
Matrix Factorization in Action

Herman Franklin Tesso Tassang¹

Abstract

Matching consumers with the most suitable products is crucial for enhancing user satisfaction and loyalty on various online platforms, such as e-commerce websites, streaming services, and social media. Recommender systems solve this problem by analyzing user interest patterns in products to provide users with personalized content and services. In this project, we will develop our own prototype of a recommender system using the Matrix Factorization technique and the MovieLens dataset with explicit user feedback.

1. Introduction

During the last few decades, the demand for recommender(Lü et al., 2012) systems has increased significantly with the rise of platforms like YouTube, Amazon, Netflix, and many others(Koenigstein et al., 2012). Recommender systems are algorithms that filter out relevant information from a large amount of dynamically generated data, based on user preferences, interests, or observed behaviors.

The purpose of recommender systems is to suggest relevant items to users. There are two major paradigms for achieving this task. The first is the "content-based" approach(Schafer et al., 2007a), which creates profiles for each user or product and builds a model based on the available "features" that explain the observed between users and items. The program then uses these profiles to associate users with products.

The second approach is called "collaborative filtering,"(Schafer et al., 2007b) which is solely based on past interactions recorded between users and items. The recommendations are stored in the so called "user-item interactions matrix". The idea behind this approach is that past user-item interactions are sufficient to detect similar users and/or items and make predictions based on this estimated proximity. There are two primary classes of collaborative filtering: **memory-based** methods and **model-based methods**.

Memory-based methods work directly with the values of recorded interactions without assuming a model. They essentially rely on nearest neighbor search(e.g find closest users from user of interest and suggest the most popular

items among these neighbours). On the other hand, model-based methods such as matrix factorization assume an underlying "generative" model that explains the user-item interactions and try to discover this model in order to make new predictions.

Some systems prompt users to provide ratings for items through the system interface, in order to construct and improve their models. These ratings are referred to as "explicit user feedback." Because it requires effort from users, explicit feedback usually comprises a sparse matrix, as each user is likely to have rated only a small percentage of possible items.

In this work, we will focus on building a recommender system using the matrix factorization approach, specifically with highly explicit feedback (ratings). The first section will discuss related work on recommender systems. The next two sections will describe the data used and the theoretical basis of matrix factorization(Koren et al., 2009), as well as the modeling process. The following two sections will be dedicated to the training, validation procedure, and recommendation process. Finally, we will discuss how to modify the model to cope with the "cold start" problem.

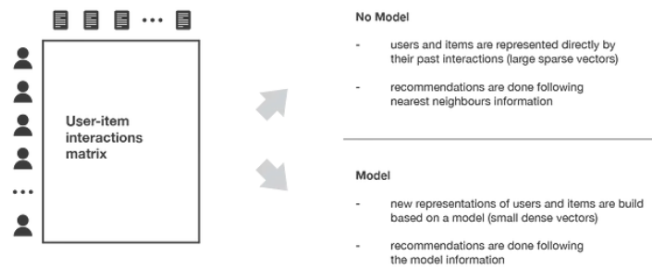


Figure 1. Overview of the Collaborative filtering paradigm

2. Data

2.1. Data Description

We will be using the [MoviesLens](#) datasets in our work. These datasets are widely recognized collections of movie ratings data provided by the [GroupLens](#) (University of Minnesota) . They are commonly used for research in recommendation systems and contain user ratings, movie meta-data, and demographic information. This data is valuable for

analyzing user preferences and developing recommendation algorithms. We will primarily focus on the 25M MovieLens dataset, which has the following characteristics:

- The dataset contains 25 million 5-star ratings.
- There are $M=162,541$ users and $N=59,047$ movies.
- Every user has rated at least 20 movies.
- each user has rated at least 20 movies.
- The dataset includes information about the movies, such as ID, title, and genres (19 genres are listed).
- There is no demographic information available.

2.2. Overall Trends in the Data

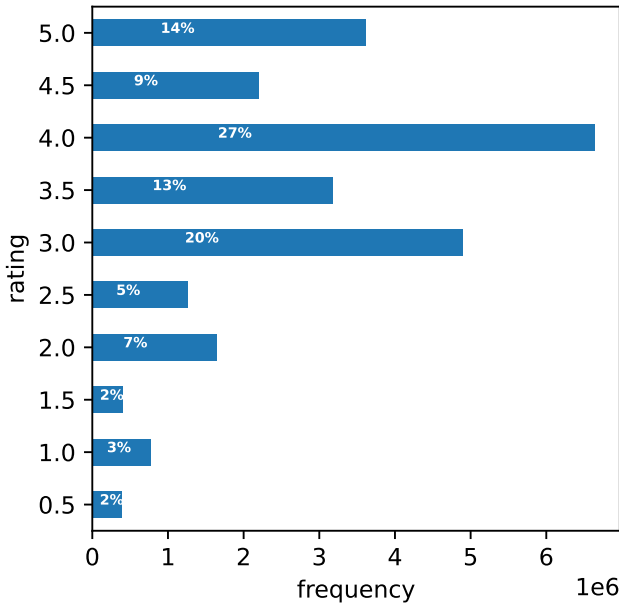


Figure 2. Rating distributions of the data

Figure 2 displays the rating distribution in the 25M MovieLens dataset. A noticeable pattern emerges where the majority of movies receive higher ratings. More specifically, ratings commonly fall within the range of 3 to 4 stars, suggesting that users generally rate movies positively. However, there is a long tail of movies with fewer ratings, encompassing both highly-rated and poorly-rated films. These observations underscore the importance of considering both movie popularity and individual user preferences when constructing our recommender system.

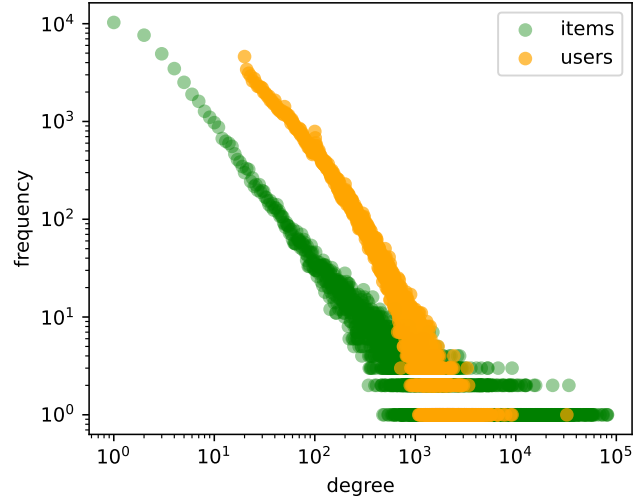


Figure 3. Rating distributions given movies/users

Figure 3 shows the distribution of ratings for movies and users. This distribution clearly follows a power-law pattern. It is observed that a small number of movies receive a disproportionately larger number of ratings compared to the majority of movies. Similarly, a small number of users have rated a large number of movies. This plot further emphasizes the significance of considering popularity bias in the following analysis.

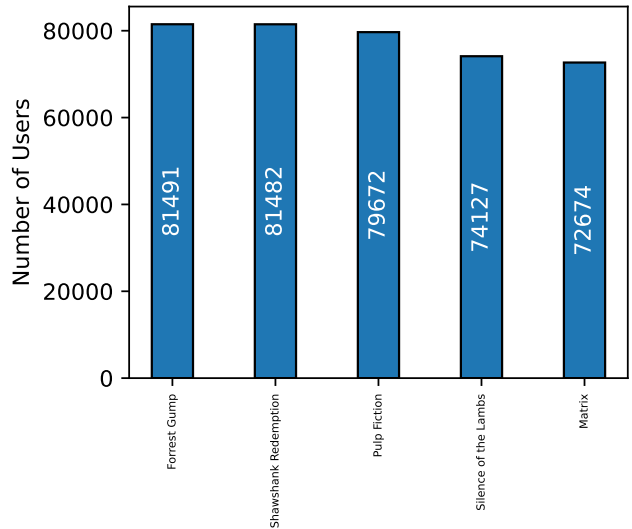


Figure 4. top popular movies

The last figure, **Figure 4**, represents the five most popular movies in our dataset. It shows the frequency of user ratings for each movie.

2.3. Data Structure

We used the dataset to create the necessary data structures for building our recommender model.

First, we mapped each user and item ID to an index, denoted by $1 \leq m \leq M$ and $1 \leq n \leq N$ respectively. We also stored the rating given by each user (indexed by m) for each movie (indexed by n) as " r_{mn} ". Finally, we created the following tree data structures.

The main list consists of multiple lists, each containing tuples. Each list within the main list is indexed by m , where m is a value between 1 and M . These lists contain the indexes of movies that have been rated by the user indexed by m , along with their corresponding ratings in the form of (item, rating) pairs.

samples list : list of tuples of the form (m, n, r_{mn}) , *user - item - rating*.

data index by user: a list of lists. The main list consists of multiple lists, each containing tuples. Each list within the main list is indexed by m ($1 \leq m \leq M$). These lists contain the indexes of movies that have been rated by the user indexed by m , along with their corresponding ratings in the form of (item, rating) pairs.

data index by movie: a list of lists.

The main list consists of a series of sublists, each containing tuples. Each sublist is indexed by n ($1 \leq n \leq N$) and contains the indices of the users who have rated the movie indexed by n , along with their corresponding ratings as pairs (*user, rating*).

3. Model Building

3.1. Matrix Factorization

Matrix factorization is a model-based collaborative filtering approach that relies only on *user-item interaction* information. It assume a latent model supposed to explains these interactions. The algorithm consist in decomposing the huge and sparse user-item interaction matrix into a product of two smaller and dense matrices: a "*user-factor matrix*" (containing user representations) that multiplies a *factor-item matrix* (containing item representations).

The main assumption is that there is a low-dimensional latent space of features in which we can represent both users and items. In this space, the interaction between a user and an item can be obtained by computing the dot product of their corresponding dense vectors. However, we do not explicitly provide these features to our model, as can be done in content-based approaches (1). Instead, we allow the system to discover these useful features on its own and

create its own representations of both users and items.

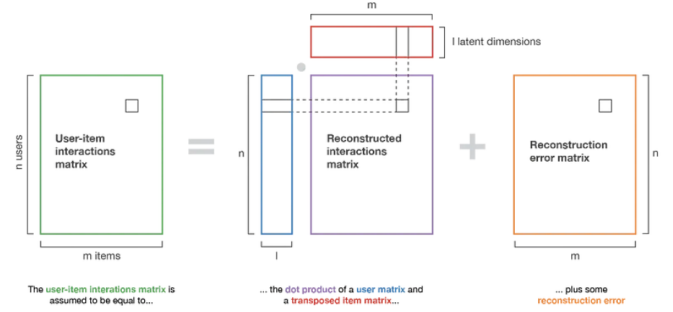


Figure 5. Graphic representation : Matrix factorization

3.2. Mathematical Formulation

Let's consider an interaction matrix $R \in \mathbb{R}^{M \times N}$ of ratings where only some items have been rated by each user (most interactions are set to *None* to express the lack of rating). We want to factorize that matrix such that:

$$R \approx UV^T \quad (1)$$

However, since most of the observed variation in rating values is due to effects associated with either users or items known as biases independent of any actions; it would be unwise to explain the full rating value by only interaction of the form UV^T . Instead, equation (1) is reformulated as:

$$R \approx UV^T + b^{(u)} + b^{(i)} \quad (2)$$

where:

- $U \in \mathbb{R}^{M \times K}$ is the "user matrix" whose rows represent the M users.
- $V \in \mathbb{R}^{N \times K}$ is the "item matrix" whose row represent the N items.
- $user_m \equiv u_m$ user latent vector $\forall m \in \{1, \dots, M\}$.
- $item_n \equiv v_n$ item latent vector $\forall n \in \{1, \dots, N\}$.
- K is the dimension of the latent space
- $b^{(u)}$ and $b^{(i)}$ are user and vector biases respectively.

this lead to the estimate:

$$\hat{r}_{mn} = u_m^T v_n + b_m^{(u)} + b_n^{(i)} \quad (3)$$

Therefore, we model it as:

$$\mathcal{P}(r_{mn} | u_m, v_n, b_m^{(u)}, b_n^{(i)}) = \mathcal{N}(r_{mn}; u_m^T v_n + b_m^{(u)} + b_n^{(i)}, \lambda^{-1}) \quad (4)$$

with:

- λ the precision factor
- $\mathcal{P}(u_m) = \mathcal{N}(u_m; 0, \tau^{-1}\mathbf{I})$, $\mathcal{P}(v_n) = \mathcal{N}(v_n; 0, \tau^{-1}\mathbf{I})$ priors for u_m and v_n respectively.
- $\mathcal{P}(b^{(u)}) = \mathcal{N}(b^{(u)}; 0, \gamma^{-1}\mathbf{I})$, $\mathcal{P}(b^{(i)}) = \mathcal{N}(b^{(i)}; 0, \gamma^{-1}\mathbf{I})$ priors for $b^{(u)}$ and $b^{(i)}$ respectively.

The objective function:

We want to find $(U, V, b^{(u)}, b^{(i)})$ that maximize the regularized log likelihood:

$$\begin{aligned} \mathcal{L} = & \log \mathcal{P}(R|U, V, b^{(u)}, b^{(i)}) + \log \mathcal{P}(U) \\ & + \log \mathcal{P}(V) + \log \mathcal{P}(b^{(u)}) + \log \mathcal{P}(b^{(i)}) \\ = & -\frac{\lambda}{2} \sum_m \sum_{n \in \Omega(m)} \left(r_{mn} - (u_m^T v_n + b_m^{(u)} + b_n^{(i)}) \right)^2 \\ & -\frac{\tau}{2} \sum_m u_m^T u_m - \frac{\tau}{2} \sum_n v_n^T v_n \\ & -\frac{\gamma}{2} \sum_m b_m^{(u)2} - \frac{\gamma}{2} \sum_n b_n^{(i)2} \end{aligned} \quad (5)$$

Thus,

$$(U, V, b^{(u)}, b^{(i)}) = \underset{U, V, b^{(u)}, b^{(i)}}{\operatorname{argmin}} \{ -\mathcal{L}(U, V, b^{(u)}, b^{(i)}) \} \quad (6)$$

3.3. Learning Algorithm

The algorithm used to minimize equation (6) is called alternating least squares (ALS). The underlying idea is that even though equation (6) is not convex, we can solve the optimization problem optimally by fixing one unknown at a time.

Specifically, for a fixed value of m , we compute the partial derivatives $\frac{\partial \mathcal{L}}{\partial u_m}$, and $\frac{\partial \mathcal{L}}{\partial b_m^{(u)}}$ and set them to zero. This allows us to solve for the optimal value (vector) of the quadratic problem. As a result, we obtain the following updates:

User vector update:

$$u_m = \lambda \left(\sum_{n \in \Omega(m)} v_n v_n^T + \tau \mathbf{I} \right)^{-1} \left(\lambda \sum_{n \in \Omega(m)} v_n (r_{mn} - b_m^{(u)} - b_n^{(i)}) \right) \quad (7)$$

User bias update:

$$b_m^{(u)} = \frac{\lambda \sum_{n \in \Omega(m)} (r_{mn} - (u_m^T v_n + b_n^{(i)}))}{\lambda \Omega(m) + \gamma} \quad (8)$$

as previously, with fixed value of n , by computing the partial derivatives $\frac{\partial \mathcal{L}}{\partial v_n}$, and $\frac{\partial \mathcal{L}}{\partial b_n^{(i)}}$ and setting them to zero,

we obtain the following updates for the items:

Item vector update:

$$v_n = \lambda \left(\sum_{m \in \Omega(n)} u_m u_m^T + \tau \mathbf{I} \right)^{-1} \left(\lambda \sum_{m \in \Omega(n)} u_m (r_{mn} - b_m^{(u)} - b_n^{(i)}) \right) \quad (9)$$

Item bias update:

$$b_n^{(i)} = \frac{\lambda \sum_{m \in \Omega(n)} (r_{mn} - (u_m^T v_n + b_m^{(u)}))}{\lambda \Omega(n) + \gamma} \quad (10)$$

Note: $\Omega(m)$ represents the set of movie indexes that has been rated by the user indexed by m . Similarly, $\Omega(n)$ is the set of user indexes that have rated the movie indexed by n .

Algorithm 1 Alternating Least Square (ALS)

Input: data-structured, K, λ, τ, γ

Initialize $U, V, b^{(u)}, b^{(i)}$

repeat

for $m = 0$ **to** $M - 1$ **do**

for (n, r) in data-index-by-user[m] **do**

 update user vector $U[m :]$ with equation (7)

 update user bias $b_m^{(u)}$ with equation (8)

end for

end for

for $n = 0$ **to** $N - 1$ **do**

for (m, r) in data-index-by-item[n] **do**

 update item vector $V[n :]$ with equation (9)

 update item bias $b_n^{(i)}$ with equation (10)

end for

end for

until convergence

4. Training and validation

4.1. experiments

In this subsection, we will experiment the optimization algorithm using a smaller MoviesLens dataset : the 100K MoviesLens dataset consist of 100,000 ratings and 9000 movies rated by 600 users.

This step is perform in order to achieve the following goals:

- 1) Quickly compare the reliability of the model that only involves biases to explain the ratings $\hat{r}_{mn} = b_m^{(u)} + b_n^{(i)}$ with our original model which also includes the latent trait vectors (equation (3)).
- 2) Have and overview of how hyperparameters K, γ, λ and τ can affect the overall performance of the model.

The performance of our models is assessed based on the **RMSE**:

$$\sqrt{\frac{1}{L} \sum_{m=1}^M \sum_{n \in \Omega(m)} \left(r_{mn} - (u_m^T v_n + b_m^{(u)} + b_n^{(i)}) \right)^2}$$

with, L the total number of samples.

Experiment procedure: run the **ALS** algorithm (equation 1 with varying values of hyperparameters, record both the *Loss* and *RMSE* values at each iteration and use them to create plots representing the training evolution.

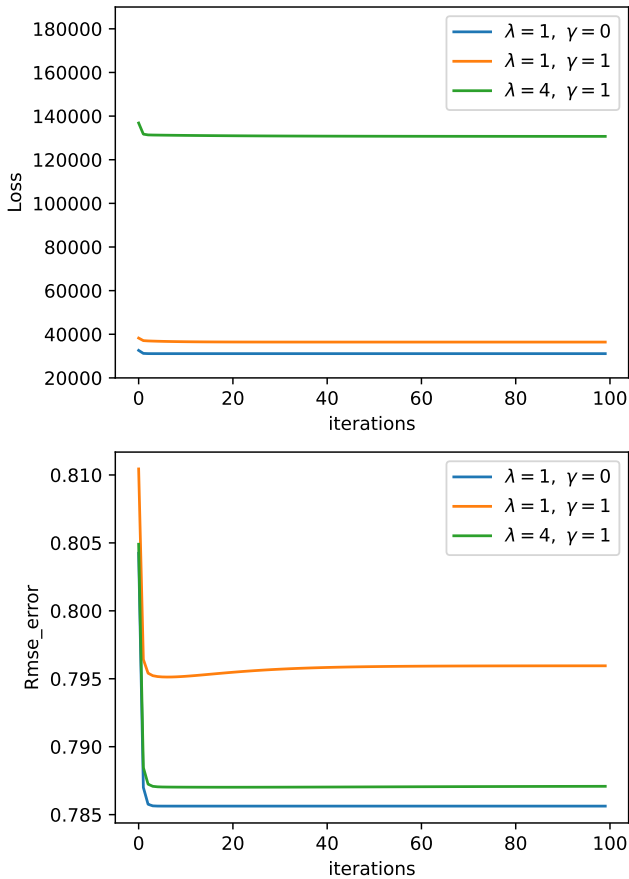


Figure 6. Training history for user+item biases update only on the 100K dataset, $\hat{r}_{mn} = b_m^{(u)} + b_n^{(i)}$.

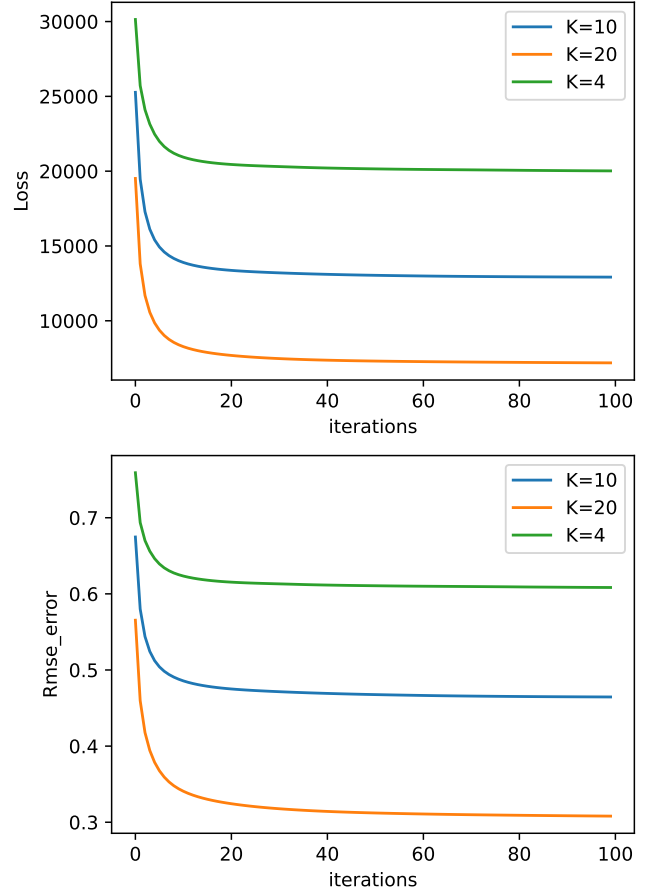


Figure 7. Training history for user+item biases+ latent vectors (equation (3)) on the 100K dataset.

Discussion:

Based on this small experiment, it appears that the models that incorporate latent factors show improvement compared to the models that rely solely on user and item biases. Different values of hyperparameters lead to significant differences in the performance of the models, as measured by RMSE and overall loss values. Specifically, it seems that for the model that includes user interaction (**Figure 7**), increasing the dimension of the latent space K results in a more flexible model. These observations emphasize the importance of carefully selecting the hyperparameters for training on the full 25M ratings dataset in order to achieve an effective model.

4.2. Training on the 25M MoviesLens Dataset.

Before training the model on the entire 25M dataset, it is necessary to determine which hyperparameters work well,

as discussed previously. These hyperparameters should allow for achieving good performance without overtraining the data.

To accomplish this, the data will be split into a training set and a test set. The training set will be used for the fitting procedure, optimizing the parameters, while the test set will be used to check for overfitting by visualizing the recorded train/test loss and RMSE over training iterations.

4.3. Splitting Strategy

To split the data, all the samples have been shuffled before being divided into 80% for the training set and 20% for the test set. The data structures, as described previously (see 2.3) have been reconstructed from the training set. Consequently, both the user and item data structures might contain empty lists, but that's also fine.

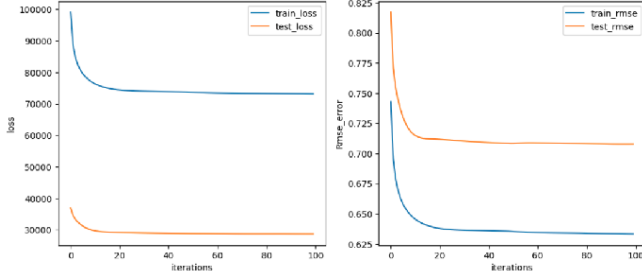


Figure 8. Training/test loss and RMSE history on the full 25 million ratings

4.4. Choose of hyperparameters

Several models have been trained using various values of K , λ , γ and τ . To evaluate each model, we plot the loss and RMSE history on the test set in order to determine if there is any overfitting occurring.

Figure 8 shows the obtained results of some trials.

After conducting multiple trials, we observed that overfitting occurs when K is set to 20 or higher. Additionally, we discovered that models with a latent dimension of $K = 10$ generally perform well, as they converge to lower test RMSE values, as long as reasonable values of K , λ , γ and τ are used. Based on these findings, we have selected the following hyperparameters for the entire 25M dataset.

$$K = 10, \lambda = 1, \gamma = 0.1, \tau = 1$$

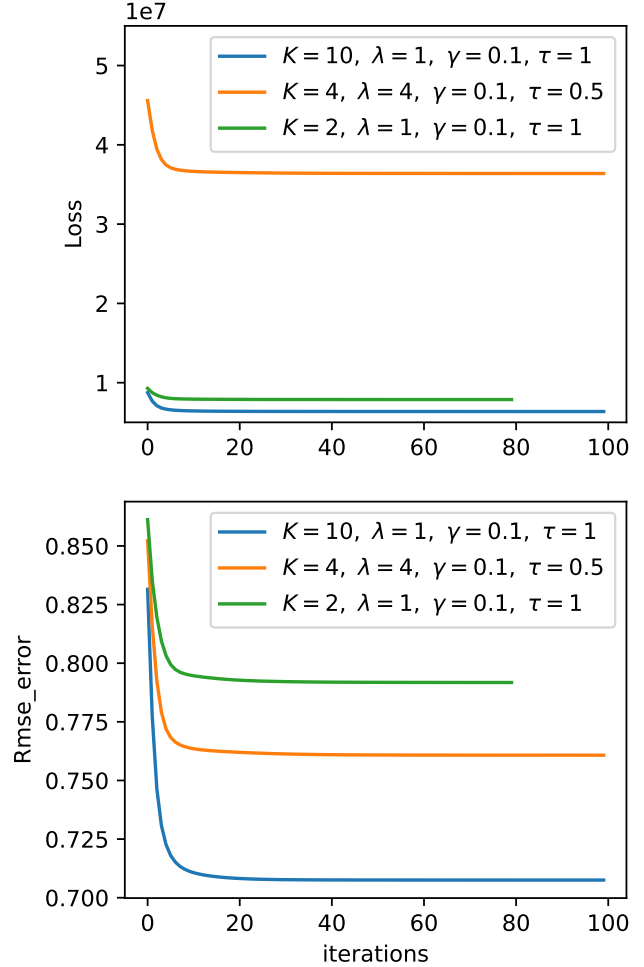


Figure 9. Errors on the test set for different values of the hyperparameters

5. Recommendations

The optimization step for all 25 million ratings has been completed. As a result, we now have access to the learning representations of both the users and items i.e optimized latent vectors as well as the optimized users and items biases. However, before we proceed, let's take a look into the latent space.

5.1. Item embeddings

In this subsection, we consider the model characterised by,

$$K = 2, \lambda = 1, \gamma = 0.1, \tau = 1, (*) \text{ (figure(9))}$$

whose the dimension of the latent space is 2. The goal is to be able to visualize the item trait vectors v_n in this latent space. We anticipate that similar movies will be embedded together and that opposite movies will be embedded in opposite parts of the space.

The plot **fig9** show the embeddings of "Children's" and "Horror" genre movies.

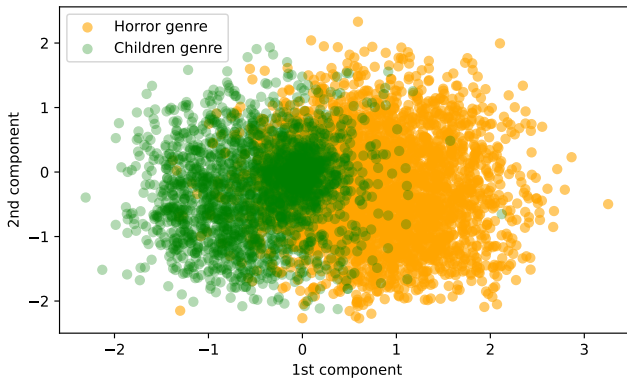


Figure 10. 2D item (movies) embeddings

5.2. Polirizing Movies

In recommender system, polarizing movies are those that generate extreme ratings from users. These movies receive a significant number of very high ratings from some users and very low ratings from others, with few ratings falling in between. Some users will rate them very highly and want similar recommendations, while others will rate them very poorly and want the opposite. Consequently, they quickly narrow down a user's taste once you know whether they like or dislike them.

It turn out that identifying polarizing movies amounts to search for movies with unusually long vectors. Mathematically this involves extracting movies whose latent vectors have the highest Euclidean norm values $\|v_n\|$.

Remember that the dimension of the latent space we used for training on the full 25 million ratings is $K = 10$ (see (*)). First, the movie data has been filtered by removing those that have received less than 50 ratings. Then, the five most polarizing movies have been extracted based on this rule.

title	genres
Twilight (2008)	Drama Fantasy Romance Thriller
Twilight Saga: New Moon, The (2009)	Drama Fantasy Horror Romance Thriller
Twilight Saga: Eclipse, The (2010)	Fantasy Romance Thriller IMAX
Twilight Saga: Breaking Dawn - Part 1, The (2011)	Adventure Drama Fantasy Romance
Twilight Saga: Breaking Dawn - Part 2, The (2012)	Adventure Drama Fantasy Romance IMAX

Figure 11. Top most polarizing movies in the dataset

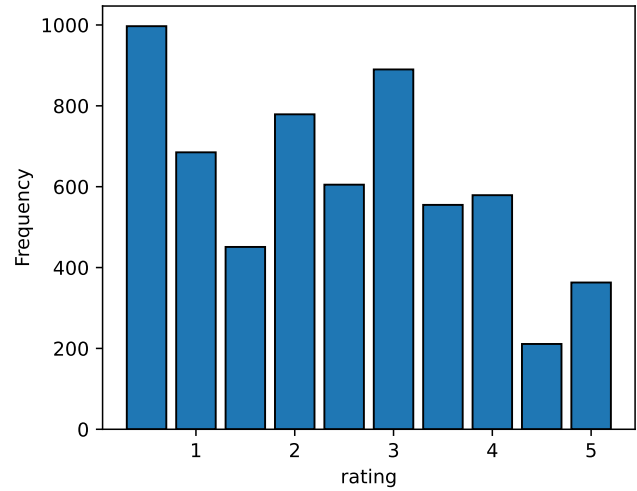


Figure 12. Rating distributions of the most polarizing movie: "Twilight Saga: New Moon"

The less polarizing movies are those that tend to have high popularity indicating that they are well-liked by a large number of users. The followings are the 5 less polarizing movies in the dataset:

title	genres
Pulp Fiction (1994)	Comedy Crime Drama Thriller
Shawshank Redemption, The (1994)	Crime Drama
Forrest Gump (1994)	Comedy Drama Romance War
Silence of the Lambs, The (1991)	Crime Horror Thriller
Matrix, The (1999)	Action Sci-Fi Thriller

Figure 13. Top less polarizing movies

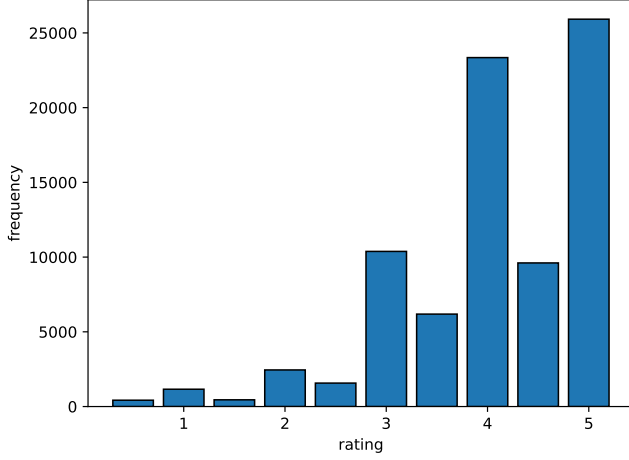


Figure 14. Rating distributions of the less polarizing movie: "Forrest Gump"

5.3. Recommendation Process

We designed this model with the specific goal of providing personalized content to users based on their ratings of certain movies. The recommendations are generated based on a kind of *user recommendation score* that is not influenced by the user's bias.

Workflow for movies recommendations to a user:

- 1) Learn the latent vector of this user using *ALS* algorithm (see 1) and the optimized item matrix V of the underlying model.
- 2) Compute the user's score for the movie "n" according to the formula:

$$Score_{item_n} = u_{user}^T v_n + \alpha \times b_n^{(i)}, \quad \forall n \in \{1, \dots, N\}$$

where:

- u_{user} is the learned user latent vector of the user
 - α is a factor use to downplay the contribution of item bias to cope with the movie popularity and/or high bias effect.
- 3) Finally, the scores are ranked and the movies associated with highest score sit at the top of the recommendations.

Figure (15) shows the recommendations for a "dummy user" that gave the movie "Lord the rings" five stars. The factor α have been choose after a several trials and the value $\alpha = 0.02$ have been shown to be an acceptable value.

title	genres
Star Wars: Episode IV - A New Hope (1977)	Action Adventure Sci-Fi
Star Wars: Episode V - The Empire Strikes Back...	Action Adventure Sci-Fi
Star Wars: Episode VI - Return of the Jedi (1983)	Action Adventure Sci-Fi
Star Wars: Episode I - The Phantom Menace (1999)	Action Adventure Sci-Fi
Lord of the Rings: The Fellowship of the Ring...	Adventure Fantasy
Star Wars: Episode II - Attack of the Clones (...)	Action Adventure Sci-Fi IMAX
Lord of the Rings: The Two Towers, The (2002)	Adventure Fantasy
Lord of the Rings: The Return of the King, The...	Action Adventure Drama Fantasy
Harry Potter and the Prisoner of Azkaban (2004)	Adventure Fantasy IMAX
Star Wars: Episode III - Revenge of the Sith (...)	Action Adventure Sci-Fi
Harry Potter and the Goblet of Fire (2005)	Adventure Fantasy Thriller IMAX
Harry Potter and the Order of the Phoenix (2007)	Adventure Drama Fantasy IMAX
Star Wars: The Clone Wars (2008)	Action Adventure Animation Sci-Fi
Harry Potter and the Half-Blood Prince (2009)	Adventure Fantasy Mystery Romance IMAX
Harry Potter and the Deathly Hallows: Part 1 (...)	Action Adventure Fantasy IMAX
Harry Potter and the Deathly Hallows: Part 2 (...)	Action Adventure Drama Fantasy Mystery IMAX
Hobbit: An Unexpected Journey, The (2012)	Adventure Fantasy IMAX
Hobbit: The Desolation of Smaug, The (2013)	Adventure Fantasy IMAX
The Hobbit: The Battle of the Five Armies (2014)	Adventure Fantasy
Star Wars: Episode VII - The Force Awakens (2015)	Action Adventure Fantasy Sci-Fi IMAX

Figure 15. Personalized recommendations provide by the recommender prototype: Top twenty.

6. Add Features information

As it only consider past interactions to make recommendations, collaborative filtering approaches like matrix factorization suffer from the "**cold start problem**": it is impossible to recommend anything to new users or to recommend a new item to any users.

This drawback can be addressed in different way:

- recommend new items (movies) to random users
- recommend new items to must active users
- recommend new items to set of various users

However, a more sophisticated approach has been utilized in this study to address the issue with new movies. This approach involves incorporating side information about the movies into the model. The concept behind this is that when new items are introduced, the model can still provide informed recommendations based on the known features of the items, even without interaction data. In the 25M MovieLens dataset, the features information includes movie genres (*Action, Horror, Children's...etc*). Therefore, we include genre information in the model, allowing the algorithm (ALS) to learn its own representation of the features (feature vectors).

6.1. Mathematical Formulation

The model (4)

$$\mathcal{P}(r_{mn}|u_m, v_n, b_m^{(u)}, b_n^{(i)}, f_i) = \mathcal{N}(r_{mn}; u_m^T v_n + b_m^{(u)} + b_n^{(i)}, \lambda^{-1}) \quad (11)$$

with:

- $\mathcal{P}(u_m) = \mathcal{N}(u_m; 0, \tau \mathbf{I})$, prior for u_m .
- $\mathcal{P}(f_i) = \mathcal{N}(f_i; 0, \beta \mathbf{I})$, prior for f_i
- $\mathcal{P}(v_n) = \mathcal{N}(v_n; \frac{1}{\sqrt{F_n}} \sum_{i \in \text{features}(n)} f_i, \tau \mathbf{I})$, prior for v_n .
- $\mathcal{P}(b^{(u)}) = \mathcal{N}(b^{(u)}; 0, \gamma \mathbf{I})$, $\mathcal{P}(b^{(i)}) = \mathcal{N}(b^{(i)}; 0, \gamma \mathbf{I})$ priors for $b^{(u)}$ and $b^{(i)}$ respectively.

The objective function:

$$\begin{aligned} \mathcal{L} &= \log \mathcal{P}(R|U, V, b^{(u)}, b^{(i)}) + \log \mathcal{P}(U) + \log \mathcal{P}(V) \\ &\quad + \log \mathcal{P}(F) + \log \mathcal{P}(b^{(u)}) + \log \mathcal{P}(b^{(i)}) \\ &= -\frac{\lambda}{2} \sum_m \sum_{n \in \Omega(m)} \left(r_{mn} - (u_m^T v_n + b_m^{(u)} + b_n^{(i)}) \right)^2 \\ &\quad - \frac{\tau}{2} \sum_m u_m^T u_m - \frac{\tau}{2} \sum_m f_i^T f_i \\ &\quad - \frac{\tau}{2} \sum_n \left(v_n - \frac{1}{\sqrt{F_n}} \sum_{i=1}^{num_f} f_i \right)^T \left(v_n - \frac{1}{\sqrt{F_n}} \sum_{i=1}^{num_f} f_i \right) \\ &\quad - \frac{\gamma}{2} \sum_m b_m^{(u)2} - \frac{\gamma}{2} \sum_n b_n^{(i)2} \end{aligned} \quad (12)$$

proceeding the same way as in 3.3, we get the following updates at each step:

User vector update: No change happen, see equation (7).

User bias update: see equation (8).

Item vector update:

$$v_n = \lambda \left(\sum_{m \in \Omega(n)} u_m u_m^T + \tau \mathbf{I} \right)^{-1} \left(\lambda \sum_{m \in \Omega(n)} u_m (r_{mn} - b_m^{(u)} - b_n^{(i)}) + \tau f \right) \quad (13)$$

with $f = \frac{1}{\sqrt{F_n}} \sum_{i \in \text{features}(n)} f_i$.

Item bias update: refer to equation (10)

feature vector update:

$$f_i = \frac{\sum_{n \in \Omega(\mathcal{F}(i))} \left(v_n \frac{1}{\sqrt{F_n}} - \frac{1}{F_n} \sum_{j \in \mathcal{F}_n, j \neq i} f_j \right)}{\left(\sum_{n \in \Omega(\mathcal{F}(i))} \frac{1}{F_n} - 1 \right)} \quad (14)$$

where:

- $\Omega(\mathcal{F}(i))$ is the set of index of movies containing the feature i .
- \mathcal{F}_n is the set of features of the movie index by n .
- f_j is the feature vector of the feature j .

Note: "features(n)" is the set of feature factors of the movie indexed by n , F_n is the number of feature factors of the movie indexed by n and num_f is the total number of features (movie genres) in the dataset.

Algorithm 2 ALS (with feature information)

Input: data-structured, K, λ, τ, γ

Initialize $U, V, F, b^{(u)}, b^{(i)}$

repeat

for $m = 0$ **to** $M - 1$ **do**

for (n,r) in data-index-by-user[m] **do**

 update user vector $U[m :]$ with equation (7)

 update user bias $b_m^{(u)}$ with equation (8)

end for

end for

for $n = 0$ **to** $N - 1$ **do**

for (m,r) in data-index-by-item[n] **do**

 update item vector $V[n :]$ with equation (13)

 update item bias $b_n^{(i)}$ with equation (10)

end for

end for

for $i = 0$ **to** $num_f - 1$ **do**

 update feature $f[i :]$ with equation (14)

end for

until convergence

6.2. Feature Embeddings

The parameters of our model, including item genres information, have been optimized using $K = 2$ as the dimension of the latent space on the full 25 million ratings. The **figure 16** represents the **2D** embedding of the learned representations for "Children's" and "Horror" movies genres.

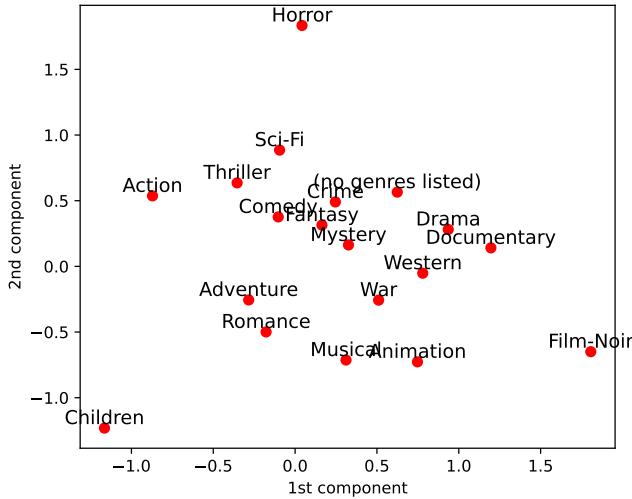


Figure 16. 2D feature embeddings

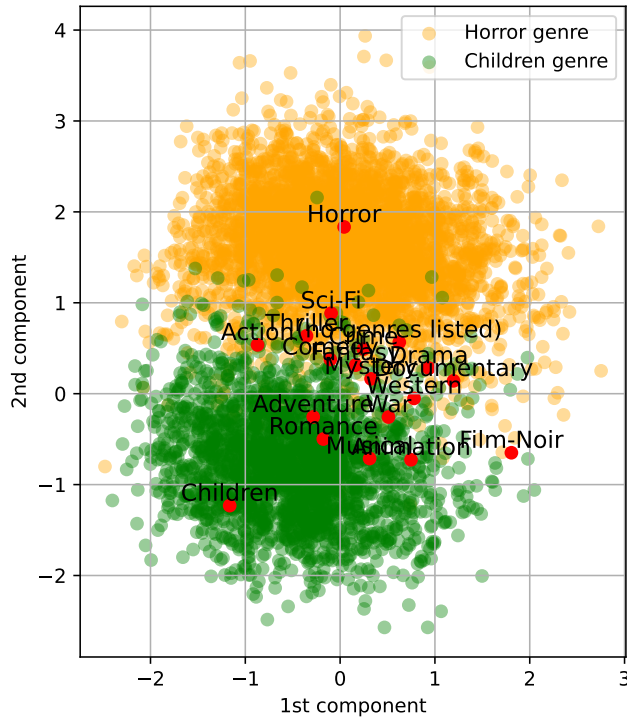


Figure 17. 2D item + feature embeddings

We can notice how better the embedding are, compared to the one with no featured information include (figure 10).

6.3. Recommendations

We have considered the hyperparameters used earlier ($K = 10$, $\lambda = 1$, $\gamma = 0.1$, $\tau = 1$) to optimize model parameters, including feature vectors. Then a "dummy

user" that gave the "Children's genre movie named "Toy Story" five stars have been created. Below are the results of the first twenty movies recommended by our recommender prototype.

title	genres
Toy Story (1995)	Adventure Animation Children Comedy Fantasy
Babe (1995)	Children Drama
Lion King, The (1994)	Adventure Animation Children Drama Musical IMAX
Aladdin (1992)	Adventure Animation Children Comedy Musical
White and the Seven Dwarfs (1937)	Animation Children Drama Fantasy Musical
Beauty and the Beast (1991)	Animation Children Fantasy Musical Romance IMAX
Pinocchio (1940)	Animation Children Fantasy Musical
Cinderella (1950)	Animation Children Fantasy Musical Romance
Dumbo (1941)	Animation Children Drama Musical
Little Mermaid, The (1989)	Animation Children Comedy Musical Romance
Babe: Pig in the City (1998)	Adventure Children Drama
Tarzan (1999)	Adventure Animation Children Drama
Toy Story 2 (1999)	Adventure Animation Children Comedy Fantasy
Monsters, Inc. (2001)	Adventure Animation Children Comedy Fantasy
Finding Nemo (2003)	Adventure Animation Children Comedy
Incredibles, The (2004)	Action Adventure Animation Children Comedy
Winnie the Pooh, The (1977)	Animation Children Musical
Toy Story 3 (2010)	Adventure Animation Children Comedy Fantasy IMAX
Winnie the Pooh and Tigger Too (1974)	Animation Children
Inside Out (2015)	Adventure Animation Children Comedy Drama Fantasy

Figure 18. Top twenty recommended movies

7. Conclusion

Recommendation systems are a fascinating subject, but they can quickly become overly complicated, especially when implemented on a large scale with millions of users and items. Matrix factorization provides a straightforward and convenient method for making recommendations when ratings are available. However, it has its pitfalls, one of which we have already encountered in our implementation: the "cold start problem." To tackle this issue, we added additional information about items. Generally, it is challenging to obtain feedback from users, as most users only rate items when they either really like them or absolutely hate them. In such cases, we often need to find ways to measure implicit feedback and employ negative sampling techniques to create a reasonable training set.

Find the source code: [Github](#)

¹AUTHORERR: Missing \icmlaffiliation. . Correspondence to: Herman Tesso <herman@aims.ac.za>.

References

- Koenigstein, N., Nice, N., Paquet, U., and Schleyen, N. The xbox recommender system. In *Proceedings of the sixth ACM conference on Recommender systems*, pp. 281–284, 2012.
- Koren, Y., Bell, R., and Volinsky, C. Matrix factorization techniques for recommender systems. *Computer*, 42(8): 30–37, 2009.
- Lü, L., Medo, M., Yeung, C. H., Zhang, Y.-C., Zhang, Z.-K., and Zhou, T. Recommender systems. *Physics reports*, 519(1):1–49, 2012.
- Schafer, J. B., Frankowski, D., Herlocker, J., and Sen, S. Collaborative filtering recommender systems. In *The adaptive web: methods and strategies of web personalization*, pp. 291–324. Springer, 2007a.
- Schafer, J. B., Frankowski, D., Herlocker, J., and Sen, S. Collaborative filtering recommender systems. In *The adaptive web: methods and strategies of web personalization*, pp. 291–324. Springer, 2007b.



AIMS

African Institute for
Mathematical Sciences
SOUTH AFRICA

Variational Bayes approach

$$R \approx UV^\top$$

We would like to average Over posterior distribution:

$$\begin{aligned} p(r_{45} \mid \mathbf{R}) &= \int p(r_{45}, \mathbf{u}_4, \mathbf{v}_5 \mid \mathbf{R}) d\mathbf{u}_4 d\mathbf{v}_5 \\ &= \int p(r_{45} \mid \mathbf{u}_4, \mathbf{v}_5) p(\mathbf{u}_4, \mathbf{v}_5 \mid \mathbf{R}) d\mathbf{u}_4 d\mathbf{v}_5 \\ &= \int p(r_{45} \mid \mathbf{u}_4, \mathbf{v}_5) p(\mathbf{U}, \mathbf{V} \mid \mathbf{R}) d\mathbf{U} d\mathbf{V} \end{aligned}$$

Problem:

With a “point mass” on U and V , like we had from the Alternating Least Square approach (ALS), this is still analytically tractable. Now it is not.

We want to approximate

$$P(U, V \mid R) = \frac{P(R \mid U, V)P(U)P(V)}{P(R)}$$

with factorized distributions, so that we can store two vectors (mean and variances) of numbers per **user** and **item** resulting in a small memory footprint on disk.

Method: Mean Field variational Inference:

A variational approximation

$$P(U, V \mid R) \approx q(U)q(V)$$

, where

$$\begin{aligned} q(V) &= \prod_n q(v_n); \quad q(v_n) = \prod_{k=1}^K \mathcal{N}(v_{n_k}; \mu_{v_{n_k}}, \sigma_{v_{n_k}}^2) \\ q(U) &= \prod_m q(u_m); \quad q(u_m) = \prod_{k=1}^K \mathcal{N}(u_{m_k}; \mu_{u_{m_k}}, \sigma_{u_{m_k}}^2) \end{aligned}$$

Note: m is the user index; n is the item index; K is the dimension of the latent space.

0.0.1 Objective function: minimize the KL-divergence

$$\min_{q(U)q(V)} D_{KL}(q(U)q(V) \parallel P(U, V \mid R))$$

$$\begin{aligned} D_{KL}(q(U)q(V) \parallel P(U, V \mid R)) &= \int_{q(U)q(V)} q(U)q(V) \log \frac{q(U)q(V)}{P(U, V \mid R)} dU dV \\ &= \mathbb{E}[\log q(U)q(V)] - \mathbb{E}[\log p(U, V \mid R)] \\ &= \mathbb{E}[\log q(U)q(V)] - \mathbb{E} \left[\log \frac{P(U, V, R)}{P(R)} \right] \quad \text{Baye's formula} \\ &= \mathbb{E}[\log q(U)q(V)] - \mathbb{E}[\log p(U, V, R)] + \log P(R) \\ &= - \left(\mathbb{E}[\log P(U, V, R)] - \mathbb{E}[\log q(U)q(V)] \right) + \log P(R) \quad (*) \end{aligned}$$

$$\begin{aligned} D_{KL}(q(U)q(V) \parallel P(U, V \mid R)) \geq 0 &\implies - \left(\mathbb{E}[\log P(U, V, R)] - \mathbb{E}[\log q(U)q(V)] \right) + \log P(R) \geq 0 \\ &\implies \log P(R) \geq \left(\mathbb{E}[\log P(U, V, R)] - \mathbb{E}[\log q(U)q(V)] \right) \quad (**) \end{aligned}$$

(*) and (**) \Rightarrow **Minimizing** D_{KL} = **Maximizing the lower bound of the evidence (ELBO)**, where

$$ELBO(q) = \mathbb{E}[\log P(U, V, R)] - \mathbb{E}[\log q(U)q(V)] \quad (1)$$

Note: the expectation here is under $q(U)q(V)$ ($\mathbb{E} \equiv \mathbb{E}_{q(U)q(V)}$)

0.0.2 Variational update

$$q(U)q(V) = \prod_m q(u_m) \prod_n q(v_n)$$

Let's ' s ' be the index of one a single latent variable, we denote by ' s ' the set of indexes of all the others latent variables except ' s '. Then, we can write

$$q(U)q(V) = \prod_s q_s(z_s), \quad P(U, V, R) = P(z_s, z_{-s}, R), \quad z_s \in \{u_1, \dots, u_M, v_1, \dots, v_N\}$$

In each step, we want to Maximize the **ELBO** (1) but we only have control over one q_s function.

$$\begin{aligned}
ELBO(q_s) &= \mathbb{E}[\log P(z_s, z_{-s}, R)] - \mathbb{E}[\log \prod_l q_l(z_l)] \\
&= \mathbb{E}_s \left[\mathbb{E}_{-s} [\log P(z_s, z_{-s}, R)] \right] - \mathbb{E}_s \left[\mathbb{E}_{-s} \left[\sum_l \log q_l(z_l) \right] \right] \\
&= \mathbb{E}_s \left[\mathbb{E}_{-s} [\log P(z_s, z_{-s}, R)] \right] - \mathbb{E}_s \left[\log q_s(z_s) + \underbrace{\mathbb{E}_{-s} \left[\sum_{l \neq s} \log q_l(z_l) \right]}_{\text{constant with regard to } q_s} \right] \\
&= \mathbb{E}_s \left[\mathbb{E}_{-s} [\log P(z_s, z_{-s}, R)] \right] - \mathbb{E}_s [\log q_s(z_s)] + \text{const} \quad (I \text{ don't care about the constant.}) \\
&= - \left(\mathbb{E}_s [\log q_s(z_s)] - \mathbb{E}_s \left[\mathbb{E}_{-s} [\log P(z_s, z_{-s}, R)] \right] \right) + \text{const} \\
&= - \left(\mathbb{E}_s [\log q_s(z_s)] - \mathbb{E}_s \left[\log e^{\mathbb{E}_{-s} [\log P(z_s, z_{-s}, R)]} \right] \right) + \text{const}, \quad (I \text{ notice the } D_{KL}) \\
&= -D_{KL} \left(q_s(z_s) \parallel e^{\mathbb{E}_{-s} [\log P(z_s, z_{-s}, R)]} \right) + \text{const} \quad (***)
\end{aligned}$$

the const in (*) also account for the normalization constant in the right term of the D_{KL} .

$$\text{maximize } ELBO(q_s) \implies \text{maximize } -D_{KL} (***) \implies \text{minimize } D_{KL} (***) \implies$$

So if we want to minimize D_{KL} in (***), best would be to put $q_s(z_s)^* \propto e^{\mathbb{E}_{-s} [\log P(z_s, z_{-s}, R)]}$.
The variational update (optimum distribution) $q_s(z_s)^*$ for fixed s is given by:

$$q_s(z_s)^* \propto e^{\mathbb{E}_{-s} [\log P(U, V, R)]}, \quad \forall z_s \in \{u_1, \dots, u_M, v_1, \dots, v_N\}^* \quad (2)$$

Note: \mathbb{E}_{-s} is the expectation with regard to all the other latent variables except z_s .

0.0.3 Mean and variance vectors of the optimum distributions

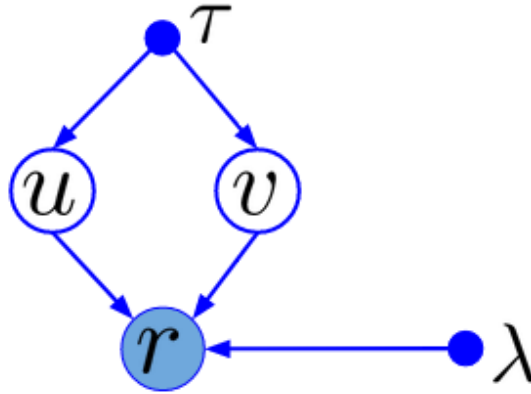


Figure 1: Bayesian Network of the model

Joint density:

From Fig. 1,

$$P(U, V, R) = p(U)p(V)p(R | U, V) = \prod_m p(u_m) \prod_n p(v_n) \prod_m \prod_{n \in \Omega(m)} p(r_{mn} | u_m, v_n)$$

$$\log P(U, V, R) = \sum_m \log p(u_m) + \sum_n \log p(v_n) + \sum_m \sum_{n \in \Omega(m)} \log p(r_{mn} | u_m, v_n)$$

for a fixed $m \in 1, \dots, M$, using 2, we have

$\log q(u_m) \approx \mathbb{E}_m[\log P(U, V, R)]$ only care about terms related to u_m all the others are consider constant

$$\begin{aligned} &\approx \log p(u_m) + \sum_{n \in \Omega(m)} \mathbb{E}_m[\log p(r_{mn} | u_m, v_n)] \\ &\approx -\frac{\tau}{2} u_m^T u_m - \frac{\lambda}{2} \sum_{n \in \Omega(m)} \mathbb{E}_m[(r_{mn} - u_m^T v_n)^2] \quad (*) \\ &\approx -\frac{\tau}{2} u_m^T u_m - \frac{\lambda}{2} \sum_{n \in \Omega(m)} (r_{mn}^2 + \mathbb{E}_m[(v_n^T u_m)^2] - 2r_{mn} u_m^T \mathbb{E}_m[v_n]) \\ &\approx -\frac{\tau}{2} u_m^T u_m - \frac{\lambda}{2} \sum_{n \in \Omega(m)} (r_{mn}^2 + \mathbb{E}_m[(v_n^T u_m)^T (v_n^T u_m)] - 2r_{mn} u_m^T \mathbb{E}_m[v_n]) \\ &\approx -\frac{\tau}{2} u_m^T u_m - \frac{\lambda}{2} \sum_{n \in \Omega(m)} (r_{mn}^2 + u_m^T \mathbb{E}_m[v_n v_n^T] u_m - 2r_{mn} u_m^T \mathbb{E}_m[v_n]) \quad u_m \text{ is const with regard to } \mathbb{E}_m \\ &\approx -\frac{\tau}{2} u_m^T u_m - \frac{\lambda}{2} \sum_{n \in \Omega(m)} (r_{mn}^2 + u_m^T (\Gamma_{v_n} + \mu_{v_n} \mu_{v_n}^T) u_m - 2r_{mn} u_m^T \mu_{v_n}) \\ &\approx -\frac{1}{2} \left[u_m^T \left(\tau \mathbb{I} + \lambda \sum_{n \in \Omega(m)} (\Gamma_{v_n} + \mu_{v_n} \mu_{v_n}^T) \right) u_m - 2u_m^T \left(\sum_{n \in \Omega(m)} \lambda r_{mn} \mu_{v_n} \right) + \dots \right] \end{aligned}$$

We notice the likelihood function developed for a multivariate normal distribution, up to a constant. We can then identify its **mean** $\mu_{u_m} = \Gamma_{u_m} (\sum_{n \in \Omega(m)} \lambda r_{mn} \mu_{v_n})$ and its **covariance**

matrix $\Gamma_{u_m} = \left(\tau \mathbb{I} + \lambda \sum_{n \in \Omega(m)} (\Gamma_{v_n} + \mu_{v_n} \mu_{v_n}^T) \right)^{-1}$.

$\forall m, u_m \sim \mathcal{N}(\mu_{u_m}, \Gamma_{u_m}), \text{ i.e}$

$$q(u_m) = \mathcal{N} \left(u_m; \left(\tau \mathbb{I} + \lambda \sum_{n \in \Omega(m)} (\Gamma_{v_n} + \mu_{v_n} \mu_{v_n}^T) \right)^{-1} \left(\sum_{n \in \Omega(m)} \lambda r_{mn} \mu_{v_n} \right), \left(\tau \mathbb{I} + \lambda \sum_{n \in \Omega(m)} (\Gamma_{v_n} + \mu_{v_n} \mu_{v_n}^T) \right)^{-1} \right) \quad (3)$$

Similarly, $\forall n, v_n \sim \mathcal{N}(\mu_{v_n}, \Gamma_{v_n}), \text{ i.e}$

$$q(v_n) = \mathcal{N} \left(v_n; \left(\tau \mathbb{I} + \lambda \sum_{m \in \Omega(n)} (\Gamma_{u_m} + \mu_{u_m} \mu_{u_m}^T) \right)^{-1} \left(\sum_{m \in \Omega(n)} \lambda r_{mn} \mu_{u_m} \right), \left(\tau \mathbb{I} + \lambda \sum_{m \in \Omega(n)} (\Gamma_{u_m} + \mu_{u_m} \mu_{u_m}^T) \right)^{-1} \right) \quad (4)$$

Remark:

- $q(u_m)$ depends on μ_{v_n} and Γ_{v_n} for $m \in \Omega(n)$ (set of users who rate the item " n).
- $q(v_n)$ depends on μ_{u_m} and Γ_{u_m} for $n \in \Omega(m)$ (set of items that have been rate by the user " m).

There is a non-linear circular dependency between the latent variables of the user and item. As a result, we need to iterate starting from random values until we find the optimum, i.e until the **ELBO** 1 has converged.

0.0.4 Calculating the ELBO

$$ELBO(q) = \mathbb{E}[\log P(U, V, R) - \log q(U)q(V)]$$

$$= \mathbb{E} \left[\sum_m \log P(u_m) + \sum_n \log P(v_n) + \sum_m \sum_{n \in \Omega(m)} \log P(r_m | u_m, v_n) - \sum_m \log q(u_m) - \sum_n \log q(v_n) \right]$$

$$\begin{aligned} \text{a) } \mathbb{E}[\log P(u_m)] &= -\frac{\tau}{2} \mathbb{E}[u_m^\top u_m] + \text{const} \\ &= -\frac{\tau}{2} [\text{trace}(\Gamma_{u_m}) + \mu_{u_m}^\top \mu_{u_m}] + \text{const} \end{aligned}$$

$$\begin{aligned} \text{b) } \mathbb{E}[\log p(v_n)] &= -\frac{\tau}{2} \mathbb{E}[v_n^\top v_n] + \text{const} \\ &= -\frac{\tau}{2} [\text{trace}(\Gamma_{v_n}) + \mu_{v_n}^\top \mu_{v_n}] + \text{const} \end{aligned}$$

$$\begin{aligned} \text{c) } \mathbb{E}[\log P(r_{mm} | u_m, v_m)] &= -\frac{\lambda}{2} \mathbb{E}[(r_{mm} - u_m^\top v_m)^2] + \text{const} \\ &= -\frac{\lambda}{2} (r_{mn}^2 + \mathbb{E}[(u_m^\top v_n)^2] - 2r_{mn} \mathbb{E}[u_m^\top v_n]) + \text{const} \\ &= -\frac{\lambda}{2} (r_{mn}^2 + \mathbb{E}[(u_m^\top v_n)^\top (u_m^\top v_n)] - 2r_{mn} \mathbb{E}[u_m^\top v_n]) + \text{const} \\ &= -\frac{\lambda}{2} (r_{mn}^2 + \mathbb{E}[\text{trace}(v_m^\top u_m u_m^\top v_n)] - 2r_{mn} \mathbb{E}[u_m^\top v_n]) + \text{const} \\ &= -\frac{\lambda}{2} (r_{mn}^2 + \mathbb{E}[\text{trace}(u_m u_m^\top v_n v_m^\top)] - 2r_{mn} \mathbb{E}[u_m^\top v_n]) + \text{const}, \text{ (permutation inside the "trace")} \\ &= -\frac{\lambda}{2} (r_{mn}^2 + \text{trace}(\mathbb{E}[u_m u_m^\top] \mathbb{E}[v_n v_n^\top]) - 2r_{mn} \mathbb{E}[u_m^\top] \mathbb{E}[v_n]) + \text{const} \\ &= -\frac{\lambda}{2} (r_{mn}^2 + \text{trace}(\mathbb{E}[u_m u_m^\top] \mathbb{E}[v_n v_n^\top]) - 2r_{mn} \mathbb{E}[u_m^\top] \mathbb{E}[v_n]) + \text{const}, \text{ (} u_m \text{ and } v_n \text{ are independent)} \\ &= -\frac{\lambda}{2} \text{trace}((\Gamma_{u_m} + \mu_{u_m} \mu_{u_m}^\top)(\Gamma_{v_n} + \mu_{v_n} \mu_{v_n}^\top)) + \lambda r_{mn} \mu_{u_m}^\top \mu_{v_n} + \text{const} \end{aligned}$$

$$\begin{aligned}
\text{d) } \mathbb{E}[\log q(u_m)] &= -\frac{1}{2} \left[\log(|\Gamma_{u_m}|) + K \log 2\pi + \mathbb{E} \left[(u_m - \mu_{u_m})^\top \Gamma_{u_m}^{-1} (u_m - \mu_{u_m}) \right] \right] \\
&= -\frac{1}{2} \left[\log(|\Gamma_{u_m}|) + \mathbb{E} \left[\text{trace} \left((u_m - \mu_{u_m})^\top \Gamma_{u_m}^{-1} (u_m - \mu_{u_m}) \right) \right] \right] \\
&= -\frac{1}{2} \left[\log(|\Gamma_{u_m}|) + \mathbb{E} \left[\text{trace} \left(\Gamma_{u_m}^{-1} (u_m - \mu_{u_m}) (u_m - \mu_{u_m})^\top \right) \right] \right] \\
&= -\frac{1}{2} \left[\log(|\Gamma_{u_m}|) + \text{trace} \left(\Gamma_{u_m}^{-1} \underbrace{\mathbb{E} \left[(u_m - \mu_{u_m}) (u_m - \mu_{u_m})^\top \right]}_{\Gamma_{u_m}} \right) \right] \\
&= -\frac{1}{2} \log(|\Gamma_{u_m}|) + \text{const}
\end{aligned}$$

$$\text{e) } \mathbb{E}[\log q(v_n)] = -\frac{1}{2} \log(|\Gamma_{v_n}|) + \text{const}$$

finally,

$$\begin{aligned}
ELBO &:= -\frac{\tau}{2} \sum_m [\text{trace}(\Gamma_{u_m}) + \mu_{u_m}^\top \mu_{u_m}] - \frac{\tau}{2} \sum_n (\text{trace}(\Gamma_{v_n}) + \mu_{v_n}^\top \mu_{v_n}) \\
&+ \lambda \sum_m \sum_{n \in \Omega(m)} \left[r_{mn} \mu_{u_m}^\top \mu_{v_n} - \frac{1}{2} \text{trace} \left((\Gamma_{u_m} + \mu_{u_m} \mu_{u_m}^\top) (\Gamma_{v_n} + \mu_{v_n} \mu_{v_n}^\top) \right) \right] \\
&+ \frac{1}{2} \sum_m \log(|\Gamma_{u_m}|) + \frac{1}{2} \sum_n \log(|\Gamma_{v_n}|) + \text{Const} \quad (5)
\end{aligned}$$

Note: As expected, the **ELBO** function depends on the distribution parameters of the latent vectors.

$$\mathbf{ELBO} = f(\Gamma_{u_m}, \Gamma_{v_n}, \mu_{u_m}, \mu_{v_n}) \quad (5)$$

Second approximation

But we don't want to save in memory the dense covariance matrices Γ_m , Γ_n for every user and item m, n .

Thus, we will approximate the learned $q(u_m)$ and $q(v_n)$ at each iteration, while ensuring that we continue to perform the variational update described in [\[5\]](#).

$$q(u_m) = \mathcal{N}(u_m; \mu_{u_m}, \Gamma_{u_m}) \approx \prod_k^K \mathcal{N}(\mu_{u_m}, c_{u_{m_k}} I) \quad \forall m,$$

Note: Given $p \sim \mathcal{N}(\mu_p, \Gamma_p)$ and $q \sim \mathcal{N}(\mu_q, \Gamma_q)$ both K dimensional, we can show that the KL -Divergence, $KL(p \parallel q)$ is given by:

$$\frac{1}{2} \left[\log \left(\frac{|\Gamma_q|}{|\Gamma_p|} \right) - K + \text{trace} \{ \Gamma_p \Gamma_q^{-1} \} + (\mu_p - \mu_q)^\top \Gamma_q^{-1} (\mu_p - \mu_q) \right].$$

Compute $KL\left(\prod_k^K \mathcal{N}(\mu_{u_m}, c_{u_{m_k}} I) \parallel q(u_m)\right)$

In our case, $\mu_p = \mu_q = \mu_m$, $\Gamma_q = \Gamma_m$ and $\Gamma_p = \text{Diag}((c_{u_{m_k}})_k)$ is a diagonal matrix. Thus, using the above formula,

$$\begin{aligned}
KL\left(\prod_k^K \mathcal{N}(\mu_{u_m}, c_{u_{m_k}} I) \parallel q(u_m)\right) &= \frac{1}{2} \left[\log(|\Gamma_m|) - \log\left(\prod_k^K c_{u_{m_k}}\right) - K + \sum_k c_{u_{m_k}} [\Gamma_m^{-1}]_k + 0 \right] \\
&= \frac{1}{2} \left[\log(|\Gamma_m|) - \sum_k \log(c_{u_{m_k}}) - K + \sum_k c_{u_{m_k}} [\Gamma_m^{-1}]_k + 0 \right] \\
&= \frac{1}{2} \left[\log(|\Gamma_m|) - \log\left(\prod_k^K c_{u_{m_k}}\right) - K + \sum_k c_{u_{m_k}} [\Gamma_m^{-1}]_k + 0 \right] \\
&= \frac{1}{2} \left[- \sum_k \log(c_{u_{m_k}}) + \sum_k c_{u_{m_k}} [\Gamma_m^{-1}]_k + \text{const} \right] \quad (\text{we account only on term})
\end{aligned}$$

$[\Gamma_m^{-1}]_k$ is the k^{th} element of the main diagonal vector of Γ_m^{-1} .

Find the Optimum $(c_{u_{m_k}})_k$ that minimize KL : Gradient

for fixed k ,

$$\frac{\partial KL}{\partial c_{u_{m_k}}} = \frac{1}{2} \left(-\frac{1}{c_{u_{m_k}}} + [\Gamma_m^{-1}]_k \right)$$

$$\frac{\partial KL}{\partial c_{u_{m_k}}} = 0 \implies c_{u_{m_k}} = \frac{1}{[\Gamma_m^{-1}]_k}$$

Recap & Algorithm

$$q(u_m) = \mathcal{N}(u_m; \mu_{u_m}, \Gamma_{u_m}) \approx \prod_k^K \mathcal{N}(\mu_{u_m}, c_{u_{m_k}} I) \quad \text{where}$$

$$\Gamma_{u_m} = \left(\tau \mathbb{I} + \lambda \sum_{n \in \Omega(m)} (\Gamma_{v_n} + \mu_{v_n} \mu_{v_n}^T) \right)^{-1}; \quad \mu_{u_m} = \Gamma_{u_m} \left(\sum_{n \in \Omega(m)} \lambda r_{mn} \mu_{v_n} \right); \quad c_{u_{m_k}} = \frac{1}{[\Gamma_{u_m}^{-1}]_k}$$

$$q(v_n) = \mathcal{N}(v_n; \mu_{v_n}, \Gamma_{v_n}) \approx \prod_k^K \mathcal{N}(\mu_{v_n}, c_{v_{n_k}} I) \quad \text{where}$$

$$\Gamma_{v_n} = \left(\tau \mathbb{I} + \lambda \sum_{m \in \Omega(n)} (\Gamma_{u_m} + \mu_{u_m} \mu_{u_m}^T) \right)^{-1}; \quad \mu_{v_n} = \Gamma_{v_n} \left(\sum_{m \in \Omega(n)} \lambda r_{mn} \mu_{u_m} \right); \quad c_{v_{n_k}} = \frac{1}{[\Gamma_{v_n}^{-1}]_k}$$

Algorithm 1 Mean Field Variational Inference

Input: data-structured, $K, \lambda, \tau, \epsilon$
Initialize: $\Gamma_{u_m}, \Gamma_{v_n}, \forall m, n$ (*positive definite matrices*)
Initialize: $\mu_{v_n}, \mu_{u_m}, \forall m, n$
repeat
 for $m = 0$ **to** $M - 1$ **do**
 for (n,r) in data-index-by-user[m] **do**
 update user mean vector μ_{u_m} (equation (3))
 update user variance Γ_{u_m} (equation (3))
 extract $Var_m = \frac{1}{Diag(\Gamma_{u_m}^{-1})}$
 end for
 end for
 for $n = 0$ **to** $N - 1$ **do**
 for (m,r) in data-index-by-item[n] **do**
 update item mean vector μ_{v_n} (equation (4))
 update item covariance Γ_{v_n} (equation (4))
 extract $Var_n = \frac{1}{Diag(\Gamma_{v_n}^{-1})}$
 end for
 end for
 Compute **ELBO** = $f(\Gamma_{u_m}, \Gamma_{v_n}, \mu_{u_m}, \mu_{v_n})$ (equation (5))
until $|ELBO_{t+1} - ELBO_t| < \epsilon$ (*converged*)

[Refer to this notebook for the implementation](#)