation task. It applies the transformer model [110], where a sequence of items purchased by a user can be viewed as a sentence in transformer, and next item prediction is equivalent to next word prediction in the language model. In Table 5.4, we perform SSE-SE on input embeddings ($p_x = 0.1, p_y = 0$), output embeddings ($p_x = 0.1, p_y = 0$) and both embeddings ($p_x = p_y = 0.1$), and observe that all of them significantly improve over state-of-the-art SASRec ($p_x = p_y = 0$). The regularization effects of SSE-SE is even more obvious when we increase the number of self-attention blocks from 2 to 6, as this will lead to a more sophisticated model with many more parameters. This leads to the model overfitting terribly even with dropout and weight decay. We can see in Table 5.4 that when both methods use dropout and weight decay, SSE-SE + SASRec is doing much better than SASRec without SSE-SE.

Table 5.5. Our proposed SSE-SE helps the Transformer achieve better BLEU scores on English-to-German in 10 out of 11 newstest data between 2008 and 2018.

Model	Test BLEU										
	2008	2009	2010	2011	2012	2013	2014	2015	2016	2017	2018
Transformer	21.0	20.7	22.7	20.6	20.6	25.3	26.2	28.4	32.1	27.2	38.8
SSE-SE + Transformer	21.4	21.1	23.0	21.0	20.8	25.2	27.2	29.2	33.1	27.9	39.9

5.4.3 Neural Machine Translation

We use the transformer model [110] as the backbone for our experiments. The baseline model is the standard 6-layer transformer architecture and we apply SSE-SE to both encoder, and decoder by replacing corresponding vocabularies’ embeddings in the source and target sentences. We trained on the standard WMT 2014 English to German dataset which consists of roughly 4.5 million parallel sentence pairs and tested on WMT 2008 to 2018 news-test sets. We use the OpenNMT implementation in our experiments. We use the same dropout rate of 0.1 and label smoothing value of 0.1 for the baseline model and our SSE-enhanced model. The only difference between the two models is whether or not we use our proposed SSE-SE with $p_0 = 0.01$ in (5.5) for both encoder and decoder embedding layers. We evaluate both models’ performances on the test datasets using BLEU scores [85].

We summarize our results in Table 5.5 and find that SSE-SE helps improving accuracy and BLEU scores on both dev and test sets in 10 out of 11 years from 2008 to 2018. In particular, on the last 5 years’ test sets from 2014 to 2018, the transformer model with SSE-SE improves BLEU scores by 0.92 on average when compared to the baseline model without SSE-SE.

5.4.4 BERT for Sentiment Classification

BERT’s model architecture [25] is a multi-layer bidirectional Transformer encoder based on the Transformer model in neural machine translation. Despite SSE-SE can be used for both pre-training and fine-tuning stages of BERT, we want to mainly focus on pre-training as fine-tuning bears more similarity to the previous section. We use SSE probability of 0.015 for embeddings (one-hot encodings) associated with labels and SSE probability of

Table 5.6. Our proposed SSE-SE applied in the pre-training stage on our crawled IMDB data improves the generalization ability of pre-trained IMDB model and helps the BERT-Base model outperform current SOTA results on the IMDB Sentiment Task after fine-tuning.

Model	IMDB Test Set		
	AUC	Accuracy	F1 Score
ULMFiT [47]	-	0.9540	-
Google Pre-trained Model + Fine-tuning	0.9415	0.9415	0.9419
Pre-training + Fine-tuning	0.9518	0.9518	0.9523
(SSE-SE + Pre-training) + Fine-tuning	0.9542	0.9542	0.9545

0.015 for embeddings (word-piece embeddings) associated with inputs. One thing worth noting is that even in the original BERT model’s pre-training stage, SSE-SE is already implicitly used for token embeddings. In original BERT model, the authors masked 15% of words for a maximum of 80 words in sequences of maximum length of 512 and 10% of the time replaced the [mask] token with a random token. That is roughly equivalent to SSE probability of 0.015 for replacing input word-piece embeddings.

We continue to pre-train Google pre-trained BERT model on our crawled IMDB movie reviews with and without SSE-SE and compare downstream tasks performances. In Table 5.6, we find that SSE-SE pre-trained BERT base model helps us achieve the state-of-the-art results for the IMDB sentiment classification task, which is better than the previous best in [47]. We report test set accuracy of 0.9542 after fine-tuning for one epoch only. For the similar SST-2 sentiment classification task in Table 5.7, we also find that SSE-SE can improve BERT pre-trains better. Our SSE-SE pre-trained model achieves 94.3% accuracy on SST-2 test set after 3 epochs of fine-tuning while the standard pre-trained BERT model only reports 93.8 after fine-tuning. Furthermore, we show that SSE-SE with SSE probability 0.01 can also improve dev and test accuracy in the fine-tuning stage. If we are using SSE-SE for both pre-training and fine-tuning stage of the BERT base model, we can achieve 94.5% accuracy on the SST-2 test set, approaching the 94.9% accuracy by the BERT large model. We are optimistic that our SSE-SE can be applied to BERT large model as well in the future.

Table 5.7. SSE-SE pre-trained BERT-Base models on IMDB datasets turn out working better on the new unseen SST-2 Task as well.

Model	SST-2 Dev Set			SST-2 Test Set
	AUC	Accuracy	F1 Score	Accuracy (%)
Google Pre-trained + Fine-tuning	0.9230	0.9232	0.9253	93.6
Pre-training + Fine-tuning	0.9265	0.9266	0.9281	93.8
(SSE-SE + Pre-training) + Fine-tuning	0.9276	0.9278	0.9295	94.3
(SSE-SE + Pre-training) + (SSE-SE + Fine-tuning)	0.9323	0.9323	0.9336	94.5

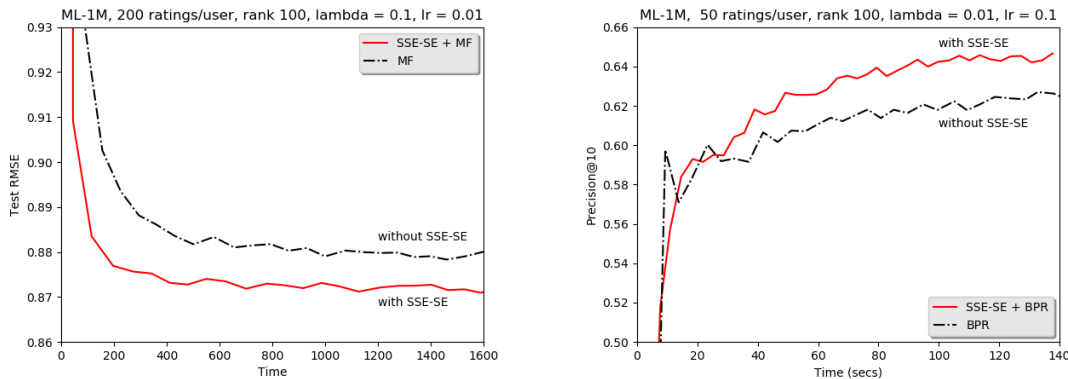


Figure 5.4. Compare Training Speed of Simple Neural Networks with One Hidden Layer, i.e. MF and BPR, with and without SSE-SE.

5.4.5 Speed and convergence comparisons.

In Figure 5.4, it is clear to see that our one-hidden-layer neural networks with SSE-SE are achieving much better generalization results than their respective standalone versions. One can also easily spot that SSE-version algorithms converge at much faster speeds with the same learning rate.

5.5 Conclusion

We have proposed Stochastic Shared Embeddings, which is a data-driven approach to regularization, that stands in contrast to brute force regularization such as Laplacian and ridge regularization. Our theory is a first step towards explaining the regularization effect of SSE, particularly, by ‘smoothing’ the Rademacher complexity. The extensive experimentation demonstrates that SSE can be fruitfully integrated into existing deep learning applications.

Chapter 6

SSE-PT: Sequential Recommendation Via Personalized Transformer

6.1 Introduction

The sequential recommendation problem has been an important open research question, yet using temporal information to improve recommendation performance has proven to be challenging. SASRec, proposed by [56] for sequential recommendation problems, has achieved state-of-the-art results and enjoyed more than 10x speed-up when compared to earlier CNN/RNN-based methods. However, the model used in SASRec is the standard Transformer which is inherently an un-personalized model. In practice, it is important to include a personalized Transformer in SASRec especially for recommender systems, but [56] found that adding additional personalized embeddings did not improve the performance of their Transformer model, and postulate that the failure of adding personalization is due to the fact that they already use the user history and the user embeddings only contribute to overfitting. In this work, we propose a novel method, Personalized Transformer (SSE-PT), that successfully introduces personalization into self-attentive neural network architectures.

Introducing user embeddings into the standard transformer model is intrinsically difficult with existing regularization techniques, as unavoidably a large number of user parameters are introduced, which is often at the same scale of the number of training

data. But we show that personalization can greatly improve ranking performance with a recent regularization technique called Stochastic Shared Embeddings (SSE) [117]. The personalized Transformer (SSE-PT) model with SSE regularization works well for all 5 real-world datasets we consider without overfitting, outperforming previous state-of-the-art algorithm SASRec by almost 5% in terms of NDCG@10. Furthermore, after examining some random users’ engagement history, we find our model is not only more interpretable but also able to focus on recent engagement patterns for each user. Moreover, our SSE-PT model with a slight modification, which we call SSE-PT++, can handle extremely long sequences and outperform SASRec in ranking results with comparable training speed, striking a balance between performance and speed requirements.

6.2 Related Work

6.2.1 Session-based and Sequential Recommendation

Both session-based and sequential (i.e., next-basket) recommendation algorithms take advantage of additional temporal information to make better personalized recommendations. The main difference between session-based recommendations [43] and sequential recommendations [56] is that the former assumes that the user ids are not recorded and therefore the length of engagement sequences are relatively short. Therefore, session-based recommendations normally do not consider user factors. On the other hand, sequential recommendation treats each sequence as a user’s engagement history [56]. Both settings, do not explicitly require time-stamps: only the relative temporal orderings are assumed known (in contrast to, for example, timeSVD++ [60] using time-stamps). Initially, sequence data in temporal order are usually modelled with Markov models, in which a future observation is conditioned on the last few observed items [92]. In [92], a personalized Markov model with user latent factors is proposed for more personalized results.

In recent years, deep learning techniques, borrowed from natural language processing (NLP) literature, are getting widely used in tackling sequential data. Like word sentences in NLP, item sequences in recommendations can be similarly modelled by recurrent neural networks (RNN) [42, 43] and convolutional neural network (CNN) [108] models.

Recently, attention models are increasingly used in both NLP [25, 110] and recommender systems [56, 68]. SASRec [56] is a recent method with state-of-the-art performance among the many deep learning models. Motivated by the Transformer model in neural machine translation [110], SASRec utilizes a similar architecture to the encoder part of the Transformer model. Our proposed model, SSE-PT, is a personalized extension of the transformer model.

6.2.2 Regularization Techniques

In deep learning, models with many more parameters than data points can easily overfit to the training data. This may prevent us from adding user embeddings as additional parameters into complicated models like the Transformer model [56], which can easily have 20 layers with millions of parameters for a medium-sized dataset like Movielens10M [38]. ℓ_2 regularization [46] is the most widely used approach and has been used in many matrix factorization models in recommender systems; ℓ_1 regularization [109] is used when a sparse model is preferred. For deep neural networks, it has been shown that ℓ_p regularizations are often too weak, while dropout [45, 106] is more effective in practice. There are many other regularization techniques, including parameter sharing [30], max-norm regularization [104], gradient clipping [82], etc. Very recently, a new regularization technique called Stochastic Shared Embeddings (SSE) [117] is proposed as a new means of regularizing embedding layers. We find that the base version SSE-SE is essential to the success of our Personalized Transformer (SSE-PT) model.

6.3 Methodology

6.3.1 Sequential Recommendation

Given n users and each user engaging with a subset of m items in a temporal order, the goal of sequential recommendation is to learn a good personalized ranking of top K items out of total m items for any given user at any given time point. We assume data in the format of n item sequences:

$$s_i = (j_{i1}, j_{i2}, \dots, j_{iT}) \text{ for } 1 \leq i \leq n. \quad (6.1)$$

Sequences s_i of length T contain indices of the last T items that user i has interacted with in the temporal order (from old to new). For different users, the sequence lengths can vary, but we can pad the shorter sequences so all of them have length T . We cannot simply randomly split data points into train/validation/test sets because they come in temporal orders. Instead, we need to make sure our training data is before validation data which is before test data temporally. We use last items in sequences as test sets, second-to-last items as validation sets and the rest as training sets. We use ranking metrics such as NDCG@ K and Recall@ K for evaluations, which are defined in the Appendix.

6.3.2 Personalized Transformer Architecture

Our model, which we call SSE-PT, is motivated by the Transformer model in [110] and [56]. It also utilizes a new regularization technique called stochastic shared embeddings [117]. In the following sections, we are going to examine each important component of our Personalized Transformer (SSE-PT) model, especially the embedding layer, and the novel application of stochastic shared embeddings (SSE) regularization technique.

Embedding Layer We define a learnable user embedding look-up table $U \in R^{n \times d_u}$ and item embedding look-up table $V \in R^{m \times d_i}$, where d_u, d_i are the number of hidden units for user and item respectively. We also specify learnable positional encoding table $P \in R^{T \times d}$, where $d = d_u + d_i$. So each input sequence $s_i \in R^T$ will be represented by the following embedding:

$$E = \begin{bmatrix} [v_{j_{i1}}; u_i] + p_1 \\ [v_{j_{i2}}; u_i] + p_2 \\ \vdots \\ [v_{j_{iT}}; u_i] + p_T \end{bmatrix} \in R^{T \times d}, \quad (6.2)$$

where $[v_{j_{it}}; u_i]$ represents concatenating item embedding $v_{j_{it}} \in R^{d_i}$ and user embedding $u_i \in R^{d_u}$ into embedding $E_t \in R^d$ for time t . Note that the main difference between our model and [56] is that we introduce the user embeddings u_i , making our model personalized.

Transformer Encoder On top of the embedding layer, we have B blocks of self-attention layers and fully connected layers, where each layer extracts features for each time step based on the previous layer’s outputs. Since this part is identical to the Transformer

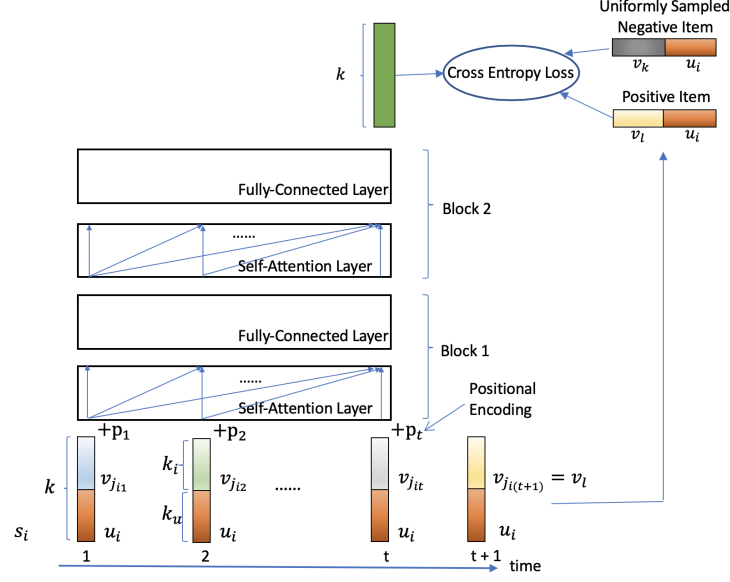


Figure 6.1. Illustration of our proposed SSE-PT model

encoder used in the original papers [56, 110], we will skip the details.

Prediction Layer At time t , the predicted probability of user i engaged item l is:

$$p_{itl} = \sigma(r_{itl}), \quad (6.3)$$

where σ is the sigmoid function and r_{itl} is the predicted score of item l by user i at time point t , defined as:

$$r_{itl} = F_{t-1}^B \cdot [v_l; u_i], \quad (6.4)$$

where F_{t-1}^B is the output hidden units associated with the transformer encoder at the last timestamp. Although we can use another set of user and item embedding look-up tables for the u_i and v_l , we find it better to use the same set of embedding look-up tables U, V as in the embedding layer. But regularization for those embeddings can be different. To distinguish the u_i and v_l in (6.4) from u_i, v_j in (6.2), we call embeddings in (6.4) output embeddings and those in (6.2) input embeddings.

The binary cross entropy loss between predicted probability for the positive item $l = j_{i(t+1)}$ and one uniformly sampled negative item $k \in \Omega$ is given as $-\log(p_{itl}) + \log(1 - p_{itk})$. Summing over s_i and t , we obtain the objective function that we want to minimize is:

$$\sum_i \sum_{t=1}^{T-1} \sum_{k \in \Omega} -[\log(p_{itl}) + \log(1 - p_{itk})]. \quad (6.5)$$

At the inference time, top- K recommendations for user i at time t can be made by sorting scores $r_{i\ell t}$ for all items ℓ and recommending the first K items in the sorted list.

Novel Application of Stochastic Shared Embeddings The most important regularization technique to SSE-PT model is the Stochastic Shared Embeddings (SSE) [117]. The main idea of SSE is to stochastically replace embeddings with another embedding with some pre-defined probability during SGD, which has the effect of regularizing the embedding layers. Without SSE, all the existing well-known regularization techniques like layer normalization, dropout and weight decay fail and cannot prevent the model from over-fitting badly after introducing user embeddings. [117] develops two versions of SSE, SSE-Graph and SSE-SE. In the simplest uniform case, SSE-SE replaces one embedding with another embedding uniformly with probability p , which is called SSE probability in [117]. Since we don't have knowledge graphs for user or items, we simply apply the SSE-SE to our SSE-PT model. We find SSE-SE makes possible training this personalized model with $O(nd_u)$ additional parameters.

There are 3 different places in our model that SSE-SE can be applied. We can apply SSE-SE to input/output user embeddings, input item embeddings, and output item embeddings with probabilities p_u , p_i and p_y respectively. Note that input user embedding and output user embedding are always replaced at the same time with SSE probability p_u . Empirically, we find that SSE-SE to user embeddings and output item embeddings always helps, but SSE-SE to input item embeddings is only useful when the average sequence length is large, e.g., more than 100 in Movielens1M and Movielens10M datasets.

Other Regularization Techniques Besides the *SSE* [117], we also utilized other widely used regularization techniques, including *layer normalization* [4], *batch normalization* [51], *residual connections* [39], *weight decay* [64], and *dropout* [106]. Since they are used in the same way in the previous paper [56], we omit the details to the Appendix.

6.3.3 Handling Long Sequences: SSE-PT++

To handle extremely long sequences, a slight modification can be made on the base SSE-PT model in terms of how input sequences s_i 's are fed into the SSE-PT neural network. We call the enhanced model SSE-PT++ to distinguish it from the previously discussed SSE-PT

model, which cannot handle sequences longer than T .

The motivation of SSE-PT++ over SSE-PT comes from: sometimes we want to make use of extremely long sequences, $s_i = (j_{i1}, j_{i2}, \dots, j_{it})$ for $1 \leq i \leq n$, where $t > T$, but our SSE-PT model can only handle sequences of maximum length of T . The simplest way is to sample starting index $1 \leq v \leq t$ uniformly and use $s_i = (j_{iv}, j_{i(v+1)}, \dots, j_{iz})$, where $z = \min(t, v + T - 1)$. Although sampling the starting index uniformly from $[1, t]$ can accommodate long sequences of length $t > T$, this does not work well in practice. Uniform sampling does not take into account the importance of recent items in a long sequence. To solve this dilemma, we introduce an additional hyper-parameter p_s which we call *sampling probability*. It implies that with probability p_s , we sample the starting index v uniformly from $[1, t - T]$ and use sequence $s_i = (j_{iv}, j_{i(v+1)}, \dots, j_{i(v+T-1)})$ as input. With probability $1 - p_s$, we simply use the recent T items $(j_{i(t-T+1)}, \dots, j_{it})$ as input. If the sequence s_i is already shorter than T , then we always use the recent input sequence for user i .

Our proposed SSE-PT++ model can work almost as well as SSE-PT with a much smaller T . One can see in Table 6.2 with $T = 100$, SSE-PT++ can perform almost as well as SSE-PT. The time complexity of the SSE-PT model is of order $O(T^2d + Td^2)$. Therefore, reducing T by one half would lead to a theoretically 4x speed-up in terms of the training and inference speeds. As to the model’s space complexity, both SSE-PT and SSE-PT++ are of order $O(nd_u + md_i + Td + d^2)$.

6.4 Experiments

In this section, we compare our proposed algorithms, Personalized Transformer (SSE-PT) and SSE-PT++, with other state-of-the-art algorithms on real-world datasets. We implement our codes in Tensorflow and conduct all our experiments on a server with 40-core Intel Xeon E5-2630 v4 @ 2.20GHz CPU, 256G RAM and Nvidia GTX 1080 GPUs.

Datasets We use 5 datasets. The first 4 have exactly the same train/dev/test splits as in [56]. The datasets are: *Beauty* and *Games* categories from Amazon product review datasets¹; *Steam* dataset introduced in [56], which contains reviews crawled from a large

¹<http://jmcauley.ucsd.edu/data/amazon/>

video game distribution platform; *Movielens1M* dataset [38], a widely used benchmark datasets containing one million user movie ratings; *Movielens10M* dataset with ten million user ratings cleaned by us. Detailed dataset statistics are given in Table 8.7. One can easily see that the first 3 datasets have short sequences (average length ≈ 12) while the last 2 datasets have very long sequences ($\approx 10\times$ longer).

Evaluation Metrics The evaluation metrics we use are standard ranking metrics, namely NDCG and Recall for top recommendations (See Appendix). We follow the same evaluation setting as the previous paper [56]: predicting ratings at time point $t + 1$ given the previous t ratings. For a large dataset with numerous users and items, the evaluation procedure would be slow because (8.14) would require computing the ranking of all items based on their predicted scores for every single user. As a means of speed-up evaluations, we sample a fixed number C (e.g., 100) of negative candidates while always keeping the positive item that we know the user will engage next. This way, both R_{ij} and Π_i will be narrowed down to a small set of item candidates, and prediction scores will only be computed for those items through a single forward pass of the neural network.

Ideally, we want both NDCG and Recall to be as close to 1 as possible, because $\text{NDCG}@K = 1$ means the positive item is always put on the top-1 position of the top- K ranking list, and $\text{Recall}@K = 1$ means the positive item is always contained by the top- K recommendations the model makes.

Baselines We include 5 non-deep-learning and 6 deep-learning algorithms in our comparisons.

Non-deep-learning Baselines The simplest baseline is *PopRec*, basically ranking items according to their popularity. More advanced methods such as matrix factorization based baselines include Bayesian personalized ranking for implicit feedback [91], namely *BPR*; Factorized Markov Chains and Personalized Factorized Markov Chains models [92] also known as *FMC* and *PFMC*; and translation based method [40] called *TransRec*.

Deep-learning Baselines Recent years have seen many advances in deep learning for sequential recommendations. *GRU4Rec* is the first RNN-based method proposed for this problem [43]; *GRU4Rec+* [42] later is proposed to address some shortcomings of the

Table 6.1. Comparing various state-of-the-art temporal collaborative ranking algorithms on various datasets. The (A) to (E) are non-deep-learning methods, the (F) to (K) are deep-learning methods and the (L) to (O) are our variants. We did not report SSE-PT++ results for beauty, games and steam, as the input sequence lengths are very short (see Table 8.7), so there is no need for SSE-PT++.

DATASET	BEAUTY		GAMES		STEAM		ML-1M	
METRIC	RECALL@10	NDCG@10	RECALL@10	NDCG@10	RECALL@10	NDCG@10	RECALL@10	NDCG@10
(A) POPREC	0.4003	0.2277	0.4724	0.2779	0.7172	0.4535	0.4329	0.2377
(B) BPR	0.3775	0.2183	0.4853	0.2875	0.7061	0.4436	0.5781	0.3287
(C) FMC	0.3771	0.2477	0.6358	0.4456	0.7731	0.5193	0.6983	0.4676
(D) FPMC	0.4310	0.2891	0.6802	0.4680	0.7710	0.5011	0.7599	0.5176
(E) TRANSREC	0.4607	0.3020	0.6838	0.4557	0.7624	0.4852	0.6413	0.3969
(F) GRU4REC	0.2125	0.1203	0.2938	0.1837	0.4190	0.2691	0.5581	0.3381
(G) STAMP	0.4607	0.3020	0.6838	0.4557	0.7624	0.4852	0.6413	0.3969
(H) GRU4REC+	0.3949	0.2556	0.6599	0.4759	0.8018	0.5595	0.7501	0.5513
(I) CASER	0.4264	0.2547	0.5282	0.3214	0.7874	0.5381	0.7886	0.5538
(J) SASREC	0.4837	0.3220	0.7434	0.5401	0.8732	0.6293	0.8233	0.5936
(K) HGN	0.4469	0.2994	0.7164	0.5209	0.7426	0.4871	0.7584	0.5241
(L) SSE-SASREC	0.4878	0.3342	0.7517	0.5535	0.8697	0.6333	0.8230	0.5995
(M) PT	0.3954	0.2449	0.6427	0.4434	0.7535	0.4853	0.7658	0.5241
(N) SSE-PT	0.5028	0.3370	0.7757	0.5660	0.8772	0.6378	0.8341	0.6281
(O) SSE-PT++	–	–	–	–	–	–	0.8389	0.6292

initial version. *Caser* is the corresponding CNN-based method [108]. *STAMP* [68] utilizes the attention mechanism without using RNN or CNN as building blocks. Very recently, *SASRec* utilizes state-of-art Transformer encoder [110] with self-attention mechanisms. Hierarchical gating networks, also known as *HGN* [69] are also proposed to solve this problem.

Experiment Setup We use the same datasets as in [56] and follow the same procedure in the paper: use last items for each user as test data, second-to-last as validation data and the rest as training data. We implemented our method in Tensorflow and solve it with Adam Optimizer [57] with a learning rate of 0.001, momentum exponential decay rates $\beta_1 = 0.9, \beta_2 = 0.98$ and a batch size of 128. In Table 6.1, since we use the same data, the performance of previous methods except STAMP have been reported in [56]. We tune the dropout rate, and SSE probabilities p_u, p_i, p_y for input user/item embeddings and

Table 6.2. Comparing SASRec, SSE-PT and SSE-PT++ on Movielens1M Dataset while varying the maximum length allowed and dimension of embeddings.

METHODS	NDCG@10	RECALL@10	MAX LEN	USER DIM	ITEM DIM
SASREC	0.5769	0.8045	100	N/A	100
SASREC	0.5936	0.8233	200	N/A	50
SASREC	0.5919	0.8202	200	N/A	100
SSE-PT	0.6142	0.8212	100	50	100
SSE-PT	0.6191	0.8358	200	50	50
SSE-PT	0.6281	0.8341	200	50	100
SSE-PT++	0.6186	0.8318	100	50	100
SSE-PT++	0.6208	0.8358	200	50	50
SSE-PT++	0.6292	0.8389	200	50	100

Table 6.3. Comparing Different Regularizations for SSE-PT on Movielens1M Dataset. NO REG stands for no regularization. PS stands for parameter sharing across all users while PS(AGE) means PS is used within each age group. SASRec is added to last row after all SSE-PT results as a baseline.

REGULARIZATION	NDCG@5	% GAIN	RECALL@5	% GAIN
NO REG (BASELINE)	0.4855	-	0.6500	-
PS	0.5065	4.3	0.6656	2.4
PS (JOB)	0.4938	1.7	0.6570	1.1
PS (GENDER)	0.5110	5.3	0.6672	2.6
PS (AGE)	0.5133	5.7	0.6743	3.7
l_2	0.5149	6.0	0.6786	4.4
DROPOUT	0.5165	6.4	0.6823	5.0
l_2 + DROPOUT	0.5293	9.0	0.6921	6.5
SSE-SE	0.5393	11.1	0.6977	7.3
l_2 + SSE-SE + DROPOUT	0.5870	20.9	0.7442	14.5
SASREC (l_2 + DROPOUT)	0.5601		0.7164	

output embeddings on validation sets and report the best NDCG and Recall for top- K recommendations on test sets. For a fair comparison, we restrict all algorithms to use up to 50 hidden units for item embeddings. For the SSE-PT and SASRec models, we use the same number of transformer encoder blocks (i.e. $B = 2$) and set the maximum

length $T = 200$ for Movielens 1M and 10M dataset and $T = 50$ for other datasets. We use top- K with $K = 10$ and the number of negatives $C = 100$ in the evaluation procedure. In practice, using a different K and C does not affect our conclusions.

Comparisons One can easily see from Table 6.1 that our proposed SSE-PT has the best performance over all previous methods on all four datasets. On most datasets, our SSE-PT improves NDCG by more than 4% when compared with SASRec [56] and more than 20% when compared to non-deep-learning methods. SSE-SE, together with dropout and weight decay, is the best choice for regularization, which is evident from Table 6.3. SSE-SE is a more effective way to regularize our neural networks than any existent techniques including parameter sharing, dropout, weight decay. In practice, these SSE probabilities, just like dropout rate, can be treated as tuning parameters and easily tuned. Movielens10M results are left to Table 8.9 in the Appendix.

Temporal Collaborative Ranking Problem			
32::Twelve Monkeys (1995)::Drama, Sci-Fi	44::Mortal Kombat (1995)::Action, Adventure	39::Clueless (1995)::Comedy, Romance	11::American President, The (1995)::Comedy, Drama, Romance
23::Assassins (1995)::Thriller	42::Dead Presidents (1995)::Action, Crime, Drama	53::Lamerica (1994)::Drama	17::Sense and Sensibility (1995)::Drama, Romance
28::Persuasion (1995)::Romance	49::When Night Is Falling (1995)::Drama, Romance	2::Jumanji (1995)::Adventure, Children's, Fantasy	30::Shanghai Triad (1995)::Drama
38::It Takes Two (1995)::Comedy	19::Ace Ventura: When Nature Calls (1995)::Comedy	14::Nixon (1995)::Drama	34::Babe (1995)::Children's, Comedy, Drama
25::Leaving Las Vegas (1995)::Drama, Romance	12::Dracula: Dead and Loving It (1995)::Comedy, Horror	50::Usual Suspects, The (1995)::Crime, Thriller	41::Richard III (1995)::Drama, War
37::Across the Sea of Time (1995)::Documentary	15::Cutthroat Island (1995)::Action, Adventure, Romance	51::Guardian Angel (1994)::Action, Drama, Thriller	5::Father of the Bride Part II (1995)::Comedy
4::Waiting to Exhale (1995)::Comedy, Drama	43::Restoration (1995)::Drama	16::Casino (1995)::Drama, Thriller	31::Dangerous Minds (1995)::Drama
8::Tom and Huck (1995)::Adventure, Children's	18::Four Rooms (1995)::Thriller	21::Get Shorty (1995)::Action, Comedy, Drama	36::Dead Man Walking (1995)::Drama
48::Pocahontas (1995)::Animation, Children's, Musical, Romance	40::Cry, the Beloved Country (1995)::Drama	47::Seven (Se7en) (1995)::Crime, Thriller	33::Wings of Courage (1995)::Adventure, Romance
1::Toy Story (1995)::Animation, Children's, Comedy	46::How to Make an American Quilt (1995)::Drama, Romance	6::Heat (1995)::Action, Crime, Thriller	35::Carrington (1995)::Drama, Romance
22::Copycat (1995)::Crime, Drama, Thriller	27::Now and Then (1995)::Drama	9::Sudden Death (1995)::Action	Predict Next:
45::To Die For (1995)::Comedy, Drama	3::Grumpier Old Men (1995)::Comedy, Romance	13::Balto (1995)::Animation, Children's	26::Othello (1995)::Drama
10::GoldenEye (1995)::Action, Adventure, Thriller	7::Sabrina (1995)::Comedy, Romance	29::City of Lost Children, The (1995)::Adventure, Sci-Fi	
52::Mighty Aphrodite (1995)::Comedy	20::Money Train (1995)::Action	24::Powder (1995)::Drama, Sci-Fi	
Top-5 Recommendations by SASRec		Top-5 Recommendations by Our PT	
480::Jurassic Park (1993)::Action, Adventure, Sci-Fi		26::Othello (1995)::Drama	
3264::Buffy the Vampire Slayer (1992)::Comedy, Horror		2181::Marnie (1964)::Thriller	
582::Metisse (1993)::Comedy		1055::Shadow Conspiracy (1997)::Thriller	
3085::The Living Dead Girl (1982)::Horror		468::The Englishman Who Went Up a Hill, But Came Down a Mountain (1995)::Comedy, Romance, Drama	
649::Cold Fever (1994)::Comedy, Drama		629::Rude (1995)::Drama	
Attention Heat Map of SASRec		Attention Heat Map of Our PT	
32 23 28 38 25 37 4 8 48 1 22 45 10 52 44 42 49 19 12 15 43 18 40 46 27 3 7 20 39 53 2 14 50 51 16 21 47 6 9 13 29 24 11 17 30 34 41 5 31 36 33 35		32 23 28 38 25 37 4 8 48 1 22 45 10 52 44 42 49 19 12 15 43 18 40 46 27 3 7 20 39 53 2 14 50 51 16 21 47 6 9 13 29 24 11 17 30 34 41 5 31 30 33 35	

Figure 6.2. Illustration of how SASRec (Left) and SSE-PT (Right) differs on utilizing the Engagement History of A Random User in Movielens1M Dataset.

6.4.1 Attention Mechanism Visualization

Apart from evaluating our SSE-PT against SASRec using well-defined ranking metrics on real-world datasets, we also visualize the differences between both methods in terms of their attention mechanisms. In Figure 6.2, a random user’s engagement history in Movielens1M dataset is given in temporal order (column-wise). We hide the last item whose index is 26 in test set and hope that a temporal collaborative ranking model can figure out item-26 is the one this user will watch next using only previous engagement history. One can see for a typical user; they tend to look at a different style of movies at different times. Earlier on, they watched a variety of movies, including Sci-Fi, animation, thriller, romance, horror, action, comedy and adventure. But later on, in the last two columns of Figure 6.2, drama and thriller are the two types they like to watch most, especially the drama type. In fact, they watched 9 drama movies out of recent 10 movies. For humans, it is natural to reason that the hidden movie should probably also be drama type. So what about the machine’s reasoning?

For our SSE-PT, the hidden item indexed 26 is put in the first place among its top-5 recommendations. Intelligently, the SSE-PT recommends 3 drama movies, 2 thriller movies and mixing them up in positions. Interestingly, the top recommendation is ‘Othello’, which like the recently watched ‘Richard III’, is an adaptation of a Shakespeare play, and this dependence is reflected in the attention weight. On the contrast, SASRec cannot provide top-5 recommendations that are personalized enough. It recommends a variety of action, Sci-Fi, comedy, horror, and drama movies but none of them match item-26. Although this user has watched all these types of movies in the past, they do not watch these anymore as one can easily tell from his recent history. Unfortunately, SASRec cannot capture this and does not provide personalized recommendations for this user by focusing more on drama and thriller movies. It is easy to see that in contrast, our SSE-PT model shares with human reasoning that more emphasis should be placed on recent movies.

6.4.2 Training Speed

In [56], it has been shown that SASRec is about 11 times faster than Caser and 17 times faster than GRU4Rec⁺ and achieves much better NDCG@10 results so we did

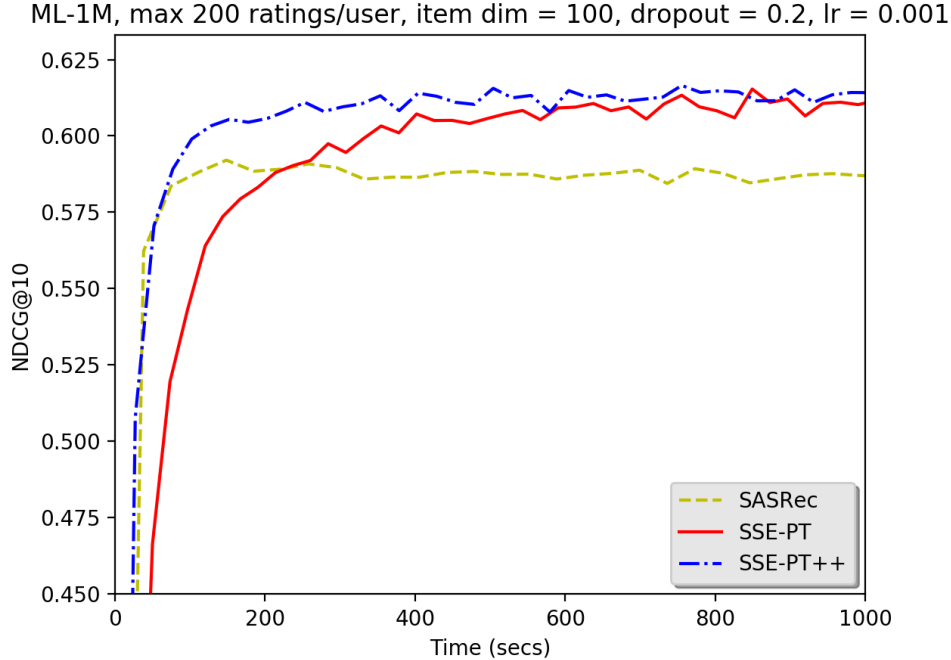


Figure 6.3. Illustration of the speed of SSE-PT

not include Caser and GRU4Rec⁺ in our comparisons. In Figure 6.3, we only compare the training speeds and ranking performances among SASRec, SSE-PT and SSE-PT++ for Movielens1M dataset. Given that we added additional user embeddings into our SSE-PT model, it is expected that it will take slightly longer to train our model than un-personalized SASRec. We find empirically that training speed of the SSE-PT and SSE-PT++ model are comparable to that of SASRec, with SSE-PT++ being the fastest and the best performing model. It is clear that our SSE-PT and SSE-PT++ achieve much better ranking performances than our baseline SASRec using the same training time.

6.4.3 Ablation Study

SSE probability Given the importance of SSE regularization for our SSE-PT model, we carefully examined the SSE probability for input user embedding in Table 8.10 in Appendix. We find that the appropriate hyper-parameter SSE probability is not very sensitive: anywhere between 0.4 and 1.0 gives good results, better than parameter sharing and not using SSE-SE. This is also evident based on comparison results in Table 6.3.

Sampling Probability Recall that the sampling probability is unique to our SSE-PT++ model. We show in Table 8.11 in Appendix using an appropriate sampling probability like $0.2 \rightarrow 0.3$ would allow it to outperform SSE-PT when the same maximum length is used.

Number of Attention Blocks We find for our SSE-PT model, a larger number of attention blocks is preferred. One can easily see in Table 8.12 in Appendix, the optimal ranking performances are achieved at $B = 4$ or 5 for Movielens1M dataset and at $B = 6$ for Movielens10M dataset.

Personalization and Number of Negatives Sampled Based on the results in Table 8.13 in Appendix, we are positive that the personalized model always outperforms the un-personalized one when we use the same regularization techniques. This holds true regardless of how many negatives sampled or what ranking metrics are used during evaluation.

6.5 Conclusion

In this chapter, we propose a novel neural network architecture called Personalized Transformer for the temporal collaborative ranking problem. It enjoys the benefits of being a personalized model, therefore achieving better ranking results for individual users than the current state-of-the-art. By examining the attention mechanisms during inference, the model is also more interpretable and tends to pay more attention to recent items in long sequences than un-personalized deep learning models.

Chapter 7

Conclusions

7.1 Summary of Contributions

The goal of this thesis is to cover some recent advances in collaborative filtering and ranking research and demonstrate that there are various types of orthogonality in which one can contribute to the field. During my PhD study, I was fortunate to explore many directions within collaborative filtering and ranking research. For the first 2 years of my research, I conducted foundational collaborative ranking research on improving the optimization procedure of the pairwise ranking loss in chapter 3 and designing new listwise loss objective functions in chapter 4 for collaborative ranking. For my last 2 years of PhD, I was exploring directions on incorporating additional side information such as graphs and temporal orderings. In chapter 5, we came up with a new graph encoding method in chapter 2 to enhance existing graph-based collaborative filtering, allowing them to encode deep graph information and therefore achieve better recommendation performances. We made the temporal collaborative ranking model personalized in chapter 7 by incorporating user embeddings. In the process, motivated by the need to prevent over-fitting caused by the additional parameters, we introduced a general regularization technique for embedding layers in deep learning in chapter 6, which was shown to be useful for many other models with lots of embedding parameters both within and outside recommendations.

7.2 Future Work

Despite that the dissertation is lengthy and has contained many important research directions, it is by no means complete and I feel there are still many interesting and important directions worth pursuing but not done within this dissertation. I imagine at the end of day, assuming the computing power catches up, it is very likely we can discard most of the feature engineering and directly learn from raw data formats such as texts, images and videos, just like what happened in computer vision and natural language processing fields with the advances of deep learning techniques. This has not yet happened in neither academics or industry. To make this happen, I believe that there are a few shortcomings in current research practices:

- There do not exist unified metrics to evaluate across papers. Some papers use root mean square errors (RMSE) as accuracy metric while others use ranking metrics. Then even within the ranking metrics, people tend to use different metrics: from AUC score to NDCG, precision and recall, not to mention that different top k can be used for NDCG, precision and recall.
- There do not exist unified datasets to test on across papers. There are many datasets outside there, from Netflix to Movielens, and Douban to Flixster, etc.. Different papers tend to use different training-test splits: some do random splits while others do splits based on temporal orderings. Moreover, the recommender systems community does not really have a shared and well-maintained leader-board on the winning solutions. The results of various proposed algorithms are very likely to perform differently on different datasets.
- Most datasets do not come with good features nor complete raw input data. To some extent, it is understandable because of the strict privacy and copyright concerns. But this has put academic researches at a disadvantage and often industry applied research will not take as much risk to pursue long-term projects. On the other hand, academics are able to take more risks to pursue longer-term research topics but unfortunately are limited by the constraints of good datasets and powerful computing resources.

While addressing these shortcomings, I think it would be very exciting to see different levels of interactions between the recommendation field and other AI fields, including natural language understanding and computer vision. News/Book recommendation, image recommendation, video recommendation and audio recommendation are some promising examples that may see breakthroughs of new models, just as what happened during Netflix competition about 10 years ago [6]. But to do that, we need a well-defined problem, dataset and metric, and lots of people participating both from academics and industry by combining strengths from both parties.

Chapter 8

Appendix

8.1 Appendix to Chapter 2

Algorithm 10. A Standard Bloom Filter

```
class BloomFilter:
    def constructor(self, c, {ht(·) : t = 1, ..., k}):
        self.b[i] = 0  ∀i = 1, ..., c
        self.ht = ht  ∀i = 1, ..., k

    def add(self, x):
        self.b[self.ht(x)] = 1  ∀t = 1, ..., k

    def union(self, bf):
        self.b[i] ← self.b[i] | bf.b[i]  ∀i = 1, ..., c

    def size(self):
        return ⌈ - $\frac{c}{k}$  log(1 -  $\frac{\text{nnz}(\text{self.b})}{c}$ ) ⌉
```

8.1.1 Simulation Study

In the simulation we carried out, we set the number of users $n = 10,000$ and the number of items $m = 2,000$. We uniformly sample 5% for training and 2% for testing out of the

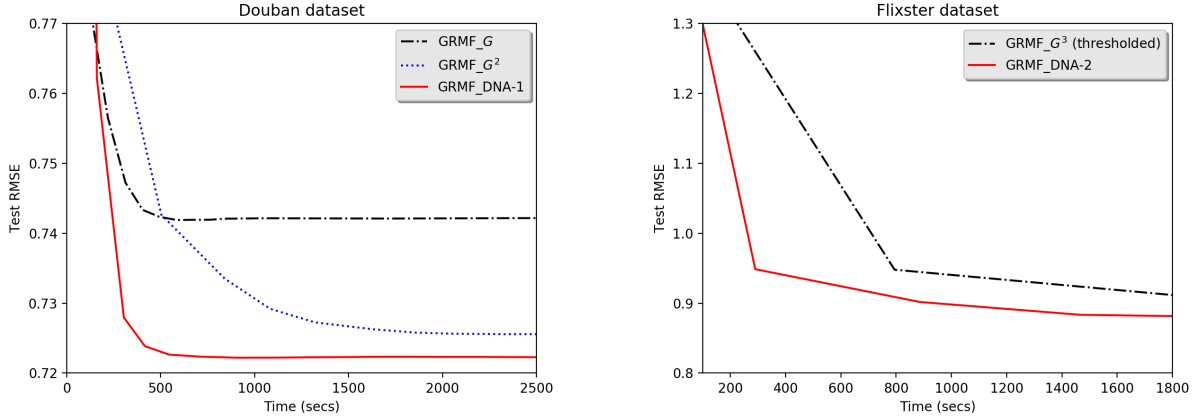


Figure 8.1. Compare Training Speed of GRMF, with and without Graph DNA.

total nm ratings. We choose $T = 3$ so the graph contains at most 6-hop information among n users. We use rank $r = 50$ for both user and item embeddings. We set influence weight $w = 0.6$, i.e. in each propagation step, 60% of one user’s preference is decided by its friends (i.e. neighbors in the friendship graph). We set $p = 0.001$, which is the probability for each of the possible edges being chosen in Erdős-Rényi graph G . A small edge probability p , influence weight $w < 1.0$, and a not too-large T is needed, because we don’t want that all users become more or less the same after T propagation steps.

8.1.2 Metrics

We omit the definitions of RMSE, Precision@ k , NDCG@ k , MAP as those can be easily found online. HLU: Half-Life Utility [11, 98] is defined as:

$$\text{HLU} = \frac{1}{n} \sum_{i=1}^n \text{HLU}_i, \quad (8.1)$$

where n is the number of users and HLU_i is given by:

$$\text{HLU}_i = \sum_{l=1}^k \frac{\max(R_{i\Pi_{il}} - d, 0)}{2^{(j-1)/(\alpha-1)}}, \quad (8.2)$$

where $R_{i\Pi_{il}}$ follows previous definition, d is the neural vote (usually the rating average), and α is the viewing halflife. The halflife is the number of the item on the list such that there is a 50-50 chance the user will review that item [11].

Algorithm 11. Simulation of Synthetic Data

Input: n users, m items, rank r , influence weight w , T propagation steps

Output: $R_{\text{tr}} \in \mathbb{R}^{n \times m}$, $R_{\text{te}} \in \mathbb{R}^{n \times m}$, $G \in \mathbb{R}^{n \times n}$

- 1: Randomly initialize $U \in \mathbb{R}^{n \times r}$, $V \in \mathbb{R}^{m \times r}$ from standard normal distribution
 - 2: Generate a random undirected Erdős-Rényi graph G with each edge being chosen with probability p
 - 3: **for** $t = 1, \dots, T$ **do**
 - 4: **for** $i = 1, \dots, n$ **do**
 - 5: $\tilde{U}_i = w \cdot \sum_{j:(i,j) \in G} U_j + (1 - w) \cdot U_i$
 - 6: Set $U = \tilde{U}$
 - 7: Generate rating matrix $R = UV^T$
 - 8: Random sample observed user/item indices in training and test data: $\Omega_{\text{tr}}, \Omega_{\text{te}}$
 - 9: Obtain $R_{\text{tr}} = \Omega_{\text{tr}} \circ R$, $R_{\text{te}} = \Omega_{\text{te}} \circ R$
 - 10: **return** rating matrices $R_{\text{tr}}, R_{\text{te}}$, user graph G
-

Table 8.1. Compare Bloom filters of different depths and sizes an on Synthesis Dataset. Note that the number of bits of Bloom filter is decided by Bloom filter’s maximum capacity and tolerable error rate (i.e. false positive error, we use 0.2 as default).

methods	max capacity	c bits	nmz ratio	RMSE ($\times 10^{-3}$)	% Relative Graph Gain
GRMF_G ²	-	-	-	2.6543	59.5903
GRMF_DNA-1	20	135	0.217	2.4303	163.8734
GRMF_DNA-1	50	336	0.093	2.4795	140.9683
GRMF_DNA-2	20	135	0.880	2.4921	135.1024
GRMF_DNA-2	50	336	0.608	2.4937	134.3575
GRMF_DNA-2	100	672	0.381	2.4510	154.2365
GRMF_DNA-2	200	1,341	0.215	2.4541	152.7933
GRMF_DNA-3	200	1,341	0.874	2.4667	146.9274
GRMF_DNA-3	600	4,020	0.525	2.4572	151.3500
GRMF_DNA-3	1,000	6,702	0.364	2.4392	159.7299
GRMF_DNA-3	1,500	10,050	0.262	2.4247	166.4804
GRMF_DNA-4	2,000	13,401	0.743	2.5532	106.6573
GRMF_DNA-4	4,000	26,799	0.499	2.4466	156.2849

Table 8.2. Compare nnz of different methods on Douban and Flixster datasets. GRMF_ G^4 and GRMF_DNA-2 are using the same 4-hop information in the graph but in different ways. Note that we do not exclude potential overlapping among columns.

Dataset	methods	R_{tr}	G	G^2	G^3	G^4	B	total nnz
Douban	MF	9,803,098	-	-	-	-	-	9,803,098
	GRMF_ G	9,803,098	1,711,780	-	-	-	-	11,514,878
	GRMF_ G^2	9,803,098	1,711,780	106,767,776	-	-	-	118,282,654
	GRMF_ G^3	9,803,098	1,711,780	106,767,776	2,313,572,544	-	-	2,431,855,198
	GRMF_ G^4	9,803,098	1,711,780	106,767,776	2,313,572,544	8,720,553,105	-	11,152,408,303
	GRMF_DNA-1	9,803,098	0	-	-	-	8,834,740	18,637,838
	GRMF_DNA-2	9,803,098	1,711,780	-	-	-	142,897,900	154,412,778
	GRMF_DNA-3	9,803,098	1,711,780	-	-	-	928,159,604	939,674,482
Flixster	MF	3,619,304	-	-	-	-	-	3,619,304
	GRMF_ G	3,619,304	2,538,746	-	-	-	-	6,158,050
	GRMF_ G^2	3,619,304	2,538,746	130,303,379	-	-	-	136,461,429
	GRMF_ G^3	3,619,304	2,538,746	130,303,379	2,793,542,551	-	-	3,060,307,359
	GRMF_ G^4	3,619,304	2,538,746	130,303,379	2,793,542,551	12,691,844,513	-	15,752,151,872
	GRMF_DNA-1	3,619,304	0	-	-	-	12,664,952	16,284,256
	GRMF_DNA-2	3,619,304	2,538,746	-	-	-	181,892,883	188,050,933
	GRMF_DNA-3	3,619,304	2,538,746	-	-	-	1,185,535,529	1,191,693,579

8.1.3 Graph Regularized Weighted Matrix Factorization for Implicit feedback

We use the rank $r = 10$, negatives' weight $\rho = 0.01$ and measure the prediction performance with metrics MAP, HLU, Precision@ k and NDCG@ k (see definitions of metrics in Appendix 8.1.2).

We follow the similar procedure to what is done before in GRMF and co-factor: we run all combinations of tuning parameters of $\lambda_l \in \{0.01, 0.1, 1, 10, 100\}$ and $\lambda_g \in \{0.01, 0.1, 1, 10, 100\}$ for each method on validation data for fixed number 40 epochs and choose the best combination as the parameters to use on test data. We then report the best prediction results during first 40 epochs on test data with the chosen parameter combination.

8.1.4 Reproducibility

To reproduce results reported in the paper, one need to download data (douban and flixster) and third-party C++ Matrix Factorization library from the link <https://www.csie.ntu.edu.tw/~cjlin/papers/ocmf-side/>. One can simply follow README there

to compile the codes in Matlab and run one-class matrix factorization library in different modes (both explicit feedback and implicit feedback works). The advantage of using this library is that the codes support multi-threading and runs quite fast with very efficient memory space allocations. It also supports with graph or other side information. All three methods' baseline can be simply run with the tuning parameters we reported in the Table 8.4, 8.5, 8.6 in Appendix.

To reproduce results of our DNA methods, one need to generate Bloom filter matrix B following Algorithm 1. We will provide our python codes implementing Algorithm 1 and Matlab codes converting into the formats the library requires.

For baselines and our DNA methods, We perform a parameter sweep for $\lambda_l \in \{0.01, 0.1, 1, 10, 100\}$, $\lambda_g \in \{0.01, 0.1, 1, 10, 100\}$, $\alpha \in \{0.0001, 0.001, 0.01, 0.1, 0.3, 0.7, 1\}$, for $\beta \in \{0.005, 0.01, 0.03, 0.05, 0.1\}$ when needed. We run all combinations of tuning parameters for each method on validation set for 40 epochs and choose the best combination as the parameters to use on test data. We then report the best test RMSE in first 40 epochs on test data with the chosen parameter combination. We provide all the chosen combinations of tuning parameters that achieves reported optimal results in results tables in the Table 8.4, 8.5, 8.6 in Appendix. One just need to exactly follow our procedures in Section 2.4 to construct new \dot{G}, \dot{U} to replace the G, U in baseline methods before feeding into Matlab.

As to simulation study, we will also provide python codes to repeat our Algorithm 11 to generate synthesis dataset. One can easily simulate the data before converting into Matlab data format and running the codes as before. The optimal parameters can be found in Table 8.3. For all the methods, we select the best parameters λ_l and λ_g from $\{0.01, 0.1, 1, 10, 100\}$. For method GRMF- G^2 , we tune an additional parameter $\alpha \in \{0.0001, 0.001, 0.01, 0.1, 0.3, 0.7, 1\}$. For the thrid-order method GRMF- G^3 , we tune $\beta \in \{0.005, 0.01, 0.03, 0.05, 0.1\}$ in addition to $\lambda_l, \lambda_g, \alpha$. Due to the speed constraint, we are not able to tune a broader range of choices for α and β as it is too time-consuming to do so especially for douban and flixster datasets. For example, it takes takes about 3 weeks using 16-cores CPU to tune both α, β on flixster dataset. We run each method with

every possible parameter combination for fixed 80 epochs on the same training data, tune the best parameter combination based on a small predefined validation data and report the best RMSE results on test data with the best tuning parameters during the first 80 epochs. Note that only on the small synthesis dataset, we calculate full G^3 and report the results. On real datasets, there is no way to calculate full G^4 to utilize the complete 4-hop information, because one can easily spot in Table 8.2 the number of non-zero elements (nnz) is growing exponentially when the hop increases by 1, which makes it impossible for one to utilize complete 3-hop and 4-hop information.

In Table 8.4, one can compare magnitude of optimal α and β to have a good idea of whether G or G^2 is more useful. G represents shallow graph information and G^2 represents deep graph information. If one already run GRMF_ G^2 , one can then use this as a preliminary test to decide whether to go deep with DNA-3 ($d = 3$) to capture deep graph information or simply go ahead with DNA-1 ($d = 1$) to fully utilize shallow information. For douban dataset, we have $\alpha = 0.05 > 0.0005 = \beta$, which implies shallow information is important and we should fully utilize it. It explains why DNA-1 is performing well both in terms of performance and speed on douban dataset. It is worth noting that GRMF_DNA-1's Bloom filter matrix B contains much more nnz than that of G in Table 8.2 though 20% less than that of G^2 . On the other hand, for flixster dataset, we have $\alpha = 0.01 < 0.1 = \beta$, which implies in this dataset deeper information is more important and we should go deeper. That explains why here GRMF_DNA-3 (6-hop) achieves about 10 times more gain than using 1-hop GRMF_ G .

8.1.5 Code

Part of our code is already made available on Github: <https://github.com/wuliwei9278/Graph-DNA>.

8.2 Appendix to Chapter 4

We include pseudo-codes for Algorithm 12, 13 and Figures 8.2, 8.3, 8.4 in the appendix to chapter 4.

Table 8.3. Compare Matrix Factorization for Explicit Feedback on Synthesis Dataset. The synthesis dataset has 10,000 users and 2,000 items with user friendship graph of size $10,000 \times 10,000$. Note that the graph only contains at most 6-hop valid information. GRMF_ G^6 means GRMF with $G + \alpha \cdot G^2 + \beta \cdot G^3 + \gamma \cdot G^4 + \epsilon \cdot G^5 + \omega \cdot G^6$. GRMF_DNA- d means depth d is used.

methods	test RMSE ($\times 10^{-3}$)	λ_l	λ_g	α	β	γ	ϵ	ω	% gain over baseline
MF	2.9971	0.01	-	-	-	-	-	-	-
GRMF_ G	2.7823	0.01	0.01	-	-	-	-	-	7.16693
GRMF_ G^2	2.6543	0.01	0.01	0.3	-	-	-	-	11.43772
GRMF_ G^3	2.5687	0.01	0.01	0.01	0.05	-	-	-	14.29382
GRMF_ G^4	2.5562	0.01	0.01	0.01	0.05	0.1	-	-	14.71088
GRMF_ G^5	2.4853	0.01	0.01	0.01	0.05	0.1	0.1	-	17.07651
GRMF_ G^6	2.4852	0.01	0.01	0.01	0.05	0.1	0.1	0.01	17.07984
GRMF_DNA-1	2.4303	0.01	0.01	-	-	-	-	-	18.91161
GRMF_DNA-2	2.4510	0.01	0.01	-	-	-	-	-	18.22095
GRMF_DNA-3	2.4247	0.01	0.01	-	-	-	-	-	19.09846
GRMF_DNA-4	2.4466	0.01	0.01	-	-	-	-	-	18.36776

Table 8.4. Compare Matrix Factorization methods for Explicit Feedback on Douban and Flixster data. We use rank $r = 10$.

Dataset	methods	test RMSE ($\times 10^{-1}$)	λ_l	λ_g	α	β	% gain over baseline
Douban	MF	7.3107	1	-	-	-	-
	GRMF_ G	7.2398	0.1	100	-	-	0.9698
	GRMF_ G^2	7.2381	0.1	100	0.001	-	0.9930
	GRMF_ G^3 (full)	7.2432	0.1	100	0.05	0.0005	0.9350
	GRMF_ G^3 (thresholded)	7.2382	0.1	100	0.05	0.0005	0.9917
	GRMF_DNA-1	7.2191	0.1	100	-	-	1.2689
	GRMF_DNA-2	7.2359	1	10	-	-	1.0232
	GRMF_DNA-3	7.2095	0.01	100	-	-	1.3843
Flixster	MF	8.8111	0.1	1	-	-	-
	GRMF_ G	8.8049	0.01	1	-	-	0.0704
	GRMF_ G^2	8.7849	0.01	1	0.05	-	0.2974
	GRMF_ G^3 (full)	8.7932	0.1	1	0.01	0.1	0.2032
	GRMF_ G^3 (thresholded)	8.7920	0.01	1	0.01	0.1	0.2168
	GRMF_DNA-1	8.8013	0.01	1	-	-	0.1112
	GRMF_DNA-2	8.8007	0.1	1	-	-	0.1180
	GRMF_DNA-3	8.7453	0.1	100	-	-	0.7468

Table 8.5. Compare Co-factor Methods for Explicit Feedback on Douban and Flixster Datasets. We use rank $r = 10$ for both methods.

Dataset	methods	test RMSE ($\times 10^{-1}$)	λ_l	% gain over baseline
Douban	co-factor_ G	7.2743	1	-
	co-factor_DNA-3	7.2674	1	0.5923
Flixster	co-factor_ G	8.7957	0.01	-
	co-factor_DNA-3	8.7354	0.01	0.8591

Table 8.6. Compare Weighted Matrix Factorization with Graph for Implicit Feedback on Douban and Flixster Datasets. We use rank $r = 10$ for both methods and all metric results are in %.

Dataset	Methods	MAP	HLU	P@1	P@5	NDCG@1	NDCG@5	λ_l	λ_g
Douban	WMF_ G	8.340	13.033	14.944	10.371	14.944	12.564	0.01	10
	WMF_DNA-3	8.400	13.110	14.991	10.397	14.991	12.619	1	1
Flixster	WMF_ G	10.889	14.909	12.303	7.9927	12.303	12.734	10	0.1
	WMF_DNA-3	11.612	15.687	12.644	8.1583	12.644	13.399	1	1

8.3 Appendix to Chapter 5

For experiments in Section 5.4.1, we use Julia and C++ to implement SGD. For experiments in Section 5.4.2, and Section 5.4.4, we use Tensorflow and SGD/Adam Optimizer. For experiments in Section 5.4.3, we use Pytorch and Adam with noam decay scheme and warm-up. We find that none of these choices affect the strong empirical results supporting the effectiveness of our proposed methods, especially the SSE-SE. In any deep learning frameworks, we can introduce stochasticity to the original embedding look-up behaviors and easily implement SSE-Layer in Figure 5.1 as a custom operator.

8.3.1 Neural Networks with One Hidden Layer

To run SSE-Graph, we need to construct good-quality knowledge graphs on embeddings. We managed to match movies in Movielens1m and Movielens10m datasets to IMDB websites, therefore we can extract plentiful information for each movie, such as the cast of the movies, user reviews and so on. For simplicity reason, we construct the knowledge graph on item-side embeddings using the cast of movies. Two items are connected by an edge when they share one or more actors/actresses. For user side, we do not have good quality graphs: we are only able to create a graph on users in Movielens1m dataset based on

Algorithm 12. Compute gradient for V when U fixed

Input: Π, U, V, λ, ρ

Output: g $\triangleright g \in \mathbb{R}^{r \times m}$ is the gradient for $f(V)$

$g = \lambda \cdot V$

for $i = 1$ **to** n **do**

Precompute $h_t = u_i^T v_{\Pi_{it}}$ for $1 \leq t \leq \bar{m}$ \triangleright For implicit feedback, it should be $(1 + \rho) \cdot \tilde{m}$ instead of \tilde{m} , since $\rho \cdot \tilde{m}$ 0's are appended to the back

Initialize $total = 0, tt = 0$

for $t = \bar{m}$ **to** 1 **do**

$total += \exp(h_t)$

$tt += 1/total$

Initialize $c[t] = 0$ for $1 \leq t \leq \bar{m}$

for $t = \bar{m}$ **to** 1 **do**

$c[t] += h_t \cdot (1 - h_t)$

$c[t] += \exp(h_t) \cdot h_t \cdot (1 - h_t) \cdot tt$

$total += \exp(h_t)$

$tt -= 1/total$

for $t = 1$ **to** \bar{m} **do**

$g[:, \Pi_{it}] += c[t] \cdot u_i$

Return g

their age groups but we do not have any side information on users in Movielens10m dataset. When running experiments, we do a parameter sweep for weight decay parameter and then fix it before tuning the parameters for SSE-Graph and SSE-SE. We utilize different ρ and p for user and item embedding tables respectively. The optimal parameters are stated in Table 5.1 and Table 5.2. We use the learning rate of 0.01 in all SGD experiments.

In the first leg of experiments, we examine users with fewer than 60 ratings in Movielens1m and Movielens10m datasets. In this scenario, the graph should carry higher importance. One can easily see from Table 5.1 that without using graph information,

Algorithm 13. Gradient update for V (Same procedure for updating U)

Input: $V, ss, rate$ $\triangleright rate$ refers to the decaying rate of the step size ss
Output: V
Compute gradient g for V \triangleright see alg 12
 $V -= ss \cdot g$
 $ss *= rate$
Return V

our proposed SSE-SE is the best performing matrix factorization algorithms among all methods, including popular ALS-MF and SGD-MF in terms of RMSE. With Graph information, our proposed SSE-Graph is performing significantly better than the Graph Laplacian Regularized Matrix Factorization method. This indicates that our SSE-Graph has great potentials over Graph Laplacian Regularization as we do not explicitly penalize the distances across embeddings but rather we implicitly penalize the effects of similar embeddings on the loss.

In the second leg of experiments, we remove the constraints on the maximum number of ratings per user. We want to show that SSE-SE can be a good alternative when graph information is not available. We follow the same procedures in [115, 116]. In Table 5.2, we can see that SSE-SE can be used with dropout to achieve the smallest RMSE across Douban, Movielens10m, and Netflix datasets. In Table 5.3, one can see that SSE-SE is more effective than dropout in this case and can perform better than STOA listwise approach SQL-Rank [116] on 2 datasets out of 3.

In Table 5.2, SSE-SE has two tuning parameters: probability p_u to replace embeddings associated with user-side embeddings and probability p_i to replace embeddings associated with item side embeddings because there are two embedding tables. But here for simplicity, we use one tuning parameter $p_s = p_u = p_i$. We use dropout probability of p_d , dimension of user/item embeddings d , weight decay of λ and learning rate of 0.01 for all experiments, with the exception that the learning rate is reduced to 0.005 when both SSE-SE and Dropout are applied. For Douban dataset, we use $d = 10, \lambda = 0.08$. For Movielens10m

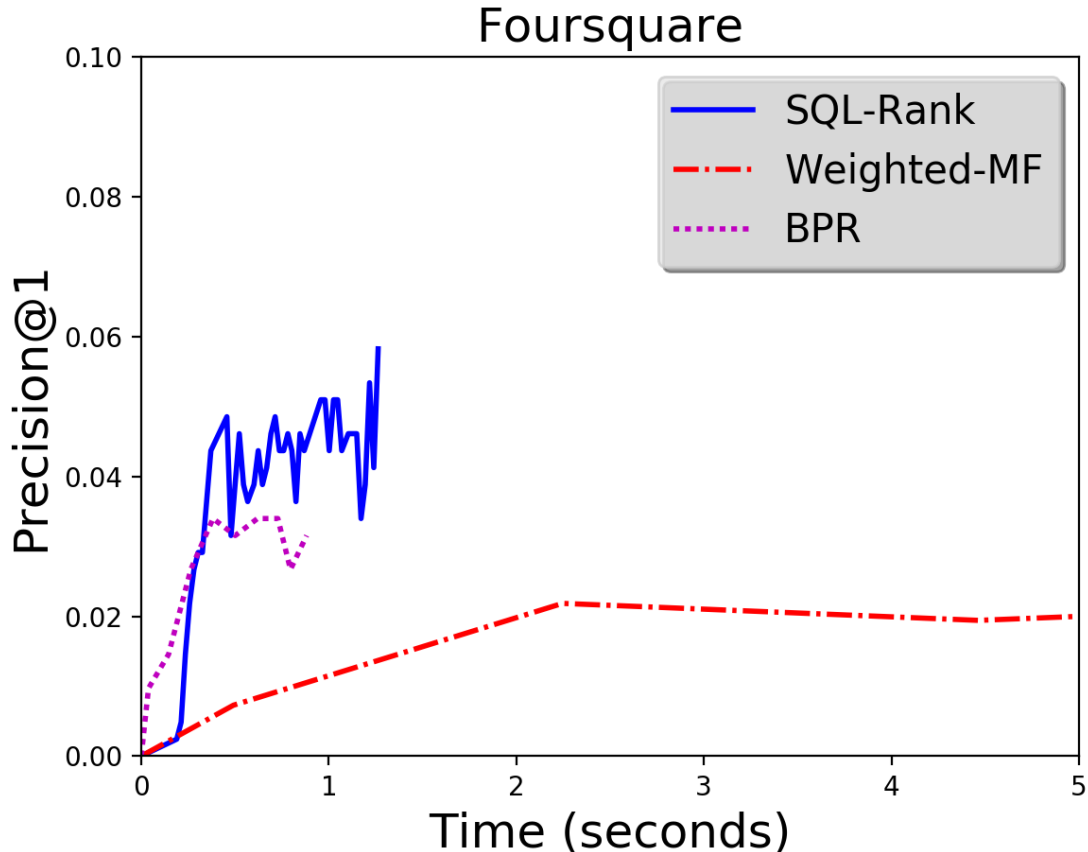


Figure 8.2. Comparing implicit feedback methods.

and Netflix dataset, we use $d = 50, \lambda = 0.1$.

8.3.2 Neural Machine Translation

We use the transformer model [110] as the backbone for our experiments. The control group is the standard transformer encoder-decoder architecture with self-attention. In the experiment group, we apply SSE-SE towards both encoder and decoder by replacing corresponding vocabularies' embeddings in the source and target sentences. We trained on the standard WMT 2014 English to German dataset which consists of roughly 4.5 million parallel sentence pairs and tested on WMT 2008 to 2018 news-test sets. Sentences were encoded into 32,000 tokens using a byte-pair encoding. We use the SentencePiece, OpenNMT and SacreBLEU implementations in our experiments. We trained the 6-layer

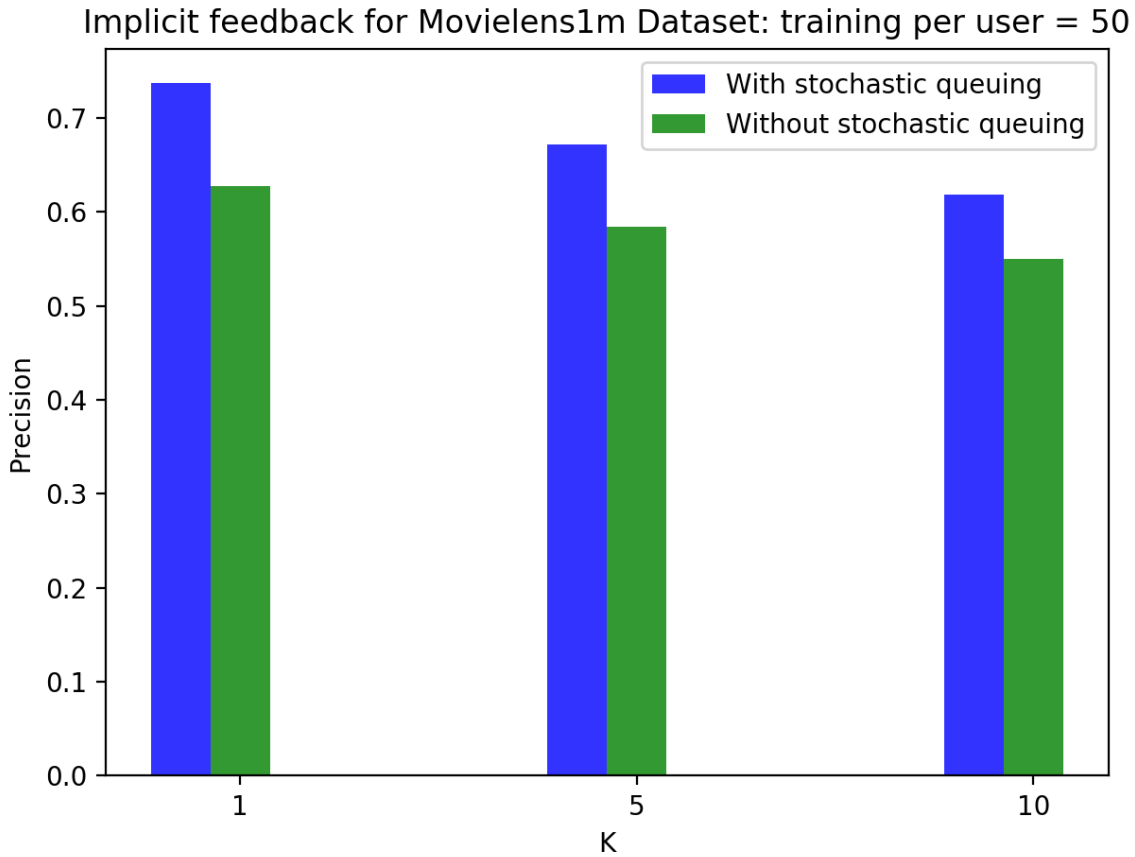


Figure 8.3. Effectiveness of Stochastic Queuing Process.

transformer base model on a single machine with 4 NVIDIA V100 GPUs for 20,000 steps. We use the same dropout rate of 0.1 and label smoothing value of 0.1 for the baseline model and our SSE-enhanced model. Both models have dimensionality of embeddings as $d = 512$. When decoding, we use beam search with the beam size of 4 and length penalty of 0.6 and replace unknown words using attention. For both models, we average last 5 checkpoints (we save checkpoints every 10,000 steps) and evaluate the model’s performances on the test datasets using BLEU scores. The only difference between the two models is whether or not we use our proposed SSE-SE with $p = 0.01$ in Equation 5.5 for both encoder and decoder embedding layers.

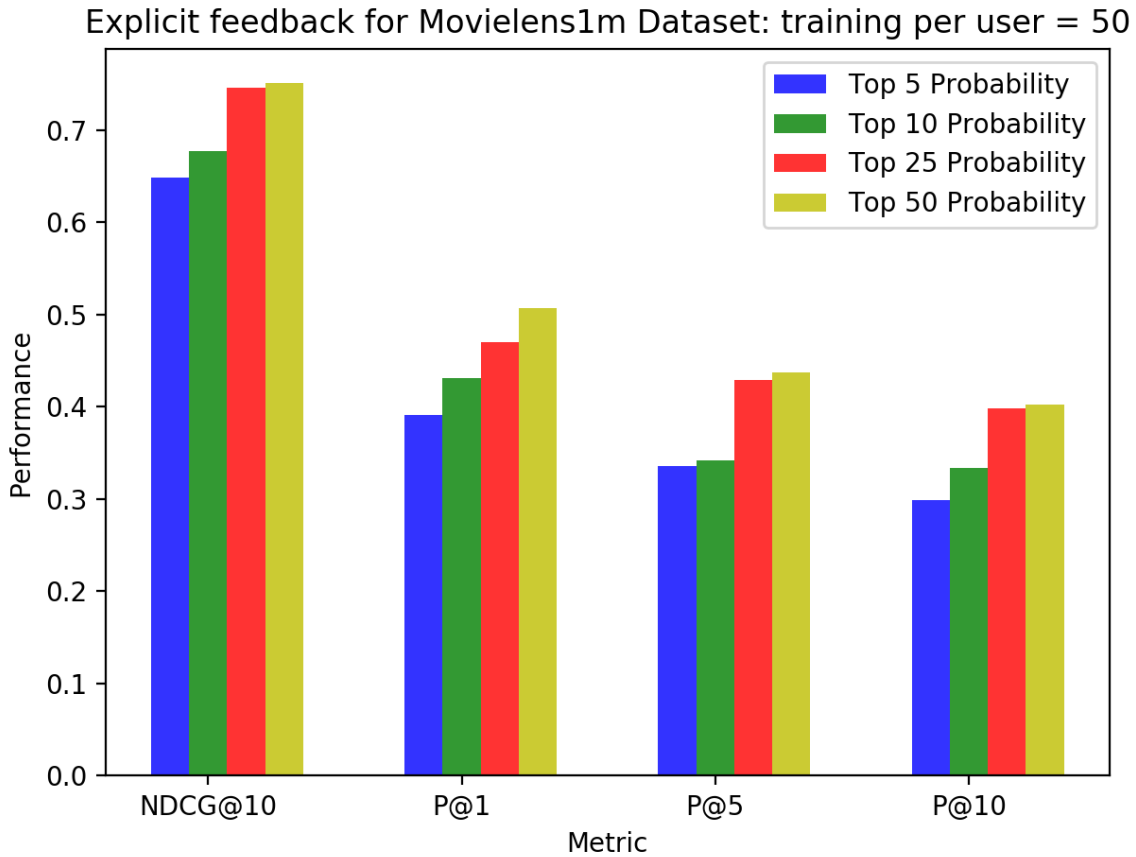


Figure 8.4. Effectiveness of using full lists.

8.3.3 BERT

In the first leg of experiments, we crawled one million user reviews data from IMDB and pre-trained the BERT-Base model (12 blocks) for 500,000 steps using sequences of maximum length 512 and batch size of 8, learning rates of $2e^{-5}$ for both models using one NVIDIA V100 GPU. Then we pre-trained on a mixture of our crawled reviews and reviews in IMDB sentiment classification tasks (250K reviews in train and 250K reviews in test) for another 200,000 steps before training for another 100,000 steps for the reviews in IMDB sentiment classification task only. In total, both models are pre-trained on the same datasets for 800,000 steps with the only difference being our model utilizes SSE-SE. In the second leg of experiments, we fine-tuned the two models obtained in the first-leg

experiments on two sentiment classification tasks: IMDB sentiment classification task and SST-2 sentiment classification task. The goal of pre-training on IMDB dataset but fine-tuning for SST-2 task is to explore whether SSE-SE can play a role in transfer learning.

The results are summarized in Table 5.6 for IMDB sentiment task. In experiments, we use maximum sequence length of 512, learning rate of $2e^{-5}$, dropout probability of 0.1 and we run fine-tuning for 1 epoch for the two pre-trained models we obtained before. For the Google pre-trained BERT-base model, we find that we need to run a minimum of 2 epochs. This shows that pre-training can speed up the fine-tuning. We find that Google pre-trained model performs worst in accuracy because it was only pre-trained on Wikipedia and books corpus while ours have seen many additional user reviews. We also find that SSE-SE pre-trained model can achieve accuracy of 0.9542 after fine-tuning for one epoch only. On the contrast, the accuracy is only 0.9518 without SSE-SE for embeddings associated with output y_i .

For the SST-2 task, we use maximum sequence length of 128, learning rate of $2e^{-5}$, dropout probability of 0.1 and we run fine-tuning for 3 epochs for all 3 models in Table 5.7. We report AUC, accuracy and F1 score for dev data. For test results, we submitted our predictions to Glue website for the official evaluation. We find that even in transfer learning, our SSE-SE pre-trained model still enjoys advantages over Google pre-trained model and our pre-trained model without SSE-SE. Our SSE-SE pre-trained model achieves 94.3% accuracy on SST-2 test set versus 93.6 and 93.8 respectively. If we are using SSE-SE for both pre-training and fine-tuning, we can achieve 94.5% accuracy on the SST-2 test set, which approaches the 94.9 score reported by the BERT-Large model. SSE probability of 0.01 is used for fine-tuning.

8.3.4 Proofs

Throughout this section, we will suppress the probability parameters, $p(., .|\Phi) = p(., .)$.

Proof: [Proof of Theorem 1] Consider the following variability term,

$$\sup_{\Theta} |S(\Theta) - S_n(\Theta)|. \tag{8.3}$$

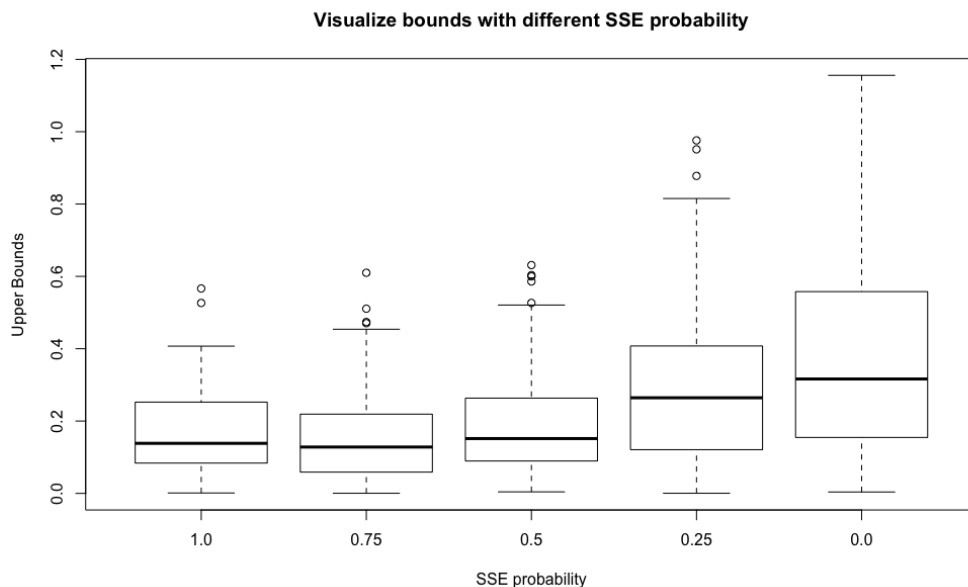


Figure 8.5. Simulation of a bound on $\rho_{L,n}$ for the movielens1M dataset. Throughout the simulation, L is replaced with ℓ (which will bound $\rho_{L,n}$ by Jensen’s inequality). The SSE probability parameter dictates the probability of transitioning. When this is 0 (box plot on the right), the distribution is that of the samples from the standard Rademacher complexity (without the sup and expectation). As we increase the transition probability, the values for $\rho_{L,n}$ get smaller.

Let us break the variability term into two components

$$\mathbb{E}_{X,Y} \sup_{\Theta} |S_n(\Theta) - \mathbb{E}_{Y|X}[S_n(\Theta)]| + \mathbb{E}_{X,Y} \sup_{\Theta} |\mathbb{E}_{Y|X}[S_n(\Theta)] - S(\Theta)|,$$

where X, Y represent the random input and label. To control the first term, we introduce a ghost dataset (x_i, y'_i) , where y'_i are independently and identically distributed according to $y_i|x_i$. Define

$$S'_n(\Theta) = \sum_i \sum_{\mathbf{k}} p(\mathbf{j}^i, \mathbf{k}) \ell(\mathbf{E}[\mathbf{k}], y'_i | \Theta) \quad (8.4)$$

be the empirical SSE risk with respect to this ghost dataset.

We will rewrite $\mathbb{E}_{Y|X}[S_n(\Theta)]$ in terms of the ghost dataset and apply Jensen’s inequality

and law of iterated conditional expectation:

$$\mathbb{E} \sup_{\Theta} |\mathbb{E}_{Y|X}[S_n(\Theta)] - S_n(\Theta)| \quad (8.5)$$

$$= \mathbb{E} \sup_{\Theta} |\mathbb{E}_{Y'|X}[S'_n(\Theta) - S_n(\Theta)]| \quad (8.6)$$

$$\leq \mathbb{E} \mathbb{E}_{Y'|X} \left[\sup_{\Theta} |S'_n(\Theta) - S_n(\Theta)| \right] \quad (8.7)$$

$$= \mathbb{E} \sup_{\Theta} |S'_n(\Theta) - S_n(\Theta)|. \quad (8.8)$$

Notice that

$$\begin{aligned} S'_n(\Theta) - S_n(\Theta) &= \sum_i \sum_{\mathbf{k}} p(\mathbf{j}^i, \mathbf{k}) (\ell(\mathbf{E}[\mathbf{k}], y'_i|\Theta) - \ell(\mathbf{E}[\mathbf{k}], y_i|\Theta)) \\ &= \sum_i \sum_{\mathbf{k}} p(\mathbf{j}^i, \mathbf{k}) (e(\mathbf{E}[\mathbf{k}], y'_i|\Theta) - e(\mathbf{E}[\mathbf{k}], y_i|\Theta)). \end{aligned}$$

Because $y_i, y'_i|X$ are independent the term $(\sum_{\mathbf{k}} p(\mathbf{j}^i, \mathbf{k}) (e(\mathbf{E}[\mathbf{k}], y'_i|\Theta) - e(\mathbf{E}[\mathbf{k}], y_i|\Theta)))_i$ is a vector of symmetric independent random variables. Thus its distribution is not effected by multiplication by arbitrary Rademacher vectors $\sigma_i \in \{-1, +1\}$.

$$\mathbb{E} \sup_{\Theta} |S'_n(\Theta) - S_n(\Theta)| = \mathbb{E} \sup_{\Theta} \left| \sum_i \sigma_i \sum_{\mathbf{k}} p(\mathbf{j}^i, \mathbf{k}) (e(\mathbf{E}[\mathbf{k}], y'_i|\Theta) - e(\mathbf{E}[\mathbf{k}], y_i|\Theta)) \right|.$$

But this is bounded by

$$2\mathbb{E} \mathbb{E}_{\sigma} \sup_{\Theta} \left| \sum_i \sigma_i \sum_{\mathbf{k}} p(\mathbf{j}^i, \mathbf{k}) e(\mathbf{E}[\mathbf{k}], y_i|\Theta) \right|.$$

For the second term,

$$\mathbb{E} \sup_{\Theta} |\mathbb{E}_{Y|X}[S_n(\Theta)] - S(\Theta)|$$

we will introduce a second ghost dataset x'_i, y'_i drawn iid to x_i, y_i . Because we are augmenting the input then this results in a new ghost encoding $\mathbf{o}j^i$. Let

$$S'_n(\Theta) = \sum_i \sum_{\mathbf{k}} p(\mathbf{j}^i, \mathbf{k}) \ell(\mathbf{E}[\mathbf{k}], y'_i|\Theta) \quad (8.9)$$

be the empirical risk with respect to this ghost dataset. Then we have that

$$S(\Theta) = \mathbb{E}_{X'} \mathbb{E}_{Y'|X'} S'_n(\Theta)$$

Thus,

$$\mathbb{E} \sup_{\Theta} |\mathbb{E}_{Y|X}[S_n(\Theta)] - S(\Theta)| \quad (8.10)$$

$$= \mathbb{E} \sup_{\Theta} |\mathbb{E}_{X'}[E_{Y|X}[S_n(\Theta)] - E_{Y'|X'}[S'_n(\Theta)]]| \quad (8.11)$$

$$\leq \mathbb{E} \mathbb{E}_{X'} \left[\sup_{\Theta} |E_{Y|X}[S_n(\Theta)] - E_{Y'|X'}[S'_n(\Theta)]| \right] \quad (8.12)$$

$$= \mathbb{E} \sup_{\Theta} |E_{Y|X}[S_n(\Theta)] - E_{Y'|X'}[S'_n(\Theta)]|. \quad (8.13)$$

Notice that we may write,

$$E_{Y|X}[S_n(\Theta)] - E_{Y'|X'}[S'_n(\Theta)] = \sum_i \sum_{\mathbf{k}} (p(\mathbf{j}^i, \mathbf{k}) - p(\mathbf{j}'^i, \mathbf{k})) L(\mathbf{E}[\mathbf{k}]|\Theta)$$

Again we may introduce a second set of Rademacher random variables σ'_i , which results in

$$\mathbb{E} \sup_{\Theta} |E_{Y|X}[S_n(\Theta)] - E_{Y'|X'}[S'_n(\Theta)]| \leq 2\mathbb{E} \mathbb{E}_{\sigma'} \sup_{\Theta} \left| \sum_i \sigma'_i \sum_{\mathbf{k}} p(\mathbf{j}^i, \mathbf{k}) L(\mathbf{E}[\mathbf{k}]|\Theta) \right|.$$

And this is bounded by

$$2\mathbb{E} \mathbb{E}_{\sigma'} \sup_{\Theta} \left| \sum_i \sigma'_i \sum_{\mathbf{k}} p(\mathbf{j}^i, \mathbf{k}) L(\mathbf{E}[\mathbf{k}]|\Theta) \right| \leq 2\mathbb{E} \sup_{\Theta} \left| \sum_i \sigma'_i \sum_{\mathbf{k}} p(\mathbf{j}^i, \mathbf{k}) \ell(\mathbf{E}[\mathbf{k}], y_i|\Theta) \right|$$

by Jensen's inequality again. \square

Proof: [Proof of Theorem 2] It is clear that $2\mathcal{B} \geq B(\hat{\Theta}) + B(\Theta^*)$. It remains to show our concentration inequality. Consider changing a single sample, (x_i, y_i) to (x'_i, y'_i) , thus resulting in the SSE empirical risk, $S_{n,i}(\Theta)$. Thus,

$$\begin{aligned} S_n(\Theta) - S_{n,i}(\Theta) &= \sum_{\mathbf{k}} p(\mathbf{j}^i, \mathbf{k}) \cdot \ell(\mathbf{E}[\mathbf{k}], y_i|\Theta) - \sum_{\mathbf{k}} p(\mathbf{j}'^i, \mathbf{k}) \cdot \ell(\mathbf{E}[\mathbf{k}], y'_i|\Theta) \\ &= \sum_{\mathbf{k}} p(\mathbf{j}^i, \mathbf{k}) \cdot (\ell(\mathbf{E}[\mathbf{k}], y_i|\Theta) - \ell(\mathbf{E}[\mathbf{k}], y'_i|\Theta)) + \sum_{\mathbf{k}} (p(\mathbf{j}'^i, \mathbf{k}) - p(\mathbf{j}^i, \mathbf{k})) \cdot \ell(\mathbf{E}[\mathbf{k}], y'_i|\Theta) \\ &\leq b \left(\sum_{\mathbf{k}} p(\mathbf{j}^i, \mathbf{k}) + \sum_{\mathbf{k}} p(\mathbf{j}'^i, \mathbf{k}) \right) \leq 2b. \end{aligned}$$

Then the result follows from McDiarmid's inequality. \square

8.4 Appendix to Chapter 6

- $\text{NDCG}@K$: defined as:

$$\text{NDCG}@K = \frac{1}{n} \sum_{i=1}^n \frac{\text{DCG}@K(i, \Pi_i)}{\text{DCG}@K(i, \Pi_i^*)}, \quad (8.14)$$

where i represents i -th user and

$$\text{DCG}@K(i, \Pi_i) = \sum_{l=1}^K \frac{2^{R_{i\Pi_{il}}} - 1}{\log_2(l+1)}. \quad (8.15)$$

In the DCG definition, Π_{il} represents the index of the l -th ranked item for user i in test data based on the learned score matrix X . R is the rating matrix and R_{ij} is the rating given to item j by user i . Π_i^* is the ordering provided by the ground truth rating.

- $\text{Recall}@K$: defined as a fraction of positive items retrieved by the top K recommendations the model makes:

$$\text{Recall}@K = \frac{\sum_{i=1}^n \mathbb{1}\{\exists 1 \leq l \leq K : R_{i\Pi_{il}} = 1\}}{n}, \quad (8.16)$$

here we already assume there is only a single positive item that user will engage next and the indicator function $\mathbb{1}\{\exists 1 \leq l \leq k : R_{i\Pi_{il}} = 1\}$ is defined to indicate whether the positive item falls into the top K position in our obtained ranked list using scores predicted in (6.4).

Layer Normalization Layer normalization [4] normalizes neurons within a layer. Previous studies [4] show it is more effective than batch normalization for training recurrent neural networks (RNNs). One alternative is the batch normalization [51] but we find it does not work as well as the layer normalization in practice even for a reasonable large batch size of 128. Therefore, our SSE-PT model adopts layer normalization.

Residual Connections Residual connections are firstly proposed in ResNet for image classification problems [39]. Recent research finds that residual connections can help training very deep neural networks even if they are not convolutional neural networks [110]. Using residual connections allows us to train very deep neural networks here. For example, the best performing model for Movielens10M dataset in Table 8.12 is the SSE-PT with 6 attention blocks, in which $1 + 6 * 3 + 1 = 20$ layers are trained end-to-end.

Weight Decay Weight decay [64], also known as l_2 regularization [46], is applied to all embeddings, including both user and item embeddings.

Dropout Dropout [106] is applied to the embedding layer E , self-attention layer and pointwise feed-forward layer by stochastically dropping some percentage of hidden units to prevent co-adaptation of neurons. Dropout has been shown to be an effective way of regularizing deep learning models.

In summary, layer normalization and dropout are used in all layers except prediction layer. Residual connections are used in both self-attention layer and pointwise feed-forward layer. SSE-SE is used in embedding layer and prediction layer.

Table 8.7. Description of Datasets Used in Evaluations.

DATASET	#USERS	#ITEMS	AVG SEQUENCE LEN	MAX SEQUENCE LEN
BEAUTY	52,024	57,289	7.6	291
GAMES	31,013	23,715	7.3	858
STEAM	334,730	13,047	11.0	1,229
ML-1M	6,040	3,416	163.5	2,275
ML-10M	69,878	65,133	141.1	7,357

- PopRec: ranking items according to their popularity.
- BPR: Bayesian personalized ranking for implicit feedback setting [91]. It is a low-rank matrix factorization model with a pairwise loss function. But it does not utilize the temporal information. Therefore, it serves as a strong baseline for non-temporal methods.
- FMC: Factorized Markov Chains: a first-order Markov Chain method, in which predictions are made only based on previously engaged item.
- PFMC: a personalized Markov chain model [92] that combines matrix factorization and first-order Markov Chain to take advantage of both users' latent long-term preferences as well as short-term item transitions.
- TransRec: a first-order sequential recommendation method [40] in which items are embedded into a transition space and users are modelled as translation vectors

Table 8.8. Comparing our SSE-PT, SSE-PT++ with SASRec on Movielens1M dataset. We use number of negatives $C = 100$, dropout probability of 0.2 and learning rate of $1e^{-3}$ for all experiments while varying others. p_u, p_i, p_o are SSE probabilities for user embedding, input item embedding and output item embedding respectively.

Model	Movielens1m		Dimensions		Number of Blocks	Sampling Probability	SSE-SE Parameters		
	NDCG@10	Recall@10	d_u	d_i	b	p_s	p_u	p_i	p_o
SASRec	0.5961	0.8195	-	50	2	-	-	-	-
SASRec	0.5941	0.8182	-	100	2	-	-	-	-
SASRec	0.5996	0.8272	-	100	6	-	-	-	-
SSE-PT	0.6101	0.8343	50	50	2	-	0.92	0.1	0
SSE-PT	0.6164	0.8336	50	50	2	-	0.92	0	0.1
SSE-PT	0.5832	0.8091	50	50	2	-	0	0.1	0.1
SSE-PT	0.6174	0.8351	50	50	2	-	0.92	0.1	0.1
SSE-PT	0.5949	0.8205	75	25	2	-	0.92	0.1	0.1
SSE-PT	0.6214	0.8359	25	75	2	-	0.92	0.1	0.1
SSE-PT	0.6281	0.8341	50	100	2	-	0.92	0.1	0.1
SSE-PT++	0.6292	0.8389	50	100	2	0.3	0.92	0.1	0.1

operating on item sequences.

SQL-Rank [116] and item-based recommendations [94] are omitted because the former is similar to BPR [91] except using the listwise loss function instead of the pairwise loss function and the latter has been shown inferior to TransRec [40].

8.4.0.1 Deep-learning baselines

- GRU4Rec: the first RNN-based method proposed for the session-based recommendation problem [43]. It utilizes the GRU structures [20] initially proposed for speech modelling.
- GRU4Rec⁺: follow-up work of GRU4Rec by the same authors: the model has a very similar architecture to GRU4Rec but has a more complicated loss function [42].
- Caser: a CNN-based method [108] which embeds a sequence of recent items in both time and latent spaces forming an ‘image’ before learning local features through horizontal and vertical convolutional filters. In [108], user embeddings are included in the prediction layer only. On the contrast, in our Personalized Transformer, user embeddings are also introduced in the lowest embedding layer so they can play an

Table 8.9. Comparing our SSE-PT with SASRec on Movielens10M dataset. Unlike Table 8.8, we use the number of negatives $C = 500$ instead of 100 as $C = 100$ is too easy for this dataset and it gets too difficult to tell the differences between different methods: Hit Ratio@10 approaches 1.

Model	Movielens1m		Dimensions		Number of Blocks	SSE-SE Parameters		
	NDCG@10	Hit Ratio@10	d_u	d_i	b	p_u	p_i	p_y
SASRec	0.7268	0.9429	-	50	2	-	-	-
SASRec	0.7413	0.9474	-	100	2	-	-	-
SSE-PT	0.7199	0.9331	50	100	2	PS	0.01	0.01
SSE-PT	0.7169	0.9296	50	100	2	0.0	0.01	0.01
SSE-PT	0.7398	0.9418	50	100	2	0.2	0.01	0.01
SSE-PT	0.7500	0.9500	50	100	2	0.4	0.01	0.01
SSE-PT	0.7484	0.9480	50	100	2	0.6	0.01	0.01
SSE-PT	0.7529	0.9485	50	100	2	0.8	0.01	0.01
SSE-PT	0.7503	0.9505	50	100	2	1.0	0.01	0.01

important role in self-attention mechanisms as well as in prediction stages.

- STAMP: a session-based recommendation algorithm [68] using attention mechanism. [68] only uses fully connected layers with one attention block that is not self-attentive.
- SASRec: a self-attentive sequential recommendation method [56] motivated by Transformer in NLP [110]. Unlike our method SSE-PT, SASRec does not incorporate user embedding and therefore is not a personalized method. SASRec paper [56] also does not utilize SSE [117] for further regularization: only dropout and weight decay are used.
- HGN: hierarchical gating networks method to solve the sequential recommendation problem [69], which incorporates the user embeddings and gating networks for better personalization than the SASRec model.

Table 8.10. Comparing Different SSE probability for user embeddings for SSE-PT on Movielens1M Dataset. Embedding hidden units of 50 for users and 100 for items, attention blocks of 2, SSE probability of 0.01 for item embeddings, dropout probability of 0.2 and max length of 200 are used.

USER-SIDE SSE-SE PROBABILITY	NDCG@10	RECALL@10
PARAMETER SHARING	0.6188	0.8294
1.0	0.6258	0.8346
0.9	0.6275	0.8321
0.8	0.6244	0.8359
0.6	0.6256	0.8341
0.4	0.6237	0.8369
0.2	0.6163	0.8281
0.0	0.5908	0.8048

Table 8.11. Comparing Different Sampling Probability, p_s , of SSE-PT++ on Movielens1M Dataset. Hyper-parameters the same as Table 8.10, except that the max length T allowed is set 100 instead of 200 to show effects of sampling sequences.

SAMPLING PROBABILITY	NDCG@10	RECALL@10
SASREC ($T = 100$)	0.5769	0.8045
SSE-PT ($T = 100$)	0.6142	0.8212
1.0	0.5697	0.7977
0.8	0.5735	0.7801
0.6	0.6062	0.8242
0.4	0.6113	0.8273
0.3	0.6186	0.8318
0.2	0.6193	0.8233
0.0	0.6142	0.8212

Table 8.12. Comparing Different Number of Blocks for SSE-PT while Keeping The Rest Fixed on Movielens1M and Movielens10M Datasets.

DATASETS	# OF BLOCKS	NDCG@10	RECALL@10
	SASREC (6 BLOCKS)	0.5984	0.8207
MOVIELENS1M	1	0.6162	0.8301
	2	0.6280	0.8365
	3	0.6293	0.8376
	4	0.6270	0.8401
	5	0.6308	0.8361
	6	0.6270	0.8397
	SASREC (6 BLOCKS)	0.7531	0.9490
MOVIELENS10M	1	0.7454	0.9478
	2	0.7512	0.9522
	3	0.7543	0.9491
	4	0.7608	0.9485
	5	0.7619	0.9524
	6	0.7683	0.9537

Table 8.13. Varying number of negatives C in evaluation on Movielens1M dataset. Other hyper-parameters are fixed for a fair comparison.

METRIC	NDCG@10	RECALL@10	C
UN-PERSONALIZED	0.3787	0.6119	500
PERSONALIZED	0.3846	0.6171	500
UN-PERSONALIZED	0.2791	0.4781	1000
PERSONALIZED	0.2860	0.4929	1000
UN-PERSONALIZED	0.1939	0.3515	2000
PERSONALIZED	0.1993	0.3667	2000

REFERENCES

- [1] Zeinab Abbassi and Vahab S Mirrokni. A recommender system based on local random walks and spectral methods. In Proceedings of the 9th WebKDD and 1st SNA-KDD 2007 workshop on Web mining and social network analysis, pages 102–108. ACM, 2007.
- [2] Shivani Agarwal. Ranking on graph data. In Proceedings of the 23rd international conference on Machine learning, pages 25–32. ACM, 2006.
- [3] Paulo Sérgio Almeida, Carlos Baquero, Nuno Preguiça, and David Hutchison. Scalable bloom filters. Information Processing Letters, 101(6):255–261, 2007.
- [4] Jimmy Lei Ba, Jamie Ryan Kiros, and Geoffrey E Hinton. Layer normalization. arXiv preprint arXiv:1607.06450, 2016.
- [5] Suhrid Balakrishnan and Sumit Chopra. Collaborative ranking. In Proceedings of the fifth ACM international conference on Web search and data mining, pages 143–152. ACM, 2012.
- [6] James Bennett, Stan Lanning, et al. The netflix prize. In In KDD Cup and Workshop in conjunction with ACM SIGKDD International Conference on Knowledge Discovery and Data Mining, volume 2007, page 35. New York, NY, USA., 2007.
- [7] James Bennett, Stan Lanning, and Netflix Netflix. The netflix prize. In In KDD Cup and Workshop in conjunction with ACM SIGKDD International Conference on Knowledge Discovery and Data Mining, 2007.
- [8] Rianne van den Berg, Thomas N Kipf, and Max Welling. Graph convolutional matrix completion. arXiv preprint arXiv:1706.02263, 2017.
- [9] Burton H Bloom. Space/time trade-offs in hash coding with allowable errors. Communications of the ACM, 13(7):422–426, 1970.
- [10] Dhruva Borthakur, Jonathan Gray, Joydeep Sen Sarma, Kannan Muthukkaruppan, Nicolas Spiegelberg, Hairong Kuang, Karthik Ranganathan, Dmytro Molkov, Aravind Menon, Samuel Rash, et al. Apache hadoop goes realtime at facebook. In Proceedings of the 2011 ACM SIGMOD International Conference on Management of data, pages 1071–1080. ACM, 2011.
- [11] John S Breese, David Heckerman, and Carl Kadie. Empirical analysis of predictive algorithms for collaborative filtering. In Proceedings of the Fourteenth conference on Uncertainty in artificial intelligence, pages 43–52. Morgan Kaufmann Publishers Inc., 1998.
- [12] Andrei Broder and Michael Mitzenmacher. Network applications of bloom filters: A survey. Internet mathematics, 1(4):485–509, 2004.

- [13] Deng Cai, Xiaofei He, Jiawei Han, and Thomas S Huang. Graph regularized nonnegative matrix factorization for data representation. IEEE transactions on pattern analysis and machine intelligence, 33(8):1548–1560, 2011.
- [14] 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. ACM, 2015.
- [15] Zhe Cao, Tao Qin, Tie-Yan Liu, Ming-Feng Tsai, and Hang Li. Learning to rank: from pairwise approach to listwise approach. In Proceedings of the 24th international conference on Machine learning, pages 129–136. ACM, 2007.
- [16] Fay Chang, Jeffrey Dean, Sanjay Ghemawat, Wilson C Hsieh, Deborah A Wallach, Mike Burrows, Tushar Chandra, Andrew Fikes, and Robert E Gruber. Bigtable: A distributed storage system for structured data. ACM Transactions on Computer Systems (TOCS), 26(2):4, 2008.
- [17] Olivier Chapelle and S Sathiya Keerthi. Efficient algorithms for ranking with svms. Information Retrieval, 13(3):201–215, 2010.
- [18] Jianfei Chen, Jun Zhu, and Le Song. Stochastic training of graph convolutional networks with variance reduction. In International Conference on Machine Learning, pages 941–949, 2018.
- [19] Kai-Yang Chiang, Cho-Jui Hsieh, and Inderjit S Dhillon. Matrix completion with noisy side information. In Advances in Neural Information Processing Systems, pages 3447–3455, 2015.
- [20] Junyoung Chung, Caglar Gulcehre, KyungHyun Cho, and Yoshua Bengio. Empirical evaluation of gated recurrent neural networks on sequence modeling. arXiv preprint arXiv:1412.3555, 2014.
- [21] Moustapha M Cisse, Nicolas Usunier, Thierry Artieres, and Patrick Gallinari. Robust bloom filters for large multilabel classification tasks. In Advances in Neural Information Processing Systems, pages 1851–1859, 2013.
- [22] Matthieu Courbariaux, Yoshua Bengio, and Jean-Pierre David. Binaryconnect: Training deep neural networks with binary weights during propagations. In Advances in neural information processing systems, pages 3123–3131, 2015.
- [23] Michaël Defferrard, Xavier Bresson, and Pierre Vandergheynst. Convolutional neural networks on graphs with fast localized spectral filtering. In Advances in Neural Information Processing Systems, pages 3844–3852, 2016.
- [24] Mukund Deshpande and George Karypis. Item-based top-n recommendation algorithms. ACM Transactions on Information Systems (TOIS), 22(1):143–177, 2004.

- [25] Jacob Devlin, Ming-Wei Chang, Kenton Lee, and Kristina Toutanova. Bert: Pre-training of deep bidirectional transformers for language understanding. arXiv preprint arXiv:1810.04805, 2018.
- [26] John Duchi, Elad Hazan, and Yoram Singer. Adaptive subgradient methods for online learning and stochastic optimization. Journal of Machine Learning Research, 12(Jul):2121–2159, 2011.
- [27] P. M. Fenwick. A new data structure for cumulative frequency tables. Software: Practice and Experience, 1994.
- [28] Yoav Freund, Raj Iyer, Robert E Schapire, and Yoram Singer. An efficient boosting algorithm for combining preferences. Journal of machine learning research, 4(Nov):933–969, 2003.
- [29] Rainer Gemulla, Erik Nijkamp, Peter J Haas, and Yannis Sismanis. Large-scale matrix factorization with distributed stochastic gradient descent. In Proceedings of the 17th ACM SIGKDD international conference on Knowledge discovery and data mining, pages 69–77. ACM, 2011.
- [30] Ian Goodfellow, Yoshua Bengio, Aaron Courville, and Yoshua Bengio. Deep learning, volume 1. MIT press Cambridge, 2016.
- [31] Marco Gori, Augusto Pucci, V Roma, and I Siena. Itemrank: A random-walk based scoring algorithm for recommender engines. In International Joint Conferences on Artificial Intelligence, volume 7, pages 2766–2771, 2007.
- [32] 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. ACM, 2016.
- [33] Suriya Gunasekar, Oluwasanmi O Koyejo, and Joydeep Ghosh. Preference completion from partial rankings. In Advances in Neural Information Processing Systems, pages 1370–1378, 2016.
- [34] Severin Hacker and Luis Von Ahn. Matchin: eliciting user preferences with an online game. In Proceedings of the SIGCHI Conference on Human Factors in Computing Systems, pages 1207–1216. ACM, 2009.
- [35] Will Hamilton, Zhitao Ying, and Jure Leskovec. Inductive representation learning on large graphs. In Advances in Neural Information Processing Systems, pages 1024–1034, 2017.
- [36] William L Hamilton, Rex Ying, and Jure Leskovec. Representation learning on graphs: Methods and applications. arXiv preprint arXiv:1709.05584, 2017.
- [37] Song Han, Huizi Mao, and William J Dally. Deep compression: Compressing deep neural networks with pruning, trained quantization and huffman coding. arXiv preprint arXiv:1510.00149, 2015.

- [38] F Maxwell Harper and Joseph A Konstan. The movielens datasets: History and context. Acm transactions on interactive intelligent systems (tiis), 5(4):19, 2016.
- [39] Kaiming He, Xiangyu Zhang, Shaoqing Ren, and Jian Sun. Deep residual learning for image recognition. In Proceedings of the IEEE conference on computer vision and pattern recognition, pages 770–778, 2016.
- [40] Ruining He, Wang-Cheng Kang, and Julian McAuley. Translation-based recommendation. In Proceedings of the Eleventh ACM Conference on Recommender Systems, pages 161–169. ACM, 2017.
- [41] Xiangnan He, Lizi Liao, Hanwang Zhang, Liqiang Nie, Xia Hu, and Tat-Seng Chua. Neural collaborative filtering. In Proceedings of the 26th International Conference on World Wide Web, pages 173–182. International World Wide Web Conferences Steering Committee, 2017.
- [42] Balázs Hidasi and Alexandros Karatzoglou. Recurrent neural networks with top-k gains for session-based recommendations. In Proceedings of the 27th ACM International Conference on Information and Knowledge Management, pages 843–852. ACM, 2018.
- [43] Balázs Hidasi, Alexandros Karatzoglou, Linas Baltrunas, and Domonkos Tikk. Session-based recommendations with recurrent neural networks. arXiv preprint arXiv:1511.06939, 2015.
- [44] Will Hill, Larry Stead, Mark Rosenstein, and George Furnas. Recommending and evaluating choices in a virtual community of use. In Proceedings of the SIGCHI conference on Human factors in computing systems, pages 194–201. ACM Press/Addison-Wesley Publishing Co., 1995.
- [45] Geoffrey E Hinton, Nitish Srivastava, Alex Krizhevsky, Ilya Sutskever, and Ruslan R Salakhutdinov. Improving neural networks by preventing co-adaptation of feature detectors. arXiv preprint arXiv:1207.0580, 2012.
- [46] Arthur E Hoerl and Robert W Kennard. Ridge regression: Biased estimation for nonorthogonal problems. Technometrics, 12(1):55–67, 1970.
- [47] Jeremy Howard and Sebastian Ruder. Universal language model fine-tuning for text classification. arXiv preprint arXiv:1801.06146, 2018.
- [48] Cho-Jui Hsieh, Nagarajan Natarajan, and Inderjit Dhillon. Pu learning for matrix completion. In International Conference on Machine Learning, pages 2445–2453, 2015.
- [49] Yifan Hu, Yehuda Koren, and Chris Volinsky. Collaborative filtering for implicit feedback datasets. In Data Mining, 2008. ICDM’08. Eighth IEEE International Conference on, pages 263–272. Ieee, 2008.

- [50] Shanshan Huang, Shuaiqiang Wang, Tie-Yan Liu, Jun Ma, Zhumin Chen, and Jari Veijalainen. Listwise collaborative filtering. In Proceedings of the 38th International ACM SIGIR Conference on Research and Development in Information Retrieval, pages 343–352. ACM, 2015.
- [51] Sergey Ioffe and Christian Szegedy. Batch normalization: Accelerating deep network training by reducing internal covariate shift. arXiv preprint arXiv:1502.03167, 2015.
- [52] Mohsen Jamali and Martin Ester. Trustwalker: a random walk model for combining trust-based and item-based recommendation. In Proceedings of the 15th ACM SIGKDD international conference on Knowledge discovery and data mining, pages 397–406. ACM, 2009.
- [53] T. Joachims. Optimizing search engines using clickthrough data. In ACM SIGKDD, 2002.
- [54] Thorsten Joachims. Optimizing search engines using clickthrough data. In Proceedings of the eighth ACM SIGKDD international conference on Knowledge discovery and data mining, pages 133–142. ACM, 2002.
- [55] Santosh Kabbur, Xia Ning, and George Karypis. Fism: factored item similarity models for top-n recommender systems. In Proceedings of the 19th ACM SIGKDD international conference on Knowledge discovery and data mining, pages 659–667. ACM, 2013.
- [56] Wang-Cheng Kang and Julian McAuley. Self-attentive sequential recommendation. arXiv preprint arXiv:1808.09781, 2018.
- [57] Diederik P Kingma and Jimmy Ba. Adam: A method for stochastic optimization. arXiv preprint arXiv:1412.6980, 2014.
- [58] Thomas N Kipf and Max Welling. Semi-supervised classification with graph convolutional networks. arXiv preprint arXiv:1609.02907, 2016.
- [59] Yehuda Koren. Factorization meets the neighborhood: a multifaceted collaborative filtering model. In Proceedings of the 14th ACM SIGKDD international conference on Knowledge discovery and data mining, pages 426–434. ACM, 2008.
- [60] Yehuda Koren. Collaborative filtering with temporal dynamics. In Proceedings of the 15th ACM SIGKDD international conference on Knowledge discovery and data mining, pages 447–456. ACM, 2009.
- [61] Yehuda Koren, Robert Bell, and Chris Volinsky. Matrix factorization techniques for recommender systems. Computer, (8):30–37, 2009.
- [62] Oluwasanmi Koyejo, Sreangsu Acharyya, and Joydeep Ghosh. Retargeted matrix factorization for collaborative filtering. In Proceedings of the 7th ACM conference on Recommender systems, pages 49–56. ACM, 2013.

- [63] Alex Krizhevsky, Ilya Sutskever, and Geoffrey E Hinton. Imagenet classification with deep convolutional neural networks. In Advances in neural information processing systems, pages 1097–1105, 2012.
- [64] Anders Krogh and John A Hertz. A simple weight decay can improve generalization. In Advances in neural information processing systems, pages 950–957, 1992.
- [65] Yanyan Lan, Tie-Yan Liu, Zhiming Ma, and Hang Li. Generalization analysis of listwise learning-to-rank algorithms. In Proceedings of the 26th Annual International Conference on Machine Learning, pages 577–584. ACM, 2009.
- [66] Jens Lehmann, Robert Isele, Max Jakob, Anja Jentzsch, Dimitris Kontokostas, Pablo N Mendes, Sebastian Hellmann, Mohamed Morse, Patrick Van Kleef, Sören Auer, et al. Dbpedia—a large-scale, multilingual knowledge base extracted from wikipedia. Semantic Web, 6(2):167–195, 2015.
- [67] Dawen Liang, Jaan Altosaar, Laurent Charlin, and David M Blei. Factorization meets the item embedding: Regularizing matrix factorization with item co-occurrence. In Proceedings of the 10th ACM conference on recommender systems, pages 59–66. ACM, 2016.
- [68] Qiao Liu, Yifu Zeng, Refuoe Mokhosi, and Haibin Zhang. Stamp: short-term attention/memory priority model for session-based recommendation. In Proceedings of the 24th ACM SIGKDD International Conference on Knowledge Discovery & Data Mining, pages 1831–1839. ACM, 2018.
- [69] Chen Ma, Peng Kang, and Xue Liu. Hierarchical gating networks for sequential recommendation. arXiv preprint arXiv:1906.09217, 2019.
- [70] Hao Ma, Dengyong Zhou, Chao Liu, Michael R Lyu, and Irwin King. Recommender systems with social regularization. In Proceedings of the fourth ACM international conference on Web search and data mining, pages 287–296. ACM, 2011.
- [71] 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.
- [72] George A Miller. Wordnet: a lexical database for english. Communications of the ACM, 38(11):39–41, 1995.
- [73] Andriy Mnih and Ruslan R Salakhutdinov. Probabilistic matrix factorization. In Advances in neural information processing systems, pages 1257–1264, 2008.
- [74] Federico Monti, Michael Bronstein, and Xavier Bresson. Geometric matrix completion with recurrent multi-graph neural networks. In Advances in Neural Information Processing Systems, pages 3697–3707, 2017.

- [75] Steven J Nowlan and Geoffrey E Hinton. Simplifying neural networks by soft weight-sharing. Neural computation, 4(4):473–493, 1992.
- [76] Lawrence Page, Sergey Brin, Rajeev Motwani, and Terry Winograd. The pagerank citation ranking: Bringing order to the web. Technical report, Stanford InfoLab, 1999.
- [77] Tapio Pahikkala, Evgeni Tsivtsivadze, Antti Airola, Jouni Järvinen, and Jorma Boberg. An efficient algorithm for learning to rank from preference graphs. Machine Learning, 75(1):129–165, 2009.
- [78] Rong Pan, Yunhong Zhou, Bin Cao, Nathan N Liu, Rajan Lukose, Martin Scholz, and Qiang Yang. One-class collaborative filtering. In Data Mining, 2008. ICDM’08. Eighth IEEE International Conference on, pages 502–511. IEEE, 2008.
- [79] Weike Pan and Li Chen. Gbpr: Group preference based bayesian personalized ranking for one-class collaborative filtering. In International Joint Conferences on Artificial Intelligence, volume 13, pages 2691–2697, 2013.
- [80] Weike Pan, Qiang Yang, Yuchao Duan, Ben Tan, and Zhong Ming. Transfer learning for behavior ranking. ACM Transactions on Intelligent Systems and Technology (TIST), 8(5):65, 2017.
- [81] Dohyung Park, Joe Neeman, Jin Zhang, Sujay Sanghavi, and Inderjit Dhillon. Preference completion: Large-scale collaborative ranking from pairwise comparisons. In International Conference on Machine Learning, pages 1907–1916, 2015.
- [82] Razvan Pascanu, Tomas Mikolov, and Yoshua Bengio. On the difficulty of training recurrent neural networks. In International Conference on Machine Learning, pages 1310–1318, 2013.
- [83] Jeffrey Pennington, Richard Socher, and Christopher 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.
- [84] 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. ACM, 2014.
- [85] Matt Post. A call for clarity in reporting BLEU scores. In Proceedings of the Third Conference on Machine Translation: Research Papers, pages 186–191, Belgium, Brussels, October 2018. Association for Computational Linguistics.
- [86] Manuel Pozo, Raja Chiky, Farid Meziane, and Elisabeth Métais. An item/user representation for recommender systems based on bloom filters. In 2016 IEEE Tenth International Conference on Research Challenges in Information Science (RCIS), pages 1–12. IEEE, 2016.

- [87] Alec Radford, Karthik Narasimhan, Tim Salimans, and Ilya Sutskever. Improving language understanding by generative pre-training. [URL https://s3-us-west-2.amazonaws.com/openai-assets/research-covers/languageunsupervised/language-understanding-paper.pdf](https://s3-us-west-2.amazonaws.com/openai-assets/research-covers/languageunsupervised/language-understanding-paper.pdf), 2018.
- [88] Alec Radford, Jeffrey Wu, Rewon Child, David Luan, Dario Amodei, and Ilya Sutskever. Language models are unsupervised multitask learners. [URL https://openai.com/blog/better-language-models](https://openai.com/blog/better-language-models), 2019.
- [89] 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.
- [90] Steffen Rendle. Factorization machines. In *Data Mining (ICDM), 2010 IEEE 10th International Conference on*, pages 995–1000. IEEE, 2010.
- [91] Steffen Rendle, Christoph Freudenthaler, Zeno Gantner, and Lars Schmidt-Thieme. Bpr: Bayesian personalized ranking from implicit feedback. In *Proceedings of the twenty-fifth conference on uncertainty in artificial intelligence*, pages 452–461. AUAI Press, 2009.
- [92] Steffen Rendle, Christoph Freudenthaler, and Lars Schmidt-Thieme. Factorizing personalized markov chains for next-basket recommendation. In *Proceedings of the 19th international conference on World wide web*, pages 811–820. ACM, 2010.
- [93] Ruslan Salakhutdinov, Andriy Mnih, and Geoffrey Hinton. Restricted boltzmann machines for collaborative filtering. In *Proceedings of the 24th international conference on Machine learning*, pages 791–798. ACM, 2007.
- [94] Badrul Munir Sarwar, George Karypis, Joseph A Konstan, John Riedl, et al. Item-based collaborative filtering recommendation algorithms. *World Wide Web Conference*, 1:285–295, 2001.
- [95] J Ben Schafer, Dan Frankowski, Jon Herlocker, and Shilad Sen. Collaborative filtering recommender systems. In *The adaptive web*, pages 291–324. Springer, 2007.
- [96] Joan Serrà and Alexandros Karatzoglou. Getting deep recommenders fit: Bloom embeddings for sparse binary input/output networks. In *Proceedings of the Eleventh ACM Conference on Recommender Systems*, pages 279–287. ACM, 2017.
- [97] Devavrat Shah et al. Gossip algorithms. *Foundations and Trends® in Networking*, 3(1):1–125, 2009.
- [98] Guy Shani, Max Chickering, and Christopher Meek. Mining recommendations from the web. In *Proceedings of the 2008 ACM conference on Recommender systems*, pages 35–42. ACM, 2008.

- [99] Qinfeng Shi, James Petterson, Gideon Dror, John Langford, Alex Smola, and SVN Vishwanathan. Hash kernels for structured data. Journal of Machine Learning Research, 10(Nov):2615–2637, 2009.
- [100] Yue Shi, Martha Larson, and Alan Hanjalic. List-wise learning to rank with matrix factorization for collaborative filtering. In Proceedings of the fourth ACM conference on Recommender systems, pages 269–272. ACM, 2010.
- [101] Anita Shinde and Ila Savant. User based collaborative filtering using bloom filter with mapreduce. In Proceedings of International Conference on ICT for Sustainable Development, pages 115–123. Springer, 2016.
- [102] Si Si, Kai-Yang Chiang, Cho-Jui Hsieh, Nikhil Rao, and Inderjit S Dhillon. Goal-directed inductive matrix completion. In ACM SIGKDD International Conference on Knowledge Discovery and Data Mining, 2016.
- [103] Ajit P Singh and Geoffrey J Gordon. Relational learning via collective matrix factorization. In Proceedings of the 14th ACM SIGKDD international conference on Knowledge discovery and data mining, pages 650–658. ACM, 2008.
- [104] Nathan Srebro, Jason Rennie, and Tommi S Jaakkola. Maximum-margin matrix factorization. In Advances in neural information processing systems, pages 1329–1336, 2005.
- [105] Nathan Srebro, Jason DM Rennie, and Tommi S Jaakkola. Maximum-margin matrix factorization. In Advances in Neural Information Processing Systems, volume 17, pages 1329–1336, 2004.
- [106] Nitish Srivastava, Geoffrey Hinton, Alex Krizhevsky, Ilya Sutskever, and Ruslan Salakhutdinov. Dropout: a simple way to prevent neural networks from overfitting. The Journal of Machine Learning Research, 15(1):1929–1958, 2014.
- [107] Christian Szegedy, Vincent Vanhoucke, Sergey Ioffe, Jon Shlens, and Zbigniew Wojna. Rethinking the inception architecture for computer vision. In Proceedings of the IEEE conference on computer vision and pattern recognition, pages 2818–2826, 2016.
- [108] Jiayi Tang and Ke Wang. Personalized top-n sequential recommendation via convolutional sequence embedding. In Proceedings of the Eleventh ACM International Conference on Web Search and Data Mining, pages 565–573. ACM, 2018.
- [109] Robert Tibshirani. Regression shrinkage and selection via the lasso. Journal of the Royal Statistical Society. Series B (Methodological), pages 267–288, 1996.
- [110] Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N Gomez, Łukasz Kaiser, and Illia Polosukhin. Attention is all you need. In Advances in Neural Information Processing Systems, pages 5998–6008, 2017.

- [111] Maksims Volkovs and Richard S Zemel. Collaborative ranking with 17 parameters. In Advances in Neural Information Processing Systems, pages 2294–2302, 2012.
- [112] Jun Wang, Lantao Yu, Weinan Zhang, Yu Gong, Yinghui Xu, Benyou Wang, Peng Zhang, and Dell Zhang. Irgan: A minimax game for unifying generative and discriminative information retrieval models. In Proceedings of the 40th International ACM SIGIR conference on Research and Development in Information Retrieval, pages 515–524. ACM, 2017.
- [113] Markus Weimer, Alexandros Karatzoglou, Quoc V Le, and Alex J Smola. Cofi rank-maximum margin matrix factorization for collaborative ranking. In Advances in neural information processing systems, pages 1593–1600, 2008.
- [114] Markus Weimer, Alexandros Karatzoglou, Quoc Viet Le, and Alex Smola. Maximum margin matrix factorization for collaborative ranking. Advances in neural information processing systems, pages 1–8, 2007.
- [115] Liwei Wu, Cho-Jui Hsieh, and James Sharpnack. Large-scale collaborative ranking in near-linear time. In Proceedings of the 23rd ACM SIGKDD International Conference on Knowledge Discovery and Data Mining, pages 515–524. ACM, 2017.
- [116] Liwei Wu, Cho-Jui Hsieh, and James Sharpnack. Sql-rank: A listwise approach to collaborative ranking. In Proceedings of Machine Learning Research (35th International Conference on Machine Learning), volume 80, 2018.
- [117] Liwei Wu, Shuqing Li, Cho-Jui Hsieh, and James Sharpnack. Stochastic shared embeddings: Data-driven regularization of embedding layers. arXiv preprint arXiv:1905.10630, 2019.
- [118] Liwei Wu, Shuqing Li, Cho-Jui Hsieh, and James Sharpnack. Temporal collaborative ranking via personalized transformer. arXiv preprint arXiv:1908.05435, 2019.
- [119] Liwei Wu, Hsiang-Fu Yu, Nikhil Rao, James Sharpnack, and Cho-Jui Hsieh. Graph dna: Deep neighborhood aware graph encoding for collaborative filtering. arXiv preprint arXiv:1905.12217, 2019.
- [120] Yao Wu, Christopher DuBois, Alice X Zheng, and Martin Ester. Collaborative denoising auto-encoders for top-n recommender systems. In Proceedings of the Ninth ACM International Conference on Web Search and Data Mining, pages 153–162. ACM, 2016.
- [121] Fen Xia, Tie-Yan Liu, Jue Wang, Wensheng Zhang, and Hang Li. Listwise approach to learning to rank: theory and algorithm. In Proceedings of the 25th international conference on Machine learning, pages 1192–1199. ACM, 2008.
- [122] Wenlei Xie, David Bindel, Alan Demers, and Johannes Gehrke. Edge-weighted personalized pagerank: breaking a decade-old performance barrier. In Proceedings

of the 21th ACM SIGKDD International Conference on Knowledge Discovery and Data Mining, pages 1325–1334. ACM, 2015.

- [123] Rex Ying, Ruining He, Kaifeng Chen, Pong Eksombatchai, William L Hamilton, and Jure Leskovec. Graph convolutional neural networks for web-scale recommender systems. In Proceedings of the 24th ACM SIGKDD International Conference on Knowledge Discovery & Data Mining, pages 974–983. ACM, 2018.
- [124] Hsiang-Fu Yu, Mikhail Bilenko, and Chih-Jen Lin. Selection of negative samples for one-class matrix factorization. In Proceedings of the 2017 SIAM International Conference on Data Mining, pages 363–371. SIAM, 2017.
- [125] Hsiang-Fu Yu, Cho-Jui Hsieh, Hyokun Yun, SVN Vishwanathan, and Inderjit S Dhillon. A scalable asynchronous distributed algorithm for topic modeling. In Proceedings of the 24th International Conference on World Wide Web, pages 1340–1350. ACM, 2015.
- [126] Hsiang-Fu Yu, Hsin-Yuan Huang, Inderjit S Dhillon, and Chih-Jen Lin. A unified algorithm for one-class structured matrix factorization with side information. In AAAI, pages 2845–2851, 2017.
- [127] R. Zafarani and H. Liu. Social computing data repository at ASU, 2009.
- [128] Tinghui Zhou, Hanhuai Shan, Arindam Banerjee, and Guillermo Sapiro. Kernelized probabilistic matrix factorization: Exploiting graphs and side information. In Proceedings of the 2012 SIAM international Conference on Data mining, pages 403–414. SIAM, 2012.

Liwei Wu
March 2020
PhD in Statistics

Advances in Collaborative Filtering and Ranking

Abstract

In this dissertation, we cover some recent advances in collaborative filtering and ranking. In chapter 1, we give a brief introduction of the history and the current landscape of collaborative filtering and ranking; chapter 2 we first talk about pointwise collaborative filtering problem with graph information, and how our proposed new method can encode very deep graph information which helps four existing graph collaborative filtering algorithms; chapter 3 is on the pairwise approach for collaborative ranking and how we speed up the algorithm to near-linear time complexity; chapter 4 is on the new listwise approach for collaborative ranking and how the listwise approach is a better choice of loss for both explicit and implicit feedback over pointwise and pairwise loss; chapter 5 is about the new regularization technique Stochastic Shared Embeddings (SSE) we proposed for embedding layers and how it is both theoretically sound and empirically effectively for 6 different tasks across recommendation and natural language processing; chapter 6 is how we introduce personalization for the state-of-the-art sequential recommendation model with the help of SSE, which plays an important role in preventing our personalized model from overfitting to the training data; chapter 7, we summarize what we have achieved so far and predict what the future directions can be; chapter 8 is the appendix to all the chapters.