
Cross-Dataset Recommender System for Overcoming Selection Bias

Anonymous Author(s)

Affiliation

Address

email

Abstract

The abstract paragraph should be indented 1/2 inch (3 picas) on both the left- and right-hand margins. Use 10 point type, with a vertical spacing (leading) of 11 points. The word **Abstract** must be centered, bold, and in point size 12. Two line spaces precede the abstract. The abstract must be limited to one paragraph. The abstract paragraph should be indented 1/2 inch (3 picas) on both the left- and right-hand margins. Use 10 point type, with a vertical spacing (leading) of 11 points. The word **Abstract** must be centered, bold, and in point size 12. Two line spaces precede the abstract. The abstract must be limited to one paragraph.

1 Introduction

Selection bias is one of the most prevalent sources of biases for recommender systems [2]. Selection bias happens when there is a pattern in the users' ratings that is unique to the training set. For example, in a recommender system for movies, users might mainly rate movies that are recommended to them, which is a small section of movies already tailored to the users' tastes [7]. However, the environment the recommender system is deployed on contains all movies regardless of personal tastes. This discrepancy produces a misalignment between training and deployed settings, which is known as a distribution shift. Tackling the selection bias in a recommender dataset has been a constant challenge in designing recommender algorithms[8].

When recommender systems are deployed in real-world platforms, it is arguably likely that the clients are able to collect rating-associated data from different source distributions. These data (or feedbacks) are either implicit or explicit [1], while implicit data often comes in a highly-quantized, binary form and explicit data is usually less quantized. For example, a user's explicit rating of a movie is usually quantized to a number between 1 and 5, while the implicit feedback, such as if a user spontaneous searches for a movie, is often binary (more quantized.) We extend our hypothesis and argue that implicit data contains less or no selection bias compared to explicit data. This is because implicit data are often from users' spontaneous actions while explicit data are prejudiced toward the output of the recommender system in the previous feedback loop and also users' innate biases.

In this paper, we attempt to take advantage of datasets from differently quantized sources. More specifically, we proposed a way to feed both a 5-quantized dataset and a binary dataset to any gradient-based recommender algorithm. To ensure both datasets do not significantly lose their values in the presence of selection bias, we first examine the susceptibility to selection bias of differently quantized datasets from a single distribution. We design experiments under a simulated environment and shows that susceptibility to selection bias is not correlated with the way a dataset is quantized.

Then, since a less-quantized dataset inherently contains more information than a more-quantized dataset [11] and thus is more suitable for training, we decided to use it as the training data of

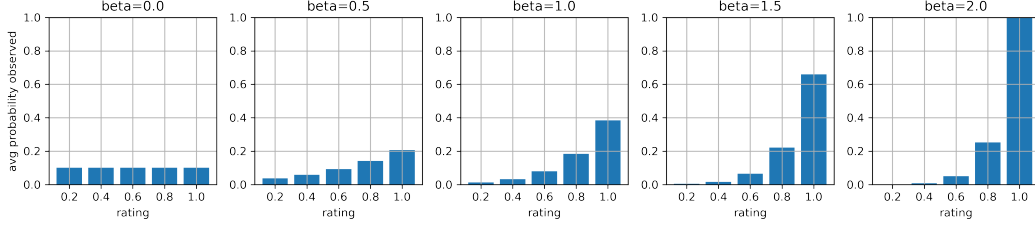


Figure 1: Visualizing the effect of controlling bias

matrix factorization and use the more-quantized dataset for propensity scoring and deriving the inverse-probability-scoring (IPS) estimator [10] [3], a causal inference approach applicable to matrix completion-based recommender algorithms. In this way, our cross-dataset learning framework empowers existing recommender algorithms to make use of the more quantized, less biased data. We carried out experiment and found our method outperforming baselines by a significant margin.

2 Related Work

Prior works on overcoming selection bias-induced distribution shift via a propensity-based approach begins with the seminal paper [8], which introduces the IPS method into recommender systems. Follow-up works aim at providing a learning-based or behavioral model of user feedbacks for propensity estimation [4] [12], which remains the central concern of this approach. While many works rely on a propensity matrix that is assumed to be known, TODO our work provides an idea of using a more quantized, implicit dataset for accurate propensity estimation.

The IPS method in our context is, essentially, a method of weighting training examples to correct the bias in the training data. Equivalent approaches such as importance weighting are widely used for domain adaptation in fields other than recommender systems [9] [13]. Discussions on the IPS-based domain adaptation for countering selection bias, which is most relevant to recommender systems, remain limited.

3 Susceptibility to Selection Bias

We first examine the susceptibility to selection bias of differently quantized data by manually introducing biased distributions of various degrees to the differently quantized training sets. It is crucial for us that the 2-quantized datasets do not exhibit particular weakness when facing selection bias so that they can be properly adopted for propensity estimation.

3.1 Simulated Environment for Controlling Bias

Since selection bias is uncontrollable in a dataset completely drawn from real-world, we have to adopt a simulated environment [6] with both semi-synthetic and synthetic datasets, which shall be explained in section 5. In our environment, we propose the **softmax observation model** and introduce a hyperparameter β to control the degree of bias. For a rating matrix R , the corresponding probability matrix of each rating being observed is $\Pr(R_{u,i} \text{ is observed}) = k \text{softmax}(\beta R_{u,i})$, where k is set so that the expected proportion of observed ratings is controlled. The effect of β on probability of being observed for different ratings is visualized in figure 1. In our experiment we assume constantly 10% of ratings are observed.

3.2 Results

Figure 2 shows our results on two dataset (see section 5) and 3 classic algorithms: user-KNN, item-KNN, and SVD matrix factorization. Although the RMSE grows consistently as β increases for all datasets and algorithms, differently quantized datasets do not exhibit significantly different growth rates. We thus conclude it is a viable approach to use the more quantized, less biased dataset for propensity estimation.

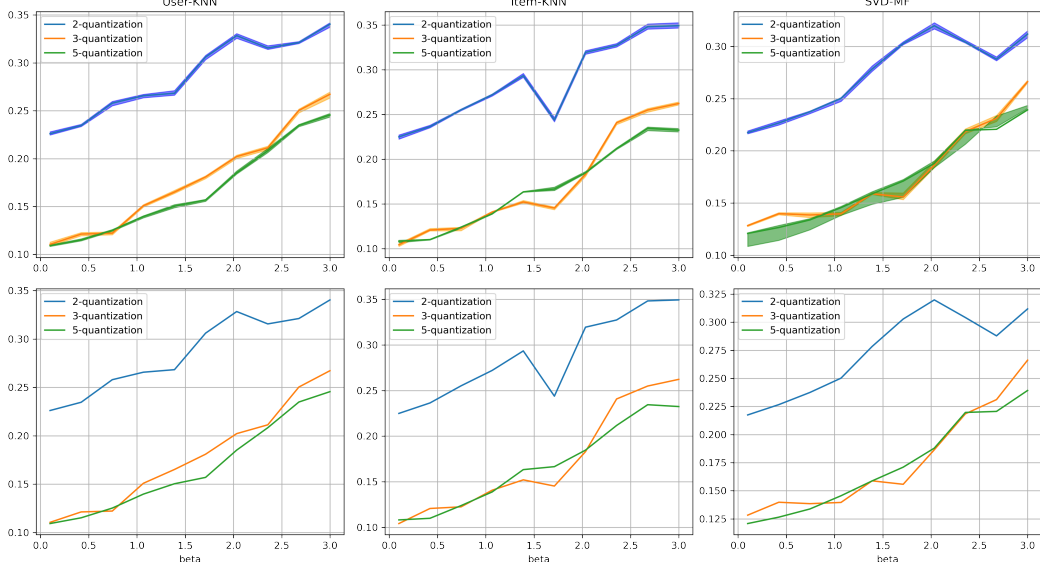


Figure 2: The effect of selection bias is not affected by quantization

4 Cross-Dataset Propensity Estimation

With the safety guarantee introduced above, we now turn to our cross-dataset matrix factorization model. Matrix factorization is a simple recommender model that decomposes the rating matrix based on known entries and then predicts the unknown entries[5]. We integrate propensities as [8] did and formulated the recommendation problem as the empirical risk minimization framework below.

$$\arg \min_{V, W, A} \frac{1}{N} \left(\frac{(Y_{u,i} - (V_u^T W_i + A))^2}{P_{u,i}} \right) + c \|A\|^2 \quad (1)$$

where $A = \{b_u, b_i, \mu\}$ represents the standard bias parameters (offset), V and W are the decomposed vectors, $\hat{Y} = V_u^T W_i + A$ is the predicted rating, N is the number of ratings, and $c \|A\|^2$ is the regularizer. The inverse propensity scores $1/P_{u,i}$ are multiplied to each rating during learning, which analogous to re-weighting ratings based on their biases.

Denote the biased training set as D . we propose the **naive-bayes propensity estimator** (NBPE-MF) from a more quantized (binary) dataset D' . Essentially,

$$P_{u,i} = \Pr(Y_{u,i} \text{ is observed} \mid Y_{u,i} = r_{u,i}) = \frac{\Pr(Y_{u,i} = r_{u,i} \mid Y_{u,i} \text{ is observed})}{\Pr(Y_{u,i} = r_{u,i})}, \quad (2)$$

where the numerator can be easily approximated from the dataset, but the denominator requires additional data from less biased sources. In this paper, we estimates the denominator from the unbiased, more quantized dataset using a categorical naive-bayes model. We use D' as a mask to hide biased ratings in the training set from propensity estimation. In other words, $Y_{u,i}$ is considered observed only when $D'_{u,i}$ is observed.

$$NB_{equation here}. \quad (3)$$

If the corresponding rating is not captured in D' , we set its propensity to the average propensity.

5 Experiment

We designed experiments to verify the performance of our cross-dataset model. We trained on two datasets with $\beta = 1$ in training set and used two baseline algorithms for comparison. We selected the root mean square error (RMSE) and the mean absolute error (MAE) as the performance metrics.

Table 1: Test set RMSE and MAE for NBPE-MF and baselines

	ML100K		Latent Factors	
	RMSE	MAE	RMSE	MAE
MF	0.1046	0.0833	0.1331	0.1045
NPE-MF	0.1048	0.0834	0.132	0.1049
MD-MF	0.1044	0.083	0.1299	0.1027
NBPE-MF	0.076	0.062	0.1065	0.0841

93 5.1 Datasets

94 **Imputed ML100K Dataset.** The ML100K dataset provides 100 thousand MNAR ratings across
 95 1683 movies rated by 944 users and is the standard large-scale dataset used for recommender systems.
 96 Since we need the ground truth ratings for controlling bias, we impute the missing ratings using
 97 standard matrix factorization. Ratings are normalized between 0 and 1 for consistency with the other
 98 dataset.

99 **Latent Factors Simulated Dataset.** The latent factors dataset is a synthetic dataset that models
 100 real-world user behavior. Users and item both have random latent vectors to simulate preferences;
 101 both also have biases. The environment is provided by [6].

102 5.2 Baselines

103 **Matrix Factorization (MF.)** As a simple baseline, we adopt the standard matrix factorization that
 104 does not adopt propensity estimation or importance weighting.

105 **Naive Propensity Estimator (NPE-MF.)** The naive propensity estimator naively estimates the
 106 propensity scores from the already biased training data and plugs in the results to equation 1.

107 **Mixing Datasets Matrix Factorization (MD-MF.)** To analyze whether our cross-dataset model
 108 more efficiently exploits all existing data, we create a larger dataset by mixing the two differently
 109 quantized datasets and train the MF algorithm on it.

110 5.3 Results

111 Table 1 shows the the experiment’s results, each entry as the average of five independent trials.
 112 Averaged across datasets and metrics, our approach gains 23.1% reduction in error compared to
 113 matrix factorization, 23.1% compared to naive propensity estimator, and 22.1% compared to mixing
 114 datasets matrix factorization.

115 Additionally, we observe that NPE-MF and MF have very similar performance across all settings. This
 116 suggests that a IPS-based model will not gain any performance boost if propensities are inaccurately
 117 estimated. Furthermore, we argue that our model is data-efficient, i.e., it beats the competitor model
 118 (MD-MF), demonstrating that utilizing less biased dataset for propensity estimation is a more viable
 119 approach than mixing it with other data, though the latter is sometimes a common engineering
 120 practice. On our benchmarks, mixing datasets MF achieves less than 2% advantage compared to
 121 using only one dataset.

122 We conclude that our cross-dataset model significantly outperforms baselines and provides a efficient
 123 way to make use of more quantized (implicit) data.

124 6 Conclusion

125 Acknowledgements

126 References

127 [1] Charu C Aggarwal et al. *Recommender systems*, volume 1. Springer, 2016.

- 128 [2] Jiawei Chen, Hande Dong, Xiang Wang, Fuli Feng, Meng Wang, and Xiangnan He. Bias and de-
129 bias in recommender system: A survey and future directions. *arXiv preprint arXiv:2010.03240*,
130 2020.
- 131 [3] Guido W Imbens and Donald B Rubin. *Causal inference in statistics, social, and biomedical*
132 *sciences*. Cambridge University Press, 2015.
- 133 [4] Thorsten Joachims, Adith Swaminathan, and Tobias Schnabel. Unbiased learning-to-rank with
134 biased feedback. In *Proceedings of the tenth ACM international conference on web search and*
135 *data mining*, pages 781–789, 2017.
- 136 [5] Yehuda Koren, Robert Bell, and Chris Volinsky. Matrix factorization techniques for recom-
137 mender systems. *Computer*, 42(8):30–37, 2009.
- 138 [6] Karl Krauth, Sarah Dean, Alex Zhao, Wenshuo Guo, Mihaela Curmei, Benjamin Recht, and
139 Michael I Jordan. Do offline metrics predict online performance in recommender systems?
140 *arXiv preprint arXiv:2011.07931*, 2020.
- 141 [7] Bruno Pradel, Nicolas Usunier, and Patrick Gallinari. Ranking with non-random missing ratings:
142 influence of popularity and positivity on evaluation metrics. In *Proceedings of the sixth ACM*
143 *conference on Recommender systems*, pages 147–154, 2012.
- 144 [8] Tobias Schnabel, Adith Swaminathan, Ashudeep Singh, Navin Chandak, and Thorsten Joachims.
145 Recommendations as treatments: Debiasing learning and evaluation. In *international conference*
146 *on machine learning*, pages 1670–1679. PMLR, 2016.
- 147 [9] Masashi Sugiyama, Matthias Krauledat, and Klaus-Robert Müller. Covariate shift adaptation by
148 importance weighted cross validation. *Journal of Machine Learning Research*, 8(5), 2007.
- 149 [10] Steven K Thompson. *Sampling*, volume 755. John Wiley & Sons, 2012.
- 150 [11] Bernard Widrow, Istvan Kollar, and Ming-Chang Liu. Statistical theory of quantization. *IEEE*
151 *Transactions on instrumentation and measurement*, 45(2):353–361, 1996.
- 152 [12] Longqi Yang, Yin Cui, Yuan Xuan, Chenyang Wang, Serge Belongie, and Deborah Estrin.
153 Unbiased offline recommender evaluation for missing-not-at-random implicit feedback. In
154 *Proceedings of the 12th ACM conference on recommender systems*, pages 279–287, 2018.
- 155 [13] Jing Zhang, Zewei Ding, Wanqing Li, and Philip Ogunbona. Importance weighted adversarial
156 nets for partial domain adaptation. In *Proceedings of the IEEE conference on computer vision*
157 *and pattern recognition*, pages 8156–8164, 2018.