

Lecture 3 - Recommender Systems

Advanced Machine Learning

Miroslav Čepek, Zdeněk Buk, **Rodrigo da Silva Alves**,
Vojtěch Rybář, Petr Šimánek

FIT CTU

9. 3. 2023

Recommender Systems

- Recommenders recommend:
 - ▶ Items to users (most common)
 - ▶ Users to Items
 - ▶ Items to Items
 - ▶ Users to Users
- Items are movies, products, news, music, books, recipes, ...

Working in pairs: try to find one example of each four recommender scenarios above.

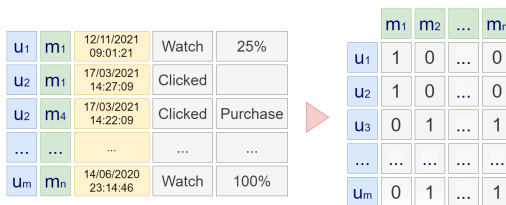
Recommender Systems

- WLOG, we will focus on recommend **users** relevant **items**
 - ▶ Predictive modelling: predict the rate of item m by the user u
 - ▶ Retrieve modelling: learning a rank systems
- Typically based on past **interactions** and (or) **attributes** (from users and items)
- **Interactions:** normally modelled as interaction matrix
 - ▶ Explicit: a user rate a song with 4 stars in a scale from 0 to 5
 - ▶ Implicit: a user watched 80% of a movie
- **Attributes:** normally modelled as attribute matrices
 - ▶ Users: gender, age, location, ...
 - ▶ Items: text, video, meta-data, ...

Modelling Interactions: Explicit feedback

	m_1	m_2	...	m_n
u_1	?	2	...	3
u_2	5	1	...	?
u_3	?	3	...	1
...
u_m	4	4	...	?

Modelling Interactions: Implicit feedback

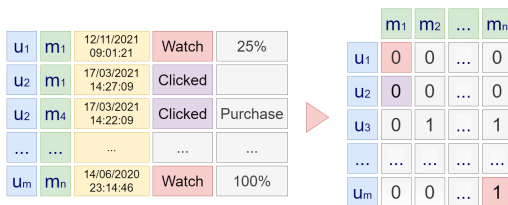


From explicit to implicit feedback

	m_1	m_2	...	m_n
u_1	?	2	...	3
u_2	5	1	...	?
u_3	?	3	...	1
...
u_m	4	4	...	?

	m_1	m_2	...	m_n
u_1	?	0	...	0
u_2	1	0	...	?
u_3	?	0	...	0
...
u_m	1	1	...	?

From implicit to explicit feedback



Modelling Attributes

	Color	Price	Category
m ₁	Black	24936	Chair
m ₂	Red and White	24944	Chair
m ₃	Red and White	1299	T-Shirt
m ₄	Black	1104	Chair



Personalised Machine Learning

- Personalisation *is not* a simple regression or classification problem
- A personalised model implies that: if the user have different interactions (or attribute) the recommendation should be different
- Suppose the vector a_u (a_m) are attribute vectors of user u (item m)
- We can use linear regression to predict how users u will like item m

$$r_{ui} = \omega^\top \times \begin{bmatrix} a_u \\ a_m \end{bmatrix}$$

- Is linear regression a personalised model for recommenders?
No!

Recommendation Algorithms

Collaborative Filtering



Day One: Joe and Julia independently read an article on police brutality



Day Two: Joe reads an article about deforestation, and then Julia is recommended the deforestation article

Content-Based Filtering



Day One: Julia watches a Drama



Day Two: Dramas are recommended

Recommender as a Matrix

- As we saw, we can model recommenders as matrices.
- The ratings can be stored in a **ranking matrix** R of dimension $m \times n$ with elements from $\mathbb{R} \cup \{?\}$.
- An example of a rating matrix for $m = 4$ users and $n = 6$ items can read

$$R = \begin{pmatrix} 1 & ? & ? & 2 & ? & 1 \\ ? & 2 & 3 & ? & 2 & 1 \\ 1 & 5 & 5 & ? & ? & 5 \\ ? & ? & 2 & ? & ? & 3 \end{pmatrix}.$$

Meaning, e.g., that user u_1 ranked items i_1 and i_6 with 1 star, item i_4 with 2 stars and had no interactions with items i_2, i_3 and i_5 .

- **Our goal is to predict the unknown ratings** $r_{u,i} = ?$ using the knowledge of the known ratings $r_{u,i} \neq ?$.

Idea of matrix factorization

- By **matrix factorization** we usually mean expressing a given matrix R as a matrix product of two (or more) matrices with some non-trivial properties. For example:

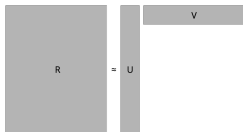
$$R = UV^{\top}$$

- These factorizations are a cornerstone of many algorithms and methods or are used to reach more numerically stable computations.

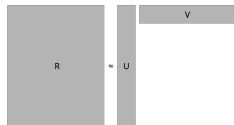
Do we need to know all the entries of a matrix R to factorize it, for example $R = UV^{\top}$?

Intuition of matrix factorization

- As for the recommendation systems, the inspiration comes mainly from the SVD as it can be used for constructing **latent features** or, in other words, **dimensionality reduction** using projections to lower dimensional space.
- The very basic idea of the lower dimensional approximation of an input matrix R of dimension $m \times n$ is based on this first-linear-algebra-lesson fact: Multiplying matrices U of dimension $m \times d$ and V of dimension $d \times n$ we get a matrix of dimension $m \times n$. This is true for any positive integer d .
- And this is the idea: **Given a rating matrix R , find lower dimensional matrices U and V so that the known elements of R are well approximated by the matrix UV^T .**



Matrix factorization for recommenders



- Let us denote:
 - ▶ The i -th row of U as u_i ; the number of rows of U equals the number of users $|\mathcal{U}|$.
 - ▶ The j -th column of V as v_j ; the number of columns of V equals the number of items $|\mathcal{I}|$.
 - ▶ Ω the subset of $\mathcal{U} \times \mathcal{I}$ of user-item pairs (i, j) such that $r_{i,j}$ is known, i.e., $r_{i,j} \neq ?$.
- The approximation of $r_{i,j}$ is given by the number $u_i^T v_j$, i.e., by the dot product of the two d -dimensional vectors.

Optimization Problem

- The **error of approximation** is usually measured by the squared residual:

$$(r_{i,j} - u_i^T v_j)^2.$$

- Hence, the matrices U and V are obtained by solving the optimization task:

$$\operatorname{argmin}_{\mathbf{U}, \mathbf{V}} \sum_{(i,j) \in \Omega} (r_{i,j} - u_i^T v_j)^2 + \lambda \left(\sum_x \|u_x\|^2 + \sum_y \|v_y\|^2 \right).$$

Sparsity and prediction

- The matrices U and V are optimized only by considering the known entries of R that are usually only a minority of entries.
- E.g. in the Netflix prize in 2006 there were $n = 17K$ movies and $m = 500K$ users, meaning that the matrix R had $8500M$ entries. But only $100M$ was given by Netflix!
- Still, the result of the matrix multiplication UV^\top is a matrix having the same dimensions as R **with all entries known!**
- The unknown rating $r_{i,j} = ?$ is estimated as $\hat{r}_{i,j} = u_i^T v_j..$

Example

- Consider our toy example matrix from above:

$$R = \begin{pmatrix} 1 & ? & ? & 2 & ? & 1 \\ ? & 2 & 3 & ? & 2 & 1 \\ 1 & 5 & 5 & ? & ? & 5 \\ ? & ? & 2 & ? & ? & 3 \end{pmatrix}.$$

- Assume that we chose the hyperparameter $d = 2$, i.e., we look for approximation matrices U and V with dimensions 4×2 and 2×6 , respectively.
- Let us pretend that the matrices resulting from the optimization are

$$U = \begin{pmatrix} 0.3 & 0.7 \\ 0.3 & 0.5 \\ 0.2 & 0.4 \\ 0.2 & 0.1 \end{pmatrix} \quad \text{and} \quad V^{\top} = \begin{pmatrix} 1 & 10 & 11 & 10 & 4 & 20 \\ 1 & -1 & -2 & -1 & 1 & -4 \end{pmatrix}.$$

Example

- The resulting approximation is

$$\mathbf{UV}^T = \begin{pmatrix} 0.3 & 0.7 \\ 0.3 & 0.5 \\ 0.2 & 0.4 \\ 0.2 & 0.1 \end{pmatrix} \begin{pmatrix} 1 & 10 & 11 & 10 & 4 & 20 \\ 1 & -1 & -2 & -1 & 1 & -4 \end{pmatrix} =$$
$$= \begin{pmatrix} 1 & \mathbf{2.3} & \mathbf{1.9} & 2.3 & \mathbf{1.9} & 3.2 \\ \mathbf{0.8} & 2.5 & 2.3 & \mathbf{2.5} & 1.7 & 4 \\ 0.6 & 1.6 & 1.4 & \mathbf{1.6} & \mathbf{1.2} & 2.4 \\ \mathbf{0.3} & \mathbf{1.9} & 2 & \mathbf{1.9} & \mathbf{0.9} & 3.6 \end{pmatrix},$$

where the red numbers are the desired predictions!

- E.g. the 3rd user predicted rating of the 4th item is $\hat{r}_{3,4} = 1.6$.

Supervised learning task

- The learning parameters: $U \in \mathbb{R}^{m \times d}$ and $V \in \mathbb{R}^{n \times d}$
- The hyperparameters:
 - ▶ the regularization constant $\lambda > 0$,
 - ▶ the matrix dimension d , which is a positive integer (significantly smaller than $\min\{m, n\}$).
- These hyperparameters can be tuned in the usual way via crossvalidation
- Therefore we would like to learn U and V , given d and λ by
$$\operatorname{argmin}_{\mathbf{U}, \mathbf{V}} \sum_{(i,j) \in \Omega} (r_{i,j} - u_i^T v_j)^2 + \lambda \left(\sum_x \|u_x\|^2 + \sum_y \|v_y\|^2 \right).$$

Alternating least squares (ALS)

- The idea of ALS is to fix alternately the matrix U and V . The non-fixed matrix is then considered learning variable and a subject to minimization.
- With one of the matrices fixed, the optimization problem becomes convex and very similar to the linear regression problem.
- Let's try to understand how the mechanism works

Alternating least squares (ALS)

$$R = U \times V^T$$

The diagram illustrates the matrix factorization $R = U \times V^T$ in the context of Alternating Least Squares (ALS). The matrix R is shown as a grid where the i -th row is highlighted in green and labeled u_i . This row is partitioned into two segments: one with diagonal hatching and one with a green dot pattern. The matrix U is a vertical column of boxes, with the i -th box highlighted in green and labeled u_i^T . The matrix V^T is a horizontal row of boxes, partitioned into four segments: diagonal hatching, solid blue, diagonal cross-hatching, and a blue dot pattern. A detailed view in a rounded rectangle shows the dot product calculation for the i -th row: $R_{\Omega^i} = u_i^T \times V_{\Omega^i}^T$, where R_{Ω^i} is a small green hatched box, u_i^T is a green box, and $V_{\Omega^i}^T$ is a small blue hatched box.

Alternating least squares (ALS)

$$R_{\Omega^i} = u_i^\top \times V_{\Omega^i}^\top$$

- Then we have the following optimization problem

$$\min_{u_i} \|R_{\Omega^i} - u_i^\top V_{\Omega^i}^\top\|^2 + \lambda \|u_i\|^2$$

- Convex problem with closed-form

$$\hat{u}_i = (V_{\Omega^i} V_{\Omega^i}^\top + \lambda I)^{-1} V_{\Omega^i}^\top R_{\Omega^i}$$

Alternating least squares (ALS)

Randomly initialize U and V

- WHILE** does not converge
 - ▶ $\forall i \in \mathcal{U}, \min_{u_i} \|R_{\Omega^i} - u_i^\top V_{\Omega^i}^\top\|^2 + \lambda \|u_i\|^2$
 - ▶ $\forall j \in \mathcal{I}, \min_{v_j} \|R_{\Omega^j} - v_j^\top U_{\Omega^j}^\top\|^2 + \lambda \|v_j\|^2$

MF for Implicit Feedback

- In real-world applications, we often observe more implicit feedback than explicit feedback.
- In fact, explicit feedback is sometimes considered implicit.
- Suppose user i watched 35% of movie A and 85% of movie B .

Does this mean that the user likes A more than B ? If so, does it mean that the user likes A more than twice as much as B ?

- The method we learned last class is more appropriate for explicit feedback. Why?

Modelling Implicit Feedback

- Let's understand a more appropriate method
- Assume the binary interaction matrix P :

$$P = \begin{pmatrix} 1 & 0 & 0 & 1 & 0 & 1 \\ 0 & 1 & 1 & 0 & 1 & 1 \\ 1 & 1 & 1 & 0 & 0 & 1 \\ 0 & 0 & 1 & 0 & 0 & 1 \end{pmatrix}.$$

- That is, if user- i interact with item- j , then $P_{ij} = 1$, otherwise $P_{ij} = 0$.
- Now let C be a matrix of confidence regarding the interaction:

$$C = \begin{pmatrix} 0.85 & 0 & 0 & 0.34 & 0 & 0.98 \\ 0 & 0.37 & 0.10 & 0 & 0.63 & 0.01 \\ 0.45 & 0.42 & 0.43 & 0 & 0 & 0.23 \\ 0 & 0 & 0.26 & 0 & 0 & 0.88 \end{pmatrix}.$$

Collaborative Filtering for Implicit Feedback

- Then we propose the following optimisation problem:

$$\min_{U,V} \sum_{i,j} C_{ij} (P_{ij} - u_i^\top v_j)^2 + \lambda \|u_i\|^2 + \lambda \|v_j\|^2$$

- Two main differences from previous MF method:
 - ▶ We need to account for the varying confidence levels
 - ▶ Optimization should account for all possible j, j pairs, rather than only those corresponding to observed data.
- We can use gradient descent to solve it.
- And ALS? By fixing V , can we find u_i ?

Closed form

- Assume V being fix and let's find u_i .
- Then we need to minimize the following loss

$$\mathcal{L}_i = \min_{u_i} \sum_j C_{ij} (P_{ij} - u_i^\top v_j)^2 + \lambda \|u_i\|^2$$

That is the same of:

$$\mathcal{L}_i = \min_{u_i} \sum_j (\sqrt{C_{ij}} (P_{ij} - u_i^\top v_j))^2 + \lambda \|u_i\|^2$$

Exercise: Find the closed form.

Alternating least squares (ALS)

$$\mathcal{L}_i = \min_{u_i} \sum_j (\sqrt{C_{ij}}(P_{ij} - u_i^\top v_j))^2 + \lambda \|u_i\|^2$$

The diagram illustrates the ALS optimization problem. It shows the relationship between the rating matrix C , the user latent matrix P , the item latent matrix V , and the user latent vector u_i .

- C : A grid representing the rating matrix. A specific row i is highlighted in purple, and a specific column j is highlighted in blue. The intersection is marked with a cross.
- P : A grid representing the user latent matrix. A specific row i is highlighted in red.
- V : A grid representing the item latent matrix. A specific column j is highlighted in blue.
- u_i^\top : A row vector representing the user latent vector for user i , highlighted in green.
- v_j : A column vector representing the item latent vector for item j , highlighted in blue.
- $\sqrt{C^i}$: A grid representing the square root of the rating matrix for row i . It is a diagonal matrix with the j -th element highlighted in blue.
- P_i : A column vector representing the user latent vector for user i , highlighted in red.
- V : A column vector representing the item latent vectors for all items, with elements v_1, v_2, \dots, v_8 listed. The j -th element is highlighted in blue.

The diagram shows the relationship between these matrices and vectors, illustrating the ALS optimization problem.

Closed form

- Therefore is the same of solving:

$$\mathcal{L}_i = ||\sqrt{C^i}P_i - \sqrt{C^i}Vu_i||^2 + \lambda + ||u_i||^2$$

- Taking the derivative

$$\nabla u_i = -2(\sqrt{C^i}V)^\top (\sqrt{C^i}P_i - \sqrt{C^i}Vu_i) + 2\lambda u_i$$

- Remind if D is diagonal $D = \sqrt{D} \times \sqrt{D}$ is trivial and $D = D^\top$
- Therefore, with just some algebraic derivations

$$u_i = (V^\top C^i V + \lambda I)^{-1} V^\top C^i P_i$$