

Recommender Systems: Missing Data and Statistical Model Estimation

Benjamin M. Marlin

University of British Columbia
Vancouver, Canada
bmarlin@cs.ubc.ca

Richard S. Zemel

University of Toronto
Toronto, Canada
zemel@cs.toronto.edu

Sam T. Roweis

New York University
New York, NY

Malcolm Slaney

Yahoo! Research
Sunnyvale, CA
malcolm@ieee.org

Abstract

The goal of rating-based recommender systems is to make personalized predictions and recommendations for individual users by leveraging the preferences of a community of users with respect to a collection of items like songs or movies. Recommender systems are often based on intricate statistical models that are estimated from data sets containing a very high proportion of missing ratings. This work describes evidence of a basic incompatibility between the properties of recommender system data sets and the assumptions required for valid estimation and evaluation of statistical models in the presence of missing data. We discuss the implications of this problem and describe extended modelling and evaluation frameworks that attempt to circumvent it. We present prediction and ranking results showing that models developed and tested under these extended frameworks can significantly outperform standard models.

1 Introduction

The development of the world wide web, electronic commerce and social media has led to a dramatic increase in the amount of content available through the Internet. The web has tens of billions of indexed pages. Electronic commerce web sites like Netflix and Amazon contain tens of thousands to millions of items. Social media web sites like YouTube continue to add new content at astounding rates. As a result, the problem of matching people to the content that best meets their needs and interests is of great importance.

Classical information retrieval methods solve the problem of matching people to content based on explicit queries. This approach has been highly successful when both the content and queries are text-based, as in the case of web search [Brin and Page, 1998]. Classical recommender systems and collaborative filtering algorithms take a different approach: they match people to content based on preferences [Goldberg *et al.*, 1992]. Preferences are often expressed using explicit numerical ratings for individual content items. The collaborative aspect of this approach stems from the fact that it leverages the stored preferences of a whole community of users to make personalized predictions and recommendations

for each specific user. The personalization aspect of recommender systems makes them well suited to applications in electronic commerce and entertainment, while the fact that they do not rely on text-based descriptions of items makes them well suited to content like movies and music.

In this paper, we focus on a key problem in rating-based collaborative filtering: the possibility of a basic incompatibility between the properties of recommender system data sets and the assumptions required for valid estimation and evaluation of statistical models in the presence of missing data. We describe properties of recommender system data sets and relate them to the statistical theory of model estimation in the presence of non-random missing data. We describe an extended modelling framework and a modified set of evaluation protocols for dealing with non-random missing data. We present rating prediction and ranking results showing that models developed and tested under this extended framework can significantly outperform standard models.

2 Recommender Systems and Missing Data

The data collected in a recommender system can be thought of as a matrix with one row per user and one column per item. Since the items often number in the thousands to millions, most individual users naturally rate only a small percentage of the items. The marginal rating distribution of a data set is a simple summary statistic that shows the proportion of each rating value in the observed data. Marginal rating distributions for several well known data sets and web sites are shown in Figure 1 including EachMovie, MovieLens, Netflix, and YouTube.¹ Typically, these data sets contain a very low proportion of observed ratings (1% to 5%), and we see that the first four data sets all exhibit a skew towards high rating values. The YouTube data exhibits the largest skew with approximately 90% of the ratings taking the maximum rating value.

The first question that interests us is what accounts for this seeming overabundance of high rating values? To begin, we consider two hypothetical processes that could both generate the observed marginal rating distributions. First, most people really do like most items and the probability of observing a

¹The marginal rating distribution for YouTube was taken from: <http://youtube-global.blogspot.com/2009/09/five-stars-dominate-ratings.html>

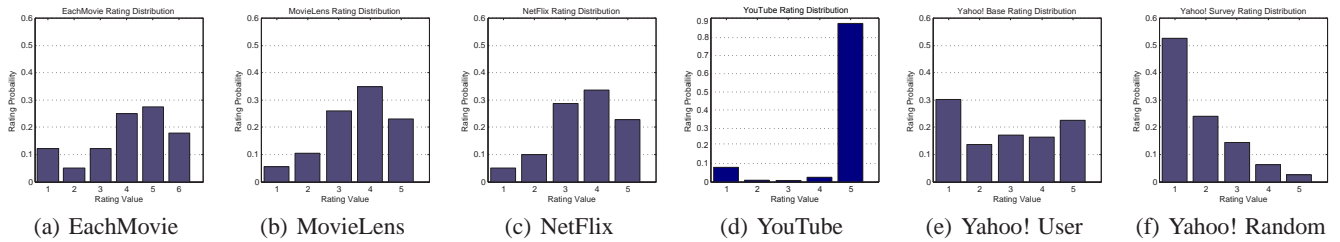


Figure 1: Distribution of rating values from several sources including EachMovie, MovieLens, Netflix, YouTube and two Yahoo! Music data sets.

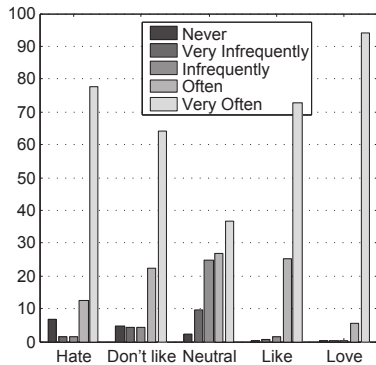


Figure 2: Results of a survey asking users to report how often they would choose to rate a song given their opinion of it. Each group of bars represents a different opinion level from *Hate it to Love it*. Each bar within a group represents a different frequency level from *Never* to *Very frequently*.

rating is independent of its value. Second, most people do not really like most items but the probability of observing a rating depends on its value. For example, the users of a movie recommender system may be more likely to watch and rate items they think they will like and less likely to watch and rate items they think they will dislike. This would create a systematic bias toward observing a disproportionate number of high rating values, explaining the skew in the first four data sets.

In previous work, we conducted a study of users of Yahoo! Music’s LaunchCast recommender system to explore these two hypotheses [Marlin *et al.*, 2007]. LaunchCast is a streaming music recommender system that generates continuous play-lists based on user ratings. The first part of the study consisted of a survey that asked users to report how often they would choose to rate a song given their opinion of it. The answers collected from 5400 users are summarized in Figure 2.²

The results show a clear dependence of rating frequency on the underlying preference level, lending support to our second hypothesis. Although high-rated items are reportedly rated

²Note that the study had over 35,000 respondents. The 5400 users included in the data set were those who, during normal use of the LaunchCast system, had rated at least 10 of the 1000 songs used in second stage of the study.

most often, the implied observation process is more complex than the example given above. The results indicate an apathy effect where users are more likely to supply ratings when their preferences are strongly positive or strongly negative, and less likely to supply ratings when their preferences are neutral. This can be explained by taking the user interaction model underlying the LaunchCast system into account. Since a user’s only control over song choice in the LaunchCast recommender system is to supply feedback by rating items, the user has a large incentive to rate items they both strongly like and strongly dislike to cause those items to be played more or less often. Items that the user feels neutral about require no action and are thus rated much less frequently.

Following the survey, users were presented with a set of ten songs selected at random from a total of 1000 songs and asked to rate them all. The artist name, song title and a thirty second audio clip were provided for each song. The marginal distribution shown in Figure 1(f) corresponds to the ratings collected during this study. It can be directly compared with the marginal distribution shown in Figure 1(e), which corresponds to existing ratings collected from the LaunchCast rating database for the same set of 1000 songs and 5400 users who participated in the study. The randomly-selected songs have a completely known (and uniformly random) observation process while the user-selected songs have an unknown observation process. We see dramatic differences between the two distributions with many fewer high rating values when songs are selected at random compared to selected by the user, again lending support to the hypothesis of a rating-value dependent observation process.

The question of how these results generalize to other recommender systems is an interesting one. As we noted above, the observation process implied by the LaunchCast survey results does appear to be rating-value dependent, but it is more complicated than the simple hypothesis that users are more likely to supply ratings for items that they like. In general, we believe that in a recommender system where users can choose what data they supply, the observation process is very likely to contain some form of rating-value dependent bias. The precise form of this bias will depend on the incentive to rate items of various quality. This incentive structure will in turn depend on the constraints and affordances built into the recommender system, as we have seen in the case of the LaunchCast system.

3 Missing Data Theory

The data presented in the previous section support the hypothesis that the probability of supplying a rating for an item is dependent on a user’s underlying rating for that item. In this section, we will formalize this idea and explore its impact on both statistical model estimation and the evaluation of rating prediction and ranking methods. We begin by defining the required notation.

We let N be the number of users in the data set, D be the number of items and V be the number of rating values. We denote the rating matrix by \mathbf{x} and the rating of user n for item d by x_{nd} . To reason about the observation process, we require a representation for missing and observed rating values. Following Little and Rubin [1987], we introduce a companion matrix of response indicators \mathbf{r} where $r_{nd} = 1$ if x_{nd} is observed, and $r_{nd} = 0$ if x_{nd} is not observed. We use \mathbf{X} and \mathbf{R} to denote random variables representing a rating matrix and matrix of response indicators.

The question of interest in this section concerns the joint probability distribution of the random variables \mathbf{X} and \mathbf{R} . This distribution can be factorized into the form shown below where μ and θ are the parameters of the joint distribution.

$$P(\mathbf{R} = \mathbf{r}, \mathbf{X} = \mathbf{x} | \mu, \theta) = P(\mathbf{R} = \mathbf{r} | \mathbf{X} = \mathbf{x}, \mu) P(\mathbf{X} = \mathbf{x} | \theta)$$

Little and Rubin refer to $P(\mathbf{R} = \mathbf{r} | \mathbf{X} = \mathbf{x}, \mu)$ as the *missing data model* (we have been referring to it as the observation process), while $P(\mathbf{X} = \mathbf{x} | \theta)$ is referred to as the *data model*. Standard maximum likelihood model estimation with missing data is based on ignoring the missing data model and optimizing the parameters of the data model given whatever elements of \mathbf{x} happen to have been observed. For this procedure to be valid, the missing data must be *missing at random* (MAR). The MAR condition asserts that the probability that a given random variable is missing depends only on the values of other random variables that are observed. The MAR condition is expressed below where the superscript *obs* indicates the observed entries in the given matrix.

$$P_{mar}(\mathbf{R} = \mathbf{r} | \mathbf{X} = \mathbf{x}, \mu) = P(\mathbf{R} = \mathbf{r} | \mathbf{X}^{obs} = \mathbf{x}^{obs}, \mu)$$

The MAR condition is best understood in the context of recommender systems in terms of the minimal circumstances where it fails to hold. Specifically, if the probability that a user will supply a rating for an item depends on the user’s underlying rating for that item, the MAR condition will fail to hold. The implications of a failure of the MAR condition are quite profound. The theory of missing data tells us that incorrectly ignoring the missing data model during parameter estimation will lead to provably biased estimates of the data model parameters [Little and Rubin, 1987].

The impact of violations of the MAR condition on model estimation and evaluation can be illustrated through a simple thought experiment. Consider a data set where a rating is only observed if it is five-stars (a slightly more extreme version of the YouTube data set). Standard empirical protocols for evaluating rating prediction and ranking in recommender systems are based on sub-sampling the observed data to form training and testing sets [Breese *et al.*, 1998]. Standard models (including non-parametric models like K-nearest neighbor

regression) will essentially learn that all items should be rated five-stars based on such a training set. Evaluating rating prediction performance on the corresponding test set (which also only contains five-star ratings) will yield zero error. However, the true task of interest is predicting ratings for all unrated items. In the worst case, all of the unrated items could actually have one-star ratings. A model that always predicts five-stars would then achieve the worst possible value of the prediction error over the set of unrated items.

In the ranking case, the true task of interest is to supply a ranking of all unrated items. Any rating prediction method can be used to produce rankings simply by sorting unrated items according to their predicted ratings. In the data set considered above, the trained model would carry no information about how to rank the items since all the predictions are the same. The test data would again not reveal the problem because all of the test items have the same maximal rating value. If the unrated items contain a small proportion of high ratings and a large proportion of low ratings, it is possible for a trained model to obtain arbitrarily poor accuracy on the true ranking task of interest.

These arguments show that violations of the missing at random assumption can significantly affect statistical model estimation as well as rating prediction and ranking evaluation when this evaluation is based on historical ratings.³ Dealing with non-random missing data in recommender systems thus requires both extended evaluation protocols and extended models. We turn to the development of models that explicitly incorporate a non-random missing data mechanism to reduce the bias in model estimation in Section 4. We then explore the question of evaluation protocols in Section 5.

4 Models for Non-Random Missing Data

The framework we consider for learning and prediction with non-random missing data follows the basic outline suggested by the theory of missing data [Little and Rubin, 1987]. We combine a probabilistic model for complete data, in this case a probabilistic clustering model, with a probabilistic model of the missing data process. We consider two missing data models that can represent a direct dependence between the probability of rating an item and its underlying rating value.

4.1 Data Model

We use a probabilistic clustering model for the data model, which has a very natural interpretation in the collaborative filtering domain. A cluster can simply be thought of as a collection of users that express similar preferences over the full set of items. In a movie recommender system, for example, clusters may reflect preferences for specific genres like action, drama, science fiction and so on.

A finite multinomial mixture model is a probabilistic clustering model for discrete data. In this model, an unobserved or latent variable Z_n is associated with every user n , indicating the cluster to which user n belongs. We assume that there are a fixed, finite number of clusters K . The parameters of

³It is worth pointing out that organizations operating their own recommender systems are not limited to evaluation based on historical ratings, as they can carry out tailored user studies.

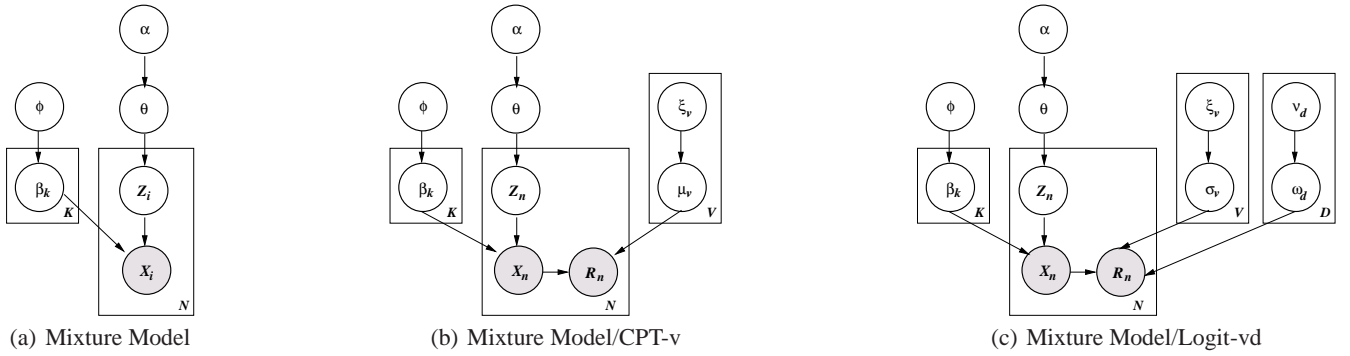


Figure 3: Graphical models illustrating the basic multinomial mixture model, the multinomial mixture/CPT-v model, and the multinomial mixture/Logit-vd model.

each cluster β_k specify the preferences of a prototypical user that belongs to cluster k . Specifically, β_{vdk} is the probability $P(x_{nd} = v | Z_n = k)$ that user n will assign rating v to item d under the assumption that user n belongs to cluster k . There is a discrete distribution over the clusters with parameters θ where $\theta_k = P(Z_n = k)$ is the prior probability that a user will belong to cluster k . We give the probabilistic model for the multinomial mixture model in Equations 1 to 4. The corresponding graphical model is pictured in Figure 3(a).

$$P(\theta | \alpha) = \mathcal{D}(\theta | \alpha) \quad (1)$$

$$P(\beta_{dk} | \phi_{dk}) = \mathcal{D}(\beta_{dk} | \phi_{dk}) \quad (2)$$

$$P(Z_n = k | \theta) = \theta_k \quad (3)$$

$$P(\mathbf{x}_n | Z_n = k, \beta) = \prod_{d=1}^D \prod_{v=1}^V \beta_{vdk}^{[x_{nd}=v]} \quad (4)$$

For the purpose of model estimation, the parameters θ and β are given prior distributions that act like regularization functions and smooth the estimated probability parameters away from extreme values. Both prior distributions are Dirichlet distributions denoted by \mathcal{D} . The square bracket notation $[s]$ represents an indicator function that takes the value 1 if the statement s is true, and 0 if the statement s is false.

The default when dealing with missing data in a mixture model is to invoke the missing at random assumption. Under the missing at random assumption, the missing data model is ignored and inference, learning, and prediction can be based on the observed data only.

4.2 Missing Data Models

The basic mixture model can be augmented with an explicit model of the missing data process when the MAR condition is not believed to hold. We consider two missing data models that we refer to as *CPT-v* and *Logit-vd* due to their parameterizations. In the CPT-v missing data model, the probability that a rating is observed depends only on its underlying value. The model can be thought of in terms of a set of biased coins, one for each rating value v . Coin v has a probability of coming up heads given by the parameter μ_v . To determine if a rating with value v will be observed, we flip coin v . We can

achieve different rating-dependent missing data processes by assigning different values to the parameters μ_v .

This simple coin flip model corresponds to a Bernoulli likelihood on each response indicator variable conditioned on the corresponding rating value, as given in Equation 5. The model is defined through the conditional probability table specified by μ , hence the name of the model. The multinomial mixture data model augmented with the CPT-v missing data model is shown in Figure 3(b). The model includes a conjugate Dirichlet prior on each μ_v .

$$P(\mathbf{r}_{nd} = 1 | \mathbf{x}_{nd} = v, \mu) = \mu_v \quad (5)$$

$$P(r_{nd} = 1 | x_{nd} = v, \sigma, \omega) = \frac{1}{1 + \exp(-(\sigma_v + \omega_d))} \quad (6)$$

The Logit-vd model shown in Equation 6 is a generalization of CPT-v that allows the observation probability to vary across items in a restricted fashion. The model includes a real-valued item non-random missing data factor σ_v and a real-valued item popularity factor ω_d . The two factors are combined through a logistic function to yield the observation probability for each item d and rating value v . The multinomial mixture model augmented with the Logit-vd missing data model is shown in Figure 3(c). The model includes an independent Gaussian prior on each parameter.

It is important to note that both of these models are highly simplified. While they can each represent a rating-value dependent missing data process and Logit-vd can model some variation across items, both models ignore the possibility that ratings for multiple items might influence whether a particular rating value is observed. They also ignore side information about users and items that might influence whether or not ratings for particular items will be observed. Neither type of information is available for the data set we consider, but an advantage of a probabilistic approach is that basic models can easily be extended to deal with additional features and side information should it be available.

4.3 Model Estimation

Locally optimal maximum likelihood estimates for the basic multinomial mixture model can be computed under the missing at random assumption using a standard Expectation Max-

imization (EM) algorithm [Dempster *et al.*, 1977]. The per-iteration complexity of the algorithm scales linearly with the number of observed ratings. In the case of the multinomial mixture model combined with the CPT-v and Logit-vd missing data models, efficient EM algorithms can be also derived where the computational complexity per iteration is dominated by the number of observed ratings. This is the main advantage of using simplified missing data models. We use the EM algorithm to simultaneously learn the parameters of both the data and missing data models in all of the experiments described in the following sections. Further details regarding the estimation and prediction algorithms for these models can be found in our previous work [Marlin *et al.*, 2007; Marlin and Zemel, 2009].

5 Evaluation Protocols

As described in Section 3, standard empirical protocols for rating prediction and ranking evaluation can lead to biased performance estimates in the presence of non-random missing data, necessitating modified empirical protocols. In the case of rating prediction, we require a test set that is as close as possible to a random selection of unrated items. The ratings for randomly selected items collected during the Yahoo! Music user study described in Section 2 provide just such a test set since the expected overlap between randomly-selected items and previously-rated items is low.

We also propose the use of ratings for randomly selected items for the evaluation of ranking accuracy, although this choice presents some issues. In particular, since we only have ten items per user and most of the items in the test set have low ratings, the ranking evaluation may unduly reflect the model’s ability to discriminate between items with low rating values. However, we feel this is preferable to measuring the ranking performance on a subset of the observed data subject to a completely unknown observation process. Whether it is possible to construct better test sets for ranking evaluation given both sources of ratings is an open question.

The full empirical protocol uses a training set containing the 5400 users who participated our study (described in Section 2), plus an additional 10000 LaunchCast users selected at random from those with at least 10 ratings on the 1000 songs used in the study data set. All of the training set ratings are ratings for items selected by the users during normal interaction with the LaunchCast music recommender system (Figure 1(e)). The validation and test sets contain ratings subsampled from the ratings for randomly-selected items collected during the user study for each of the 5400 study users (Figure 1(f)).

The models we evaluate include the multinomial mixture model under the MAR assumption (MM/MAR), as well as the multinomial mixture model combined with the CPT-v (MM/CPT-v) and Logit-vd (MM/Logit-vd) missing data models. We also evaluate two very common collaborative filtering models that implicitly assume random missing data: a matrix factorization (MF) model [Salakhutdinov and Mnih, 2008] and an item-based K-nearest neighbor method (iKNN) [Sarwar *et al.*, 2001].

We train each mixture-based model using 1, 5, 10 and 20

mixture components and select the best setting using cross validation. The prior parameters for all of the mixture-based models were set to yield broad priors. For the matrix factorization model, we considered ranks $K = 1, 5, 10, 20$ and regularization parameters 0.1, 1, 5, 10 and selected the best values by cross validation. For the item-based KNN method, we use an adjusted cosine similarity metric [Sarwar *et al.*, 2001], combined with the standard weighted nearest neighbor prediction rule.

Once the models are trained, we condition on the training set ratings for each user and predict the ratings for each of that user’s test items. We form a ranked list of test items for each user by sorting that user’s test items according to their predicted ratings.

6 Results

Rating Prediction: We evaluate rating prediction performance in terms of normalized mean absolute error (NMAE). This error measure is proportional to the average absolute difference between actual and predicted ratings. NMAE is computed as seen below. We assume there are T test items per user with indices $i(1, n)$ to $i(T, n)$. The normalizing constant (equal to 1.6) is the expected MAE assuming uniformly distributed predictions and true ratings. Note that *lower* NMAE indicates better prediction performance.

$$NMAE = \sum_{n=1}^N \sum_{t=1}^T \frac{|x_{ni(t,n)} - \hat{x}_{ni(t,n)}|}{1.6NT}$$

Figure 4(a) shows the test NMAE score for each of the five models. We select the optimal complexity for each model based on cross validation NMAE scores. Standard errors are represented on the plots using error bars. We see that the MM/Logit-vd and MM/CPT-v models, which do not assume the MAR condition, drastically outperform the MM, iKNN and MF models, which do assume random missing data, when measuring performance on ratings for randomly selected items.

Ranking: We evaluate the ranked lists of test items produced by each method using a standard ranking accuracy measure, the normalized discounted cumulative gain (NDCG). NDCG@L measures how well the predicted ranking matches the true ranking (obtained by sorting the items by their actual ratings) for a ranked list of length L . NDCG places more emphasis on ranking errors at the top of the ordering and is normalized so that the true ranking yields an accuracy of 1. The NDCG@L score is computed as seen below where $\pi(l, n)$ is the index of the item with rank l when test items are sorted in descending order by true rating x_{nd} , $\hat{\pi}(l, n)$ is the index of the item with rank l when items are sorted in descending order according to their predicted ratings \hat{x}_{nd} . When sorting by true and predicted ratings, ties can be broken arbitrarily without affecting the NDCG@L score. Note that *higher* NDCG indicates better ranking performance.

$$NDCG@L = \sum_{n=1}^N \frac{\sum_{l=1}^L (2^{x_{n\hat{\pi}(l,n)}} - 1) / \log(1+l)}{N \sum_{l=1}^L (2^{x_{n\pi(l,n)}} - 1) / \log(1+l)}$$

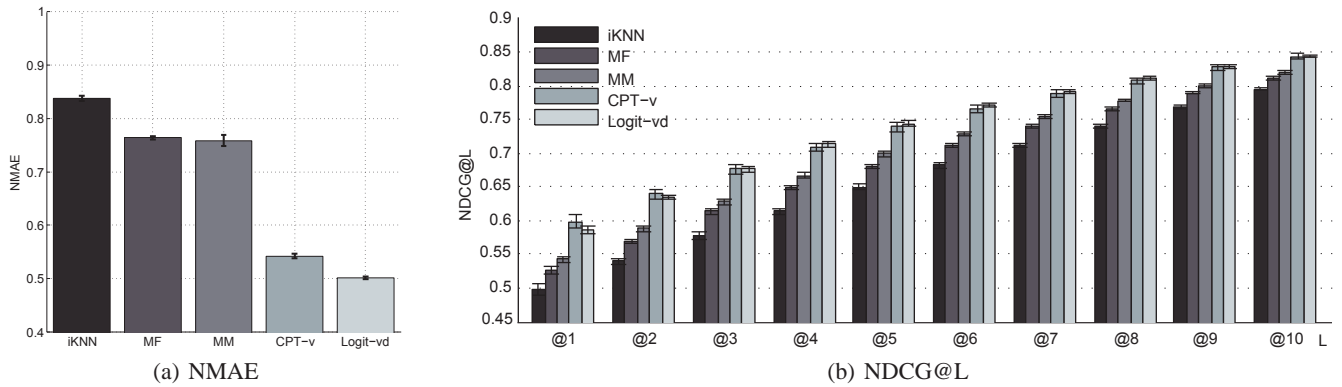


Figure 4: Figure (a) presents the test set rating prediction error on randomly selected items. Figure (b) presents the test set NDCG@L results on randomly selected items. The methods are iKNN, MF, MM/MAR, MM/CPT-v, and MM/Logit-vd.

Figure 4(b) shows the test NDCG@L performance for each model estimated on lists of ratings for the 10 randomly selected items. We select the optimal complexity for each model based on cross validation NDCG@L scores. The results again show that the MM/Logit-vd and MM/CPT-v models, which do not assume the MAR condition, outperform the MM, iKNN and MF models, which do assume random missing data.

7 Conclusions

In this paper, we have explored properties of the missing data process in recommender systems, discussed their impact on the validity of standard statistical model estimation and evaluation procedures, and described and tested extended modeling and evaluation frameworks that seek to overcome the problems caused by non-random missing data. The development of more sophisticated models within the extended modeling framework is of great interest, as is the design of better test sets for ranking. The question of how non-random missing data affects methods for learning to rank is also of great interest. The continued convergence of recommender systems, content-based search and social networks raises the question of the extent to which the detrimental effects of non-random missing ratings can be mitigated by incorporating additional sources of information including content-based features for items (including social tags) and information about both individual users and the relationships between users. This is a very interesting direction for future research.

Acknowledgments

We are indebted to Sam Roweis for his contributions to this research. He was a wonderful colleague, mentor and friend. We thank the Natural Sciences and Engineering Research Council of Canada, the Killam Trusts, the Pacific Institute for the Mathematical Sciences and the Canadian Institute for Advanced Research for funding this work. We thank Yahoo! for making the data set used in this work, *Yahoo! Music Ratings for User Selected and Randomly Selected songs, V 1.0*, available through the Webscope Program.

References

- [Breese *et al.*, 1998] J. S. Breese, D. Heckerman, and C. Kadie. Empirical Analysis of Predictive Algorithms for Collaborative Filtering. In *Proceedings of the 14th Conference on Uncertainty in Artificial Intelligence*, pages 43–52, 1998.
- [Brin and Page, 1998] S. Brin and L. Page. The anatomy of a large-scale hypertextual web search engine. In *Proceedings of the seventh international conference on World Wide Web 7*, pages 107–117, 1998.
- [Dempster *et al.*, 1977] A.P. Dempster, N.M. Laird, and D.B. Rubin. Maximum likelihood from incomplete data via the EM algorithm. *Journal of the Royal Statistical Society, Series B*, 39(1):1–38, 1977.
- [Goldberg *et al.*, 1992] D. Goldberg, D. Nichols, B. M. Oki, and D. Terry. Using collaborative filtering to weave an information tapestry. *Communications of the ACM*, 35:61–70, December 1992.
- [Little and Rubin, 1987] R. J. A. Little and D. B. Rubin. *Statistical Analysis with Missing Data*. John Wiley and Sons, Inc., 1987.
- [Marlin and Zemel, 2009] B. M. Marlin and R. S. Zemel. Collaborative prediction and ranking with non-random missing data. In *Proceedings of the third ACM conference on Recommender systems*, pages 5–12, 2009.
- [Marlin *et al.*, 2007] B. Marlin, R. Zemel, S. Roweis, and M. Slaney. Collaborative filtering and the missing at random assumption. In *Uncertainty in Artificial Intelligence 23*, 2007.
- [Salakhutdinov and Mnih, 2008] R. Salakhutdinov and A. Mnih. Probabilistic matrix factorization. In *Advances in Neural Information Processing Systems*, volume 20, 2008.
- [Sarwar *et al.*, 2001] B. Sarwar, G. Karypis, J. Konstan, and J. Reidl. Item-based collaborative filtering recommendation algorithms. In *Proceedings of the 10th international conference on World Wide Web*, pages 285–295, 2001.