# On Variational Inference for User Modeling in Attribute-Driven Collaborative Filtering

Venugopal Mani*
*Walmart Labs*
*Sunnyvale, CA, USA*
venugopal.mani@walmartlabs.com

Ramasubramanian Balasubramanian*
*Walmart Labs*
*Sunnyvale, CA, USA*
r.balasubramanian@walmartlabs.com

Sushant Kumar
*Walmart Labs*
*Sunnyvale, CA, USA*
skumar4@walmartlabs.com

Abhinav Mathur
*Walmart Labs*
*Sunnyvale, CA, USA*
amathur1@walmartlabs.com

Kannan Achan
*Walmart Labs*
*Sunnyvale, CA, USA*
kachan@walmartlabs.com

*Abstract*—Recommender Systems have become an integral part of online e-Commerce platforms, driving customer engagement and revenue. Most popular recommender systems attempt to learn from users' past engagement data to understand behavioral traits of users and use that to predict future behavior. In this work, we present an approach to use causal inference to learn user-attribute affinities through temporal contexts. We formulate this objective as a Probabilistic Machine Learning problem and apply a variational inference based method to estimate the model parameters. We demonstrate the performance of the proposed method on the next attribute prediction task on two real world datasets and show that it outperforms standard baseline methods.

*Index Terms*—Recommender Systems, Variational Methods, Collaborative Filtering, Bayesian Statistics

## I. INTRODUCTION

Recommender Systems have traditionally been studied from the lens of attempting to increase customer engagement by user modeling from past interactions. These interactions are often collected in terms of explicit user signals like ratings and item reviews. Recently, there has been a shift in literature towards building recommenders by using implicit user signals like item views, item purchases, etc. Implicit signals, while useful in increasing the coverage of user signals over items, can suffer from lack of definition of what constitutes a negative signal. This has led to a class of problems known as One Class Collaborative Filtering [1] where techniques like low rank approximation and negative sampling are used to improve user understanding by eliminating the ambiguity over the negative train samples.

A common assumption in Implicit OCCF is that all positive signals are equal. However, this assumption can fail to capture the wide ranging spectrum of user interactions in some domains. Normalization techniques do exist to scale these but there exists scope for more nuanced values for the positive samples. With modern data collection capabilities, *domain specific fine tuning* of user interactions can be achieved to

further our understanding of abstract concepts about users (like loyalty, satisfaction with the product, etc). One such idea was introduced in [2] where the concept of *long term customer satisfaction* was defined through a function to track the continuous implicit signal of the user.

While Implicit OCCF systems are quite effective, an often cited drawback of these systems for user modeling is their lack of sensitivity to temporally changing user behavior. The traits of a user from a few months prior need not necessarily model their present behavior. It follows logically to try to encode the temporal aspect of implicit signals into a user understanding objective and optimize over it. This is particularly relevant in the subset of attribute-driven collaborative filtering as users tend to develop repeat patterns on certain attributes of an item. For example, in the domain of music recommendation, affinity to artists has been studied, and the recommendation of artists similar to the ones the user is loyal to has also seen improvements in results [3]. In this work, we try to further our understanding of users through the notion of *temporal loyalty* and integrate it into the attribute-driven collaborative filtering framework. We optimize the objective from two sources : the transaction matrix which is a binary matrix that indicates past interaction (or lack of thereof) of the involved user with the item's attribute as well as a temporal loyalty matrix which attempts to capture drifting user loyalty over time.

The contributions of this work are as follows: first, we model temporal loyalty of the users to augment the transaction matrix. Then, we demonstrate that optimization using variational inference over these matrices outperforms plain collaborative filtering based methods on the next attribute prediction task, thus leading to a better understanding of user preferences. The rest of the work is organized as follows: Section II delves into the literature of related work, Section III describes our proposed system model, Section IV describes our experiments on two real world datasets , Section V analyzes the results, and Section VI concludes the work and describes possible future directions.

---

* Both authors contributed equally to this work.

## II. Related Work

The idea of optimizing over two matrices for modeling user preferences is relatively new. There is active research around the kind of domain-specific objectives to be optimized for and the corresponding data that could be augmented. The authors of [2] consider measures of satisfaction with the purchased items, such as the amount of time spent playing a game or the number of times a particular artist was heard. Other works focus on tasks like using dwell time in session-based recommendations [4], [5] or to enrich the user-item matrix [6], leveraging implicit signals such as internet browsing logs [7], etc. On the other hand, several works exist that leverage the binary transaction matrix to tackle the well-known top-k recommendation problem in large-scale datasets, such as those dealing with memory-based collaborative filtering for explicit feedback [8], item-based collaborative filtering to address scalability concerns [9], using stratified SGD to deal with large-scale matrix factorization [10], etc. However, these rely solely on explicit user signals and fail to incorporate any temporal signals.

In our work, we use the temporal loyalty to an item attribute as the second matrix, thus leveraging both explicit and temporal signals to set up an optimization over the two matrices. The use of loyalty is motivated by works such as [11], where the authors model consumers' repeat purchase behavior, as well as our experience in the domain of e-Commerce and grocery. Attribute-based collaborative filtering has been explored before in works such as [12] where the authors use categorical attributes to improve recommendation through multi-task learning or hierarchical classification, and [13] which deals with attribute-aware collaborative filtering. Our work captures the changing affinity of the users to these attributes, and thus could be used as a first stage in hierarchical classification algorithms: to predict which brands the users will buy next, before recommending particular items of that brand.

In terms of the application of variational inference and Bayesian statistics to solve collaborative filtering problems, most works focus on the use of Variational Auto Encoders. For example, [14] introduces a generative model with a multi-nomial likelihood and uses Bayesian inference for parameter estimation, [15] uses VAE to alleviate the problem of poor robustness and over-fitting caused by large-scale data, etc. Other works using Bayesian inference, such as [16], which presents a scalable inference for Variational Bayesian matrix factorization with side information, or [17], which proposes a distributed memo-free variational inference method for large-scale matrix factorization problems, address some of the well-known shortcomings of the same in recommender systems.

The experimental framework that we have adopted is called the Box's Loop [18], which is used to uncover patterns from the conditional distribution of a latent variable model and use them to model the data and make predictions. This was well suited for our problem, since we assume a particular

structure of the latent variables and use that to explain the data, and then use those variables for future recommendations. Related works in the sub-field of latent variable modeling for recommender systems include works such as [19], which uses blind regression to complete the partial user-item interaction matrix and uses the features of the users and items as the latent variables, and [20], which presents a Bayesian latent variable model for rating prediction that models ratings over each user's latent interests and each item's latent topics. Embedding based approaches to model the users and items have also been tried, in works such as [21], which replaces the inner product between user and item latent features used in classic matrix factorization by a neural architecture, and [22], which couples deep feature learning and deep interaction modeling with a rating matrix to improve recommendation performance. Our work fits the optimization task into this framework and generates user and item attribute embeddings that explain the data well by applying variational inference, and these user representations thus obtained help us develop a better understanding of the user preferences.

## III. System Model

### A. The Top-k Attribute Recommendation Problem

The classic top-k recommendation problem can be defined as follows : given a catalog of items $C$ containing items $i_1, i_2, \cdots, i_n$ and an item $a \in C$, henceforth referred to as the *anchor item*, finding a ranked list of distinct items $i_1, i_2, \cdots, i_k \in C$ to be offered alongside the anchor item, such that the user of an e-Commerce platform is most likely to engage with them. This engagement of users can be defined by a variety of metrics (in our case, the future purchase of the item).

A sub-problem of the top-k recommendation problem is the attribute recommendation problem. Rather than the recommendation of the items to cause the next most likely engagement, the task is to predict on the attribute. For instance, in the space of e-Commerce, the attribute of interest could be the brand of the item. Given a user $u$, the task would be to predict the $k$ most likely brands that they're likely to engage with. Attribute recommendation is a slightly more well defined space as, particularly with certain attributes like brands, users are very likely to develop behaviors like loyalty towards certain brands. While in the item space, repurchasing of an item is a commonly seen behavior, especially in domains like grocery, the re-engagement with attributes is a much more commonly exhibited behavior in a wide variety of recommender system domains (songs by the same artist, books by the same author, movies starring the same actor) and solving it can hence involve harvesting this richer behavioral signal.

### B. The Temporal Loyalty Matrix

Most recommender systems dealing with grocery data focus on the *binary transaction matrix*. Given a set of users

$u_1, u_2, \cdots, u_m$ and a set of items $i_1, i_2, \cdots, i_n$, the $(p,q)^{th}$ element of the transaction matrix is 1 if $u_p$ has purchased $i_q$ within the training window, and 0 otherwise. Oftentimes, user-based collaborative filtering algorithms that use such matrices fail to pick up important signals such as a user's loyalty to a brand, how their behavior changes dynamically with time, etc. In this work, we try to capture those signals by optimizing over a temporal loyalty matrix, L, in addition to the binary transaction matrix, T. Since we are dealing with an attribute recommendation problem, consider a set of users $u_1, u_2, \cdots, u_m$ and an associated set of attribute values $v_1, v_2, \cdots, v_n$ for a particular attribute (say, brand). The $(p,q)^{th}$ element of the transaction matrix,

$$
T_{pq} = \begin{cases} 1, & \begin{array}{l} u_p \text{ has bought an item with at-} \\ \text{tribute value } v_q \text{ at least once in the} \\ \text{training window} \end{array} \\ \\ 0, & \text{otherwise} \end{cases} \tag{1}
$$

The $(p,q)^{th}$ element of the temporal loyalty matrix is the time-decayed sum of all the purchases of attribute value $v_q$ made by user $u_p$, that is,

$$
L_{pq} = \begin{cases} \sum\limits_{t=t_1}^{t_k} 2^{\frac{t-t_{start}}{t_{end}-t_{start}}}, & \begin{array}{l} u_p \text{ has bought an item with} \\ \text{attribute value } v_q \ k(\geq 1) \\ \text{times in the training win-} \\ \text{dow} \end{array} \\ \\ 0, & \text{otherwise} \end{cases} \tag{2}
$$

The variables $t_{start}$ and $t_{end}$ in Equation 2 represent the start and end times of the training window, and $t_1, t_2, ..., t_k$ are the time instances when the user purchased items with that attribute value. For instance, a particular brand of beer purchased over a year ago should not get the same weight as the one purchased a week ago as user preferences might have changed. We further show that this framework works well for recommender systems that have some notion of loyalty/preference, such as readers' predilection for certain authors, etc. This optimization [2] balances data from both the transaction matrix and the temporal loyalty matrix.

### C. The Bayesian Framework

Probabilistic Machine Learning (PML) is a sub-field of Machine Learning where domain knowledge and assumptions about the hidden structure of data are leveraged to explain the observed data. PML models large, interesting, and interconnected datasets at scale.

The iterative probabilistic pipeline, coined *Box's Loop* by [18] lists the steps of modeling a PML pipeline as *positing a model* with assumptions about the hidden structure of data, *inferring* the hidden variables, and *criticizing* the model (the evaluation step). If the evaluation does not meet the standard

required, the values of the hidden variables are revised to better explain the data at hand.

The structure of collaborative filtering by Matrix Factorization effectively lends itself to Box's Loop. The entries of the transaction matrix and the temporal loyalty matrix constitute the observed data. The classic matrix factorization problem involves decomposition of a given Matrix $M$ into latent factor matrices $U$ and $V$ along with their respective biases $B_u$ and $B_v$. The matrices $U$ and $V$ are known as *embedding matrices* in the setting of collaborative filtering. The learning task then becomes to learn the embeddings and biases such that their probability, given the observed transaction and temporal loyalty matrices is maximized. This is also known as the *posterior distribution*.

### D. The Posterior Distribution

We chose to model both the matrices, T and L, as well as the priors with normal distributions. This is logical because the values in the L matrix are continuous and distributions from the exponential family have shown good results in literature [2]. Also, the normal family of distributions is conjugate to itself (or self-conjugate) with respect to a normal likelihood function, and conjugacy has desirable properties, such as yielding a closed-form expression for the posterior.

Fig. 1 is a probabilistic graphical representation of our latent variable model. It shows how the random variables depend on each other in our generative process. Thus, it helps us form the posterior by connecting the assumptions that we made about the data to the model. The components of this graphical model are the ones used in standard graphical models in the field of machine learning, such as [18]: the nodes represent random variables, the edges represent a dependence between the nodes that they connect, and the plates denote replication. Each entry in the transaction matrix, $T_{pq}$, and the temporal loyalty matrix, $L_{pq}$, depends only on its local variables $u_p$, $bu_p$, $v_q$, and $bv_q$, which are the embedding and the bias vectors of the $p^{th}$ user, and the embedding and the bias vectors of the $q^{th}$ attribute respectively, and the corresponding global variables ($\kappa$ and $\psi$), as is the case with conditionally conjugate models. $\kappa$ and $\psi$ are the scale and the location parameters, which allow the distributions of T and L to have different dynamic ranges and be centered around different means, despite sharing some parameters which model the positive correlation between the transactions and the temporal loyalty scores, as seen in works like [2].

As mentioned earlier, all these variables have normal priors with mean 0 (except the scale parameters, which have a mean of 1) and a variance that depends on a hyperparameter, denoted by $\alpha$ with a subscript corresponding to the variable, as seen in equation 3. Utilising a modeling strategy similar to [2], we write out the posterior in equation 3 as being proportional to the product of the likelihoods and the priors. H is the set of all hyperparameters: those represented by solid black circles in Figure 1, $\gamma$, and $\beta$. $\gamma$ allows us to control how much

importance we give to the two likelihoods relative to each other, and $\beta$ is used to model the variance. The matrices T and L constitute the observed data, represented by grey circles in Fig. 1. $\theta$ is the set of latent variables: the user and item embeddings and biases, and the location and scale parameters, represented by white circles in Fig. 1. Each user and item vector is of dimension d.

$$
\begin{aligned}
P(\theta|T,L,H) &\propto P(T,L|\theta,H)P(\theta|H) \\
&= \prod_{(p,q,T_{pq})\in T} P(T_{pq}|u_p,v_q,bu_p,bv_q,\kappa_t,\psi_t,H) \\
&\quad \prod_{(p,q,L_{pq})\in L} P(L_{pq}|u_p,v_q,bu_p,bv_q,\kappa_l,\psi_l,H) \\
&\quad \prod_{p=1}^{m}[P(u_p|\alpha_u)P(bu_p|\alpha_{bu})]\prod_{q=1}^{n}[P(v_q|\alpha_v)P(bv_q|\alpha_{bv})] \\
&\quad P(\kappa_t|\alpha_{\kappa_t})P(\psi_t|\alpha_{\psi_t})P(\kappa_l|\alpha_{\kappa_l})P(\psi_t|\alpha_{\psi_l}) \\
&= \prod_{(p,q,T_{pq})\in T} \mathscr{N}(\kappa_t(u_p^T v_q + bu_p + bv_q)+\psi_t,(\gamma\beta)^{-1}) \\
&\quad \prod_{(p,q,L_{pq})\in L} \mathscr{N}(\kappa_l(u_p^T v_q + bu_p + bv_q)+\psi_l,((1-\gamma)\beta)^{-1}) \\
&\quad \prod_{p=1}^{m}[\mathscr{N}(0,\alpha_u^{-1}\boldsymbol{I_d})\mathscr{N}(0,\alpha_{bu}^{-1})] \\
&\quad \prod_{q=1}^{n}[\mathscr{N}(0,\alpha_v^{-1}\boldsymbol{I_d})\mathscr{N}(0,\alpha_{bv}^{-1})] \\
&\quad \mathscr{N}(1,\alpha_{\kappa_t}^{-1})\mathscr{N}(0,\alpha_{\psi_t}^{-1})\mathscr{N}(1,\alpha_{\kappa_l}^{-1})\mathscr{N}(0,\alpha_{\psi_l}^{-1})
\end{aligned}
\tag{3}
$$

### E. Variational Inference

The posterior $P(\theta|T,L,H)$ from Equation 3 solves for the family of high dimensional latent variables $\theta$, given the initial prior distributions over the latent variables and a likelihood function $P(T,L|\theta,H)$ that we posit about the model. The direct application of Bayes' theorem has the problem of an intractable high dimensional integration in the denominator and hence an approximate Bayesian inference of the posterior is carried out. One of the most popular methods for approximate Bayesian inference is the Markov Chain Monte Carlo (MCMC). However, despite the high accuracy of MCMC, the scale of our problem means that it is virtually impossible to use due to the computational time involved.

Instead, a scalable approach is to treat the posterior approximation as an optimization problem through variational inference [23], [24]. The objective of variational inference is to find the *variational distribution* which is a proxy-posterior $q$ parametrized by $\nu$, such that the variational distribution is least-divergent from the true posterior $p$. We adopted the widely used Kullback–Leibler divergence (KL Divergence) as the divergence metric between the two distributions. The KL divergence term is an intractable one and the equivalent of minimizing the KL-Divergence is the maximization of the

Evidence Lower Bound (ELBO) [25]. The ELBO, $\mathscr{L}(\nu)$, is described in equation 4.

$$
\mathscr{L}(\nu) = \mathbb{E}_q[log(p(T,L|\theta))] - KL(q(\theta;\nu)\parallel p(\theta)) \tag{4}
$$

The terms here provide the classical Bayesian trade off between the log likelihood of the data and the prior over the parameters of the model. That is, the first term tries to maximize the likelihood of the observed transactions and the temporal loyalty scores, given the embedding vectors. The second term is the KL divergence between the the variational distribution and the prior over the embedding vectors. The second term effectively acts as a regularizer as it tries to minimize the divergence from the prior and hence prevents the optimizer from converging to the maximum likelihood estimate.

Stochastic gradient descent is the commonly used approach to optimize the ELBO objective. Some works [2] also use coordinate descent and other variants of gradient descent to compute the gradients and update the parameters. The gradients can be obtained by rewriting equation 4 in terms of the complete log likelihood and then computing the gradient, as shown in equation 5.

$$
\nabla_\nu \mathscr{L}(\nu) = \nabla_\nu \mathbb{E}_q[log(p(T,L,\theta)) - log(q(\theta;\nu))] \tag{5}
$$

In this work, we use score function gradient estimators [26], [27], by rewriting equation 5 as

$$
\nabla_\nu \mathscr{L}(\nu) = \mathbb{E}_q[\nabla_\nu log(q(\theta;\nu))(log(p(T,L,\theta)) - log(q(\theta;\nu)))] \tag{6}
$$

### F. Prediction function

Once the variational distribution is approximated, predictions are made from the *posterior predictive function*. The posterior predictive function uses the likelihood function from Equation 3, $P(T,L|\theta,H)$ to generate the predictions. After the estimation of the latent variables $\theta$, the values of the transaction entry $T_{pq}$ and temporal loyalty entry $L_{pq}$ for user $p$ and attribute $q$ are estimated from the distributions $\mathscr{N}(\kappa_t(u_p^T v_q + bu_p + bv_q)+\psi_t,(\gamma\beta)^{-1})$ and $\mathscr{N}(\kappa_l(u_p^T v_q + bu_p + bv_q)+\psi_l,((1-\gamma)\beta)^{-1})$ respectively. Once the two values are determined, a simple addition of the two values gives the overall score for that particular user-attribute pair.

## IV. EXPERIMENTS

### A. Datasets and preprocessing

We demonstrate our results on two datasets from different domains. The first is a private dataset from a large-scale e-Commerce company. We collected six months' worth of grocery transaction data. From the transaction metadata, the customer id, the brand, the transaction date, and the event epoch (exact epoch at which the transaction happened) were
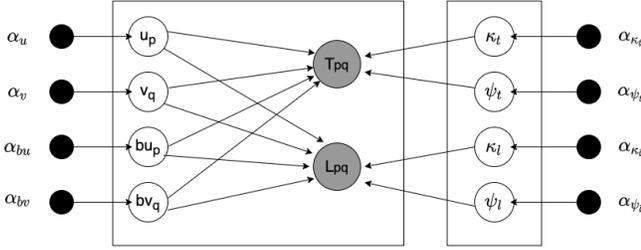
Fig. 1: A graphical representation of the proposed latent variable model

chosen. Since we are trying to understand and model loyalty, we decided to filter out customers that didn't have more than a threshold number of items in their basket, thus keeping only *engaged customers* in the dataset. From experience, we have seen that there are a few large merchants and clients, such as grocery stores, that place bulk orders. These would not be representative of a single customer and hence we decided to filter out those customers that had more than an upper threshold of transactions as well. We used the first five months of data as our train set and the final month as the test set. We also filtered out customers that weren't present in both the train and the test set, since we would not have embeddings for customers we have not seen. The other problem that both the baseline models and our model suffered from was the introduction of new brands during the test period, which happens due to change in consumer demand, seasonality effects, etc. We removed the brands not present in both the train and the test sets as well. This resulted in a dataset containing 180,000 customers and 11,200 brands.

The second dataset is the publicly available Goodreads Book Reviews dataset. This contains ratings, reviews, and a lot of other attributes of the items and the users, such as user-book interactions, metadata of the books, etc., collected in 2017 by scraping users' public shelves on Goodreads [28], [29]. The group that collected the dataset recommends using a subset (by genre) of the dataset, as the entire dataset is really large. In keeping with our theme of loyalty, we decided to go ahead with the 'fantasy and paranormal' genre. Here, we are trying to assess readers' loyalty to authors. And this genre had a high density of such interactions, as was expected due to the presence of sequels, authors that write multiple books with similar themes, etc. The relevant columns in this case were author, user id, and the time when the book was marked as read. Even though we had information about the time the book was shelved, we felt that that would be a weaker (albeit denser) signal, similar to adding an item to cart in the grocery world, and hence decided to go ahead with the time the book was marked 'read', which is analogous to a transaction. We again filtered the data in a fashion similar to the one described for the first dataset, and ended up with a 150,000 users and 11,400

authors. One thing to note is that the interactions in this dataset weren't as dense as the first dataset.

### B. Comparison methods

To compare our model, we chose the standard baselines in literature [28] : Popularity Model (Pop) and classic Matrix Factorization model (MF). The popularity model captures the popularity of attributes across each customer and recommends the most popular attributes for them.

The second model is a standard implicit OCCF Matrix Factorization. The observed transaction data is used to learn latent factors for the users and the attributes and to predict user-attribute interactions. This was done using the standard Alternating Least Squares (ALS) optimization.

We studied these under two settings : the first setting is a more realistic setting with the notion of *explore-exploit (EE)* built into the recommender system. Most real life recommenders employ a strategy to diversify their recommendations in the hope of increasing exposure of items which do not have much user interaction. The second setting removes the explore-exploit strategy from the two baselines to give a sterner test to our model. We also included a weighting to favor attributes that the user has prior interactions with. The baseline models and our model are compared across the metrics that are described in subsection IV-C.

### C. Evaluation metrics

The ground truth dataset was the list of brands bought by the users in the grocery dataset or the list of authors whose books were read by the readers in the Goodreads dataset, in the test window. The predictions from the model were a list of brands/authors, ordered by the probability that the given user would buy/read the given brand/author in the test window. We compare our model with the baseline models on five different evaluation metrics, most of them well-known in the collaborative filtering literature [30]–[34].

*1) NDCG@k:* As is known, DCG works on the idea that highly relevant entries appearing lower in the predictions list returned by the models should be penalized. In our case, the relevance for a brand/author is 1 if it appears in the top k predictions for a user and is present in the ground truth, and 0 otherwise. Ideal DCG (IDCG) is used to normalize this score to account for the varying lengths of the recommendation lists returned for different users. Finally, we take a mean of the NDCG values over all the queries, which are the users in the test set, to get a measure of the performance. In the following formulae, $rel_i$ denotes the relevance of the entry at the $i^{th}$ position in the predictions list returned by the models.

$$DCG_k = \sum_{i=1}^{k} \frac{rel_i}{log_2(i+1)}, NDCG_k = \frac{DCG_k}{IDCG_k}$$

*2) MAP@k:* The area under the precision-recall curve, which is obtained by plotting the precision and recall at every position in a ranked list of predictions, is called the average precision. Mean of the average precision scores over a set of queries i.e. users, gives the MAP. In the following formulae, AP is the average precision, MAP is the mean average precision, P(i) is the precision at position i, $\Delta r(i)$ is the change in recall from position i-1 to i, #rel is the total number of relevant brands/authors for that user (up to k), and $|U|$ is the number of users in the test set.

$$AP = \sum_{i=1}^{k} P(i)\Delta r(i) = \frac{\sum_{i=1}^{k} P(i) rel_i}{\#\text{rel}}, MAP = \frac{\sum_{j=1}^{|U|} AP_j}{|U|}$$

*3) Hit Rate@k:* Essentially the true positive rate, where a true positive is a brand/author predicted in the top k that is present in the ground truth for that query (user). We take a mean of the hit rate values over all the queries (users) and report that in section V.

$$\text{Hit rate} = \frac{\text{Number of True Positives}}{\text{Number of Positives}}$$

*4) MRR@k:* The reciprocal rank of a query response, i.e. predictions for a user, is the inverse of the position of the first item in the predictions list that is present in the ground truth for that query. Here, we consider only the first k predictions, and average the reciprocal ranks over all the users. In the following formula, $pos_i$ represents the position of the first prediction for the $i^{th}$ user that is present in the ground truth list for that user.

$$MRR = \frac{1}{|U|} \sum_{i=1}^{|U|} \frac{1}{pos_i}$$

*5) Limited AUC@k:* The ROC curve is a plot of the true positive rate against the false positive rate at various threshold values, and the general objective in recommender systems is to maximize the area under the ROC curve. But, in most such settings, the entries at the top of a list are more impactful than those at the bottom, but AUC is equally affected by swaps at different places in the returned list. To address this, we use limited AUC [35], which basically is the area under the part of the curve formed by the top k recommendations. This assumes that all the other relevant recommendations (apart from the top k) are distributed uniformly over the rest of the ranking list until all entries are retrieved. Thus, a straight line is drawn between the end point of the curve formed by these k recommendations and (1,1), the upper-right point of any ROC curve, and the area thus obtained is measured. This addresses some of the issues mentioned before, since swaps below the top k don't affect the AUC. This also has a few other good properties, such as a top-k list that contains more relevant entries will yield a higher AUC score, with the order mattering if the length of the list is close to the total number of brands/authors. We take a mean over all the queries (users) to get a mean LAUC.

## D. Implementation details

To generate the baselines, we used Turicreate [36], an open source toolkit for generating core machine learning models including recommenders. For the popularity recommender, we used the popularity recommender class and for the Matrix Factorization based model, we utilized the factorization recommender class. K-Fold cross validation was performed on both classes of models using the in-built capability to tune the models and finally the best models of each were selected for comparison.

For our model, we wrote a custom training loop and used Edward2 [37], [38] to do black-box variational inference [26]. Edward2 is a low-level language for specifying probabilistic models as programs and performing computations. We fed the models/distributions as functions whose inputs were the random variables that we were conditioning on and the outputs were the random variables that the probabilistic program was over. In the training loop, we first computed the log-likelihood using samples from the variational distribution. We used Edward's and TensorFlow's tracing functionalities (in steps 8 and 12, Algorithm 1) to record the model's computations for automatic differentiation. We then computed the KL divergence between the variational distribution and the prior distribution using the attributes of TensorFlow's distributions, and combined that with the log likelihood obtained from the posterior predictive function to get the ELBO. We tried different optimizers, learning schedules, and hyperparameter settings. A pseudocode of Edward's custom training loop adapted to our problem setting has been presented in Algorithm 1. This loop is called a certain number of times (to ensure convergence) for each batch in each epoch and the values of the variational parameters used to build the variational distribution (step 2, algorithm 1) are the updated values (step 13, algorithm 1) from the previous run.

---

**Algorithm 1** Variational Inference Training Loop

---

**INPUT:** Batch from Transaction matrix $T_b$, batch from Temporal Loyalty matrix $L_b$, Transaction matrix T, Temporal Loyalty matrix L, set of hyperparameters H, set of latent variables $\theta$, set of prior variables $\{u, v, bu, bv, \kappa_t, \psi_t, \kappa_l, \psi_l\}$

1: **procedure** CUSTOM TRAINING LOOP($T_b$, $L_b$)
2:      variational_family, trainable_parameters ← Build variational distribution
3:      $qu, qv, qbu, qbv, q\kappa_t, q\psi_t, q\kappa_l, q\psi_l$ ← Sample posterior variables from the variational_family
4:      $PP_T, PP_L$ ← Obtain posterior predictive functions, $P(T|\theta, H)$ and $P(L|\theta, H)$, from equation 3 by setting prior variables to the sample posterior values
5:      $LL_{T_b}, LL_{L_b}$ ← Compute the log likelihood of $T_b$ and $L_b$ from $PP_T$ and $PP_L$ respectively
6:      Initialize KL ← 0
7:      **for** prior_variable, variational_variable in [(u, qu), (v, qv), (bu, qbu), (bv, qbv), $(\kappa_t, q\kappa_t)$, $(\psi_t, q\psi_r)$,$(\kappa_s, q\kappa_s)$, $(\psi_s, q\psi_s)$] **do**
8:          KL ← KL + KL divergence between the distributions of the variational_variable and the prior_variable
9:      **end for**
10:      ELBO ← Compute ELBO using KL, $LL_{T_b}$, and $LL_{L_b}$ from equation 4
11:      Loss ← -ELBO
12:      Get the gradients using the loss and the trainable_parameters obtained
13:      Update the parameter values
14: **end procedure**

---

| Metric \ Method | | Pop + EE | MF + EE | Pop | MF | VI-MF |
|---|---|---|---|---|---|---|
| NDCG | @5 | 0.047 | 0.054 | 0.144 | 0.210 | 0.212 |
| | @10 | 0.031 | 0.036 | 0.096 | 0.140 | 0.141 |
| | @15 | 0.026 | 0.031 | 0.081 | 0.118 | 0.120 |
| | @20 | 0.025 | 0.028 | 0.077 | 0.112 | 0.114 |
| MAP | @5 | 0.016 | 0.021 | 0.049 | 0.098 | 0.099 |
| | @10 | 0.008 | 0.011 | 0.026 | 0.051 | 0.053 |
| | @15 | 0.006 | 0.008 | 0.020 | 0.040 | 0.040 |
| | @20 | 0.006 | 0.007 | 0.018 | 0.037 | 0.038 |
| HR | @5 | 0.064 | 0.064 | 0.197 | 0.196 | 0.198 |
| | @10 | 0.033 | 0.033 | 0.101 | 0.101 | 0.102 |
| | @15 | 0.024 | 0.025 | 0.075 | 0.075 | 0.076 |
| | @20 | 0.022 | 0.022 | 0.067 | 0.066 | 0.068 |
| MRR | @5 | 0.080 | 0.108 | 0.246 | 0.491 | 0.492 |
| | @10 | 0.080 | 0.108 | 0.246 | 0.491 | 0.492 |
| | @15 | 0.080 | 0.108 | 0.246 | 0.491 | 0.492 |
| | @20 | 0.080 | 0.108 | 0.246 | 0.491 | 0.492 |
| LAUC | @5 | 0.532 | 0.532 | 0.598 | 0.598 | 0.599 |
| | @10 | 0.516 | 0.516 | 0.551 | 0.551 | 0.552 |
| | @15 | 0.512 | 0.512 | 0.539 | 0.540 | 0.540 |
| | @20 | 0.511 | 0.511 | 0.536 | 0.537 | 0.537 |

Table I: Comparison of evaluation metrics across models on e-Commerce grocery data

| Metric \ Method | | Pop + EE | MF + EE | Pop | MF | VI-MF |
|---|---|---|---|---|---|---|
| NDCG | @5 | 0.020 | 0.027 | 0.047 | 0.068 | 0.069 |
| | @10 | 0.014 | 0.019 | 0.034 | 0.049 | 0.051 |
| | @15 | 0.012 | 0.016 | 0.030 | 0.043 | 0.044 |
| | @20 | 0.010 | 0.015 | 0.028 | 0.040 | 0.041 |
| MAP | @5 | 0.007 | 0.013 | 0.017 | 0.033 | 0.034 |
| | @10 | 0.004 | 0.008 | 0.011 | 0.021 | 0.022 |
| | @15 | 0.003 | 0.006 | 0.009 | 0.018 | 0.018 |
| | @20 | 0.002 | 0.005 | 0.008 | 0.016 | 0.017 |
| HR | @5 | 0.031 | 0.028 | 0.070 | 0.069 | 0.071 |
| | @10 | 0.018 | 0.017 | 0.041 | 0.041 | 0.042 |
| | @15 | 0.013 | 0.014 | 0.031 | 0.031 | 0.033 |
| | @20 | 0.011 | 0.012 | 0.026 | 0.027 | 0.028 |
| MRR | @5 | 0.033 | 0.061 | 0.075 | 0.148 | 0.150 |
| | @10 | 0.034 | 0.061 | 0.075 | 0.148 | 0.151 |
| | @15 | 0.034 | 0.061 | 0.075 | 0.149 | 0.151 |
| | @20 | 0.034 | 0.061 | 0.076 | 0.149 | 0.152 |
| LAUC | @5 | 0.514 | 0.512 | 0.533 | 0.533 | 0.534 |
| | @10 | 0.508 | 0.507 | 0.521 | 0.521 | 0.522 |
| | @15 | 0.506 | 0.505 | 0.517 | 0.517 | 0.518 |
| | @20 | 0.505 | 0.505 | 0.516 | 0.516 | 0.516 |

Table II: Comparison of evaluation metrics across models on Goodreads data

## V. RESULTS AND ANALYSIS

The results on the e-Commerce data and the open source Goodreads data have been presented in Table I and Table II respectively. The metrics for our model are shown in the final column, titled VI-MF (Variational Inference Matrix Factorization).

The first two baselines, with the explore-exploit strategy, suffer from trading off accuracy for diversity and hence do not perform as well as the other models. In both settings, with and without explore-exploit, MF-based models outperform the pop models because the pop models simply recommend attributes of the items that the user has bought most in the past whereas the latent factors capture user affinities well as they learn better representations from the interactions.

Our model shows a clear 1 to 3 percent increase in all metrics across the various ranks as compared to the best performing baseline model (that is, the classic Matrix Factorization) in both the e-Commerce grocery dataset as well as the Goodreads dataset. The size and scale of the datasets mean that these gains are significant. Quantitatively, in the e-Commerce setting, for a business with tens of billions of dollars in revenue, a 1 to 3 percent increase translates to hundreds of millions of dollars. This indicates that incorporating temporal loyalty leads to a better understanding of the user preferences, thus having an effect on the prediction of user behavior and subsequently revenue.

Overall, the metrics on the e-Commerce grocery data are higher than the ones on the Goodreads data. This can be explained by the higher density of grocery data leading to stronger user affinities to attributes. Interestingly, the trends across the models seem to hold across the domains of grocery and 'Fantasy and Paranormal' genre. In other words, the notion of brand loyalty in grocery seems similar to the notion of author loyalty in the 'Fantasy and Paranormal' genre of books.

## VI. CONCLUSION AND FUTURE WORK

In this work, we leverage a customer's temporal loyalty to an item attribute in addition to the engagement behavior to model their preferences and subsequently tackle the top-k attribute recommendation problem. We model this as an optimization problem over two matrices and use the Box's Loop framework and variational inference to estimate the parameter values and train the user embeddings that best explain the observed explicit and temporal signals. We demonstrate the effectiveness of the user embeddings learnt by showing that the proposed approach outperforms standard baselines for this task on a private e-Commerce grocery dataset as well as the publicly available Goodreads dataset, which also supports the hypothesis that capturing a customer's temporally changing interests can lead to better recommendations.

In terms of future directions, one could explore other ways to come up with the loyalty scores in the Temporal Loyalty matrix, L. Some works, such as [39], also focus on enriching the transaction matrix, T, to address issues that arise due to sparsity; we plan to investigate coupling those with our current approach. Another direction to explore would be to model the priors and the likelihoods with other distributions, informed by domain knowledge and the type of data one is dealing with.

## VII. Abbreviations and Acronyms

ALS: Alternating Least Squares
AP: Average Precision
DCG: Discounted Cumulative Gain
EE: Explore-Exploit
ELBO: Evidence Lower Bound
HR: Hit-Rate
IDCG: Ideal Discounted Cumulative Gain
KL: Kullback-Leibler Divergence
LAUC: Limited Area Under the Curve
MAP: Mean Average Precision
MF: Matrix Factorization
MCMC: Markov Chain Monte Carlo
MRR : Mean Reciprocal Rank
NDCG: Normalized Discounted Cumulative Gain
OCCF: One Class Collaborative Filtering
PML: Probabilistic Machine Learning
ROC: Receiver Operating Characteristic
SGD: Stochastic Gradient Descent
VAE: Variational Auto Encoder
VI : Variational Inference

## References

[1] R. Pan, Y. Zhou, B. Cao, N. N. Liu, R. Lukose, M. Scholz, and Q. Yang, "One-class collaborative filtering," in *Proceedings of the 2008 Eighth IEEE International Conference on Data Mining*, ICDM '08, (USA), p. 502–511, IEEE Computer Society, 2008.

[2] G. Lavee, N. Koenigstein, and O. Barkan, "When actions speak louder than clicks: A combined model of purchase probability and long-term customer satisfaction," in *Proceedings of the 13th ACM Conference on Recommender Systems*, RecSys '19, (New York, NY, USA), p. 287–295, Association for Computing Machinery, 2019.

[3] N. Lin, P. Tsai, Y. Chen, and H. H. Chen, "Music recommendation based on artist novelty and similarity," in *2014 IEEE 16th International Workshop on Multimedia Signal Processing (MMSP)*, pp. 1–6, 2014.

[4] V. Bogina and T. Kuflik, "Incorporating dwell time in session-based recommendations with recurrent neural networks.," 2017.

[5] X. Yi, L. Hong, E. Zhong, N. N. Liu, and S. Rajan, "Beyond clicks: Dwell time for personalization," in *Proceedings of the 8th ACM Conference on Recommender Systems*, RecSys '14, (New York, NY, USA), p. 113–120, Association for Computing Machinery, 2014.

[6] P. Yin, P. Luo, W.-C. Lee, and M. Wang, "Silence is also evidence: Interpreting dwell time for recommendation from psychological perspective," in *Proceedings of the 19th ACM SIGKDD International Conference on Knowledge Discovery and Data Mining*, KDD '13, (New York, NY, USA), p. 989–997, Association for Computing Machinery, 2013.

[7] R. Ronen, E. Yom-Tov, and G. Lavee, "Recommendations meet web browsing: enhancing collaborative filtering using internet browsing logs," pp. 1230–1238, 05 2016.

[8] F. Aiolli, "Efficient top-n recommendation for very large scale binary rated datasets," 10 2013.

[9] M. Deshpande and G. Karypis, "Item-based top-n recommendation algorithms," *ACM Trans. Inf. Syst.*, vol. 22, p. 143–177, Jan. 2004.

[10] R. Gemulla, E. Nijkamp, P. J. Haas, and Y. 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*, KDD '11, (New York, NY, USA), p. 69–77, Association for Computing Machinery, 2011.

[11] R. Bhagat, S. Muralidharan, A. Lobzhanidze, and S. Vishwanath, "Buy it again: Modeling repeat purchase recommendations," in *Proceedings of the 24th ACM SIGKDD International Conference on Knowledge Discovery & Data Mining*, KDD '18, (New York, NY, USA), p. 62–70, Association for Computing Machinery, 2018.

[12] Q. Zhao, J. Chen, M. Chen, S. Jain, A. Beutel, F. Belletti, and E. H. Chi, "Categorical-attributes-based item classification for recommender systems," in *Proceedings of the 12th ACM Conference on Recommender Systems*, RecSys '18, (New York, NY, USA), p. 320–328, Association for Computing Machinery, 2018.

[13] W.-H. Chen, C.-C. Hsu, Y.-A. Lai, V. Liu, M.-Y. Yeh, and S.-D. Lin, "Attribute-aware recommender system based on collaborative filtering: Survey and classification," *Frontiers in Big Data*, vol. 2, p. 49, 2020.

[14] D. Liang, R. Krishnan, M. Hoffman, and T. Jebara, "Variational autoencoders for collaborative filtering," 02 2018.

[15] K. Zheng, X. Yang, Y. Wang, Y. Wu, and X. Zheng, "Collaborative filtering recommendation algorithm based on variational inference," *International Journal of Crowd Science*, vol. ahead-of-print, 01 2020.

[16] Y.-D. Kim and S. Choi, "Scalable Variational Bayesian Matrix Factorization with Side Information," vol. 33 of *Proceedings of Machine Learning Research*, (Reykjavik, Iceland), pp. 493–502, PMLR, 22–25 Apr 2014.

[17] G. Chen, F. Zhu, and P. A. Heng, "Large-scale bayesian probabilistic matrix factorization with memo-free distributed variational inference," *ACM Trans. Knowl. Discov. Data*, vol. 12, Jan. 2018.

[18] D. M. Blei, "Build, compute, critique, repeat: Data analysis with latent variable models," *Annual Review of Statistics and Its Application*, vol. 1, no. 1, pp. 203–232, 2014.

[19] D. Song, C. E. Lee, Y. Li, and D. Shah, "Blind regression: Nonparametric regression for latent variable models via collaborative filtering," in *Advances in Neural Information Processing Systems 29* (D. D. Lee, M. Sugiyama, U. V. Luxburg, I. Guyon, and R. Garnett, eds.), pp. 2155–2163, Curran Associates, Inc., 2016.

[20] M. Harvey, M. J. Carman, I. Ruthven, and F. Crestani, "Bayesian latent variable models for collaborative item rating prediction," in *Proceedings of the 20th ACM International Conference on Information and Knowledge Management*, CIKM '11, (New York, NY, USA), p. 699–708, Association for Computing Machinery, 2011.

[21] X. He, L. Liao, H. Zhang, L. Nie, X. Hu, and T.-S. Chua, "Neural collaborative filtering," in *Proceedings of the 26th International Conference on World Wide Web*, WWW '17, (Republic and Canton of Geneva, CHE), p. 173–182, International World Wide Web Conferences Steering Committee, 2017.

[22] W. Chen, F. Cai, H. Chen, and M. D. Rijke, "Joint neural collaborative filtering for recommender systems," *ACM Trans. Inf. Syst.*, vol. 37, Aug. 2019.

[23] D. M. Blei, "Probabilistic topic models," *Commun. ACM*, vol. 55, p. 77–84, Apr. 2012.

[24] D. M. Blei, A. Kucukelbir, and J. D. McAuliffe, "Variational inference: A review for statisticians," *Journal of the American Statistical Association*, vol. 112, no. 518, pp. 859–877, 2017.

[25] M. Braun and J. McAuliffe, "Variational inference for large-scale models of discrete choice," *Journal of the American Statistical Association*, vol. 105, p. 324–335, Mar 2010.

[26] R. Ranganath, S. Gerrish, and D. Blei, "Black box variational inference," in *AISTATS*, 2014.

[27] J. Paisley, D. Blei, and M. Jordan, "Variational bayesian inference with stochastic search," 2012.

[28] M. Wan and J. J. McAuley, "Item recommendation on monotonic behavior chains," in *Proceedings of the 12th ACM Conference on Recommender Systems, RecSys 2018, Vancouver, BC, Canada, October 2-7, 2018* (S. Pera, M. D. Ekstrand, X. Amatriain, and J. O'Donovan, eds.), pp. 86–94, ACM, 2018.

[29] M. Wan, R. Misra, N. Nakashole, and J. J. McAuley, "Fine-grained spoiler detection from large-scale review corpora," in *Proceedings of the 57th Conference of the Association for Computational Linguistics, ACL 2019, Florence, Italy, July 28- August 2, 2019, Volume 1: Long Papers* (A. Korhonen, D. R. Traum, and L. Màrquez, eds.), pp. 2605–2610, Association for Computational Linguistics, 2019.

[30] M. Zhu, "Recall, precision and average precision," 09 2004.

[31] A. Tharwat, "Classification assessment methods," *Applied Computing and Informatics*, 2018.

[32] O. Vechtomova, "Introduction to information retrieval christopher d. manning, prabhakar raghavan, and hinrich schütze (stanford university, yahoo! research, and university of stuttgart) cambridge: Cambridge university press, 2008, xxi+ 482 pp; hardbound, isbn 978-0-521-86571-5, $60.00," 2009.

[33] K. Järvelin and J. Kekäläinen, "Ir evaluation methods for retrieving highly relevant documents," in *ACM SIGIR Forum*, vol. 51, pp. 243–250, ACM New York, NY, USA, 2017.

[34] C. Lioma, J. G. Simonsen, and B. Larsen, "Evaluation measures for relevance and credibility in ranked lists," in *Proceedings of the ACM SIGIR International Conference on Theory of Information Retrieval*, ICTIR '17, (New York, NY, USA), p. 91–98, Association for Computing Machinery, 2017.

[35] G. Schröder, M. Thiele, and W. Lehner, "Setting goals and choosing metrics for recommender system evaluations," vol. 811, 01 2011.

[36] Apple, *Turicreate (https://github.com/apple/turicreate)*, 2014 (accessed May 11, 2020).

[37] D. Tran, A. Kucukelbir, A. B. Dieng, M. Rudolph, D. Liang, and D. M. Blei, "Edward: A library for probabilistic modeling, inference, and criticism," *arXiv preprint arXiv:1610.09787*, 2016.

[38] D. Tran, M. D. Hoffman, R. A. Saurous, E. Brevdo, K. Murphy, and D. M. Blei, "Deep probabilistic programming," in *International Conference on Learning Representations*, 2017.

[39] Y. He, H. Chen, Z. Zhu, and J. Caverlee, "Pseudo-implicit feedback for alleviating data sparsity in top-k recommendation," in *2018 IEEE International Conference on Data Mining (ICDM)*, pp. 1025–1030, 2018.