

*annotated
version*

Machine Learning Course - CS-433

Matrix Factorizations

Nov 19, 2019

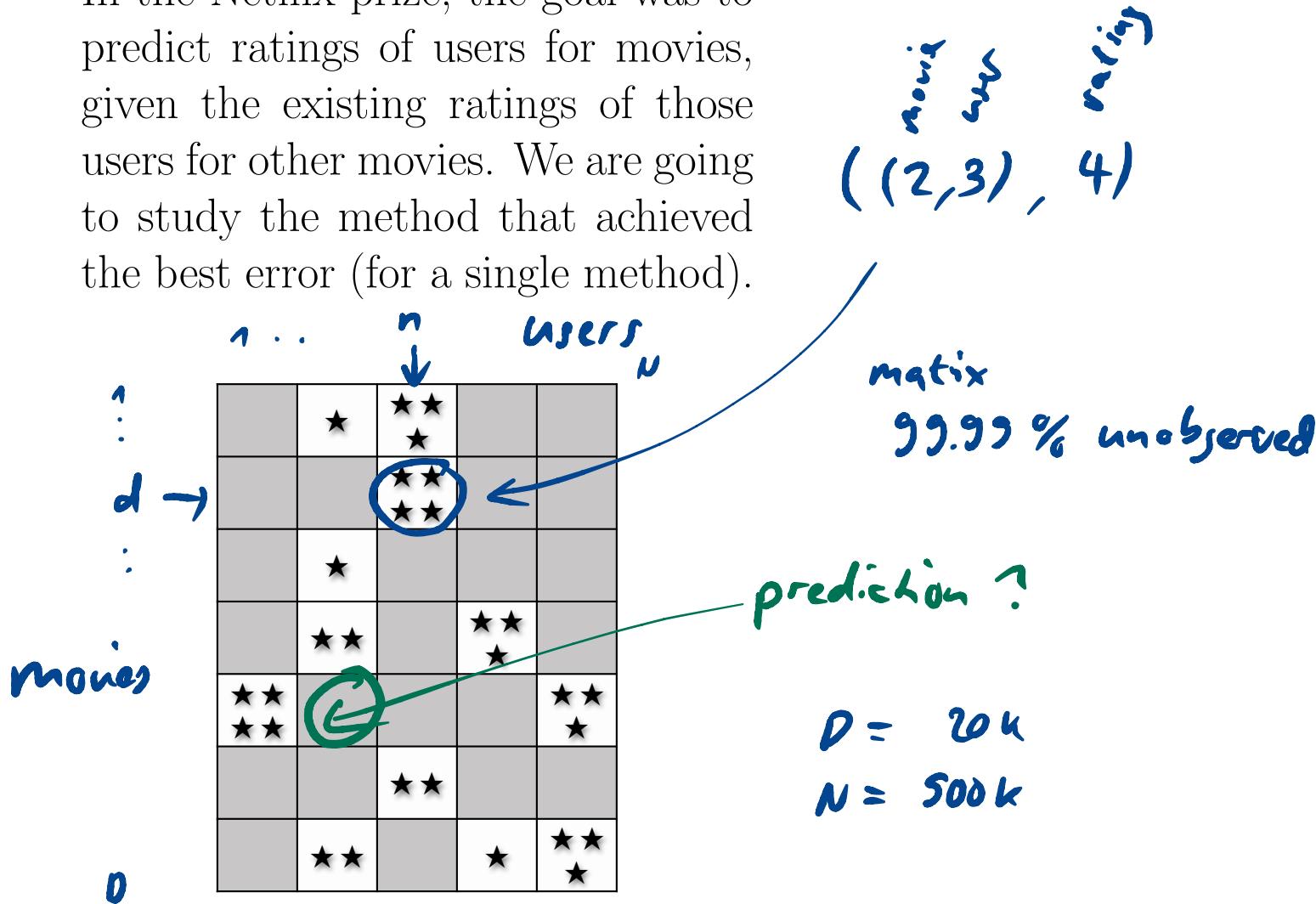
changes by Martin Jaggi 2019, 2018, 2017, ©Martin Jaggi and Mohammad Emtiyaz Khan 2016

Last updated on: November 19, 2019



Motivation

In the Netflix prize, the goal was to predict ratings of users for movies, given the existing ratings of those users for other movies. We are going to study the method that achieved the best error (for a single method).



The Movie Ratings Data

Given $\text{movies } d = 1, 2, \dots, D$ and $\text{users } n = 1, 2, \dots, N$, we define \mathbf{X} to be the $D \times N$ matrix containing all rating entries. That is, x_{dn} is the rating of n -th user for d -th movie.

Note that most ratings x_{dn} are missing and our task is to predict those missing ratings accurately.

SVD Special case $\mathcal{L} = \text{all entries } \|x - wz^\top\|_{\text{Fro}}^2$ $\min_{w,z} (\mathcal{L}(w,z) = f(wz^\top))$

Prediction Using a Matrix Factorization

We will aim to find \mathbf{W}, \mathbf{Z} s.t.

$$\mathbf{X} \approx \mathbf{W}\mathbf{Z}^\top.$$

So we hope to ‘explain’ each rating x_{dn} by a numerical representation of the corresponding movie and user

- in fact by the inner product of a movie feature vector with the user feature vector.

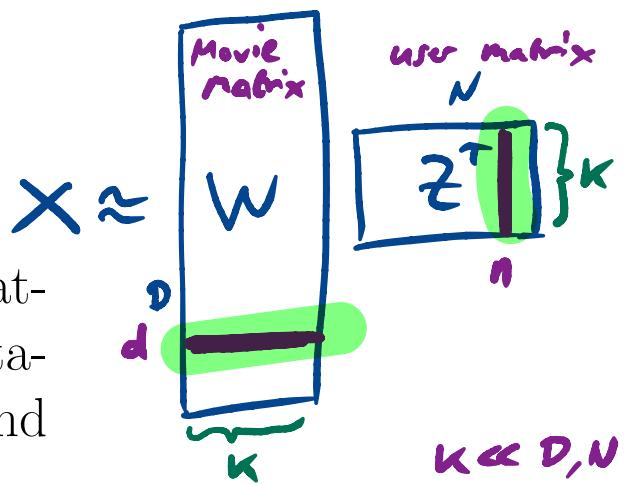
Loss function

$$\min_{\mathbf{W}, \mathbf{Z}} \mathcal{L}(\mathbf{W}, \mathbf{Z}) := \frac{1}{2} \sum_{(d,n) \in \Omega} (x_{dn} - (\mathbf{W}\mathbf{Z}^\top)_{dn})^2$$

where $\mathbf{W} \in \mathbb{R}^{D \times K}$ and $\mathbf{Z} \in \mathbb{R}^{N \times K}$ are tall matrices, having only $K \ll D, N$ columns.

The set $\Omega \subseteq [D] \times [N]$ collects the indices of the observed ratings of the input matrix \mathbf{X} .

Each row of those matrices is the feature representation of a movie (rows of \mathbf{W}) or a user (rows of \mathbf{Z}) respectively.

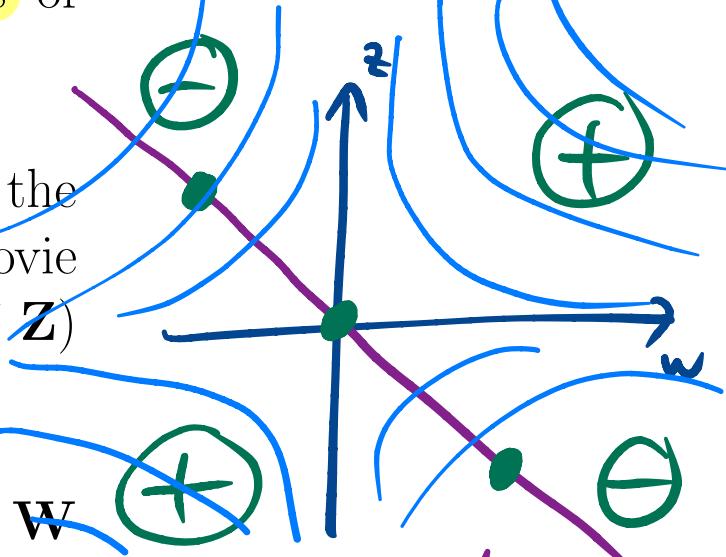


model = $w_d: \cdot z_n:$

① convex? Example $f = id$ no!

$$\mathcal{L}(w, z) = f(wz^\top)$$

Example = $w \cdot z \in \mathbb{R}$



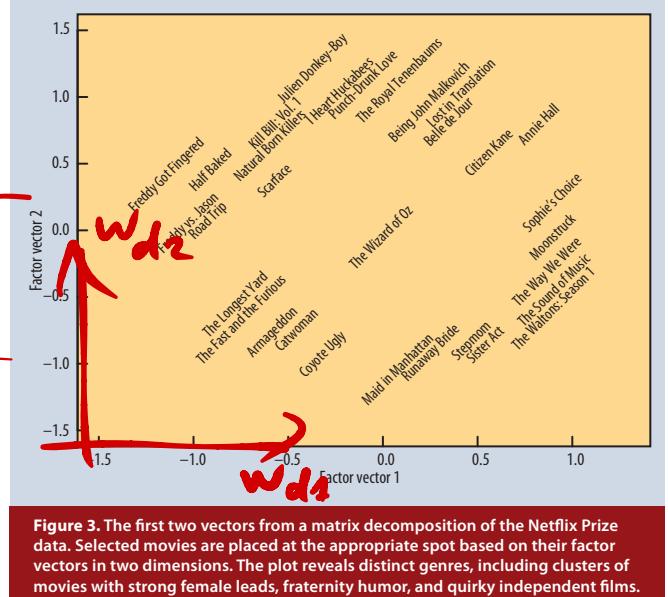
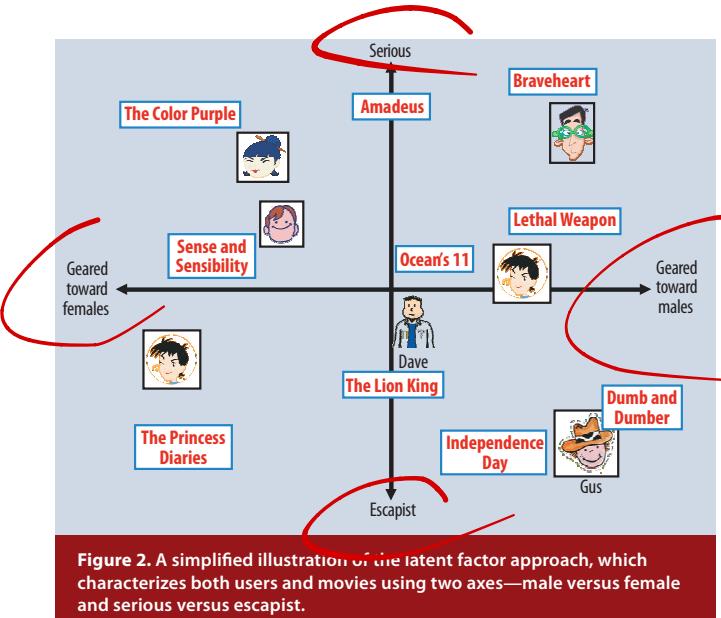
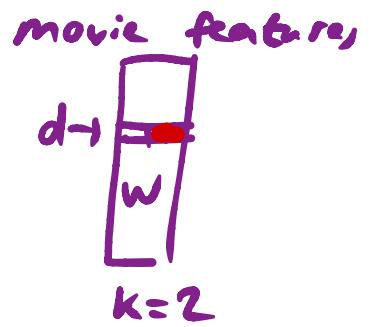
Is this cost jointly convex w.r.t. \mathbf{W} and \mathbf{Z} ? Is the model identifiable?

Given w^*, z^* also $(\beta w^*, \frac{1}{\beta} z^*)$ optimal
no!

Q: $f(wz^\top) := (wz)^\top$

Choosing K

K is the number of *latent* features.



Recall that for K -means, K was the number of clusters. (Similarly for GMMs, K was the number of latent variable dimensions).

Large K facilitates overfitting.

$$W \quad Z^T \\ X \quad 1 \\ 1 \quad X$$

if $K \geq \max(D, N)$

Regularization

We can add a regularizer and minimize the following cost:

$$\mathcal{L}(w, z) = \frac{1}{2} \sum_{(d,n) \in \Omega} [x_{dn} - (WZ^T)_{dn}]^2 + \frac{\lambda_w}{2} \|W\|_{\text{Frob}}^2 + \frac{\lambda_z}{2} \|Z\|_{\text{Frob}}^2$$

Loss

Regularizer

where $\lambda_w, \lambda_z > 0$ are scalars.

Stochastic Gradient Descent (SGD)

The training objective is a sum over $|\Omega|$ terms (one per rating):

$$L(w, z) = \frac{1}{|\Omega|} \sum_{(d,n) \in \Omega} \underbrace{\frac{1}{2} [x_{dn} - (\mathbf{WZ}^\top)_{dn}]^2}_{f_{dn}}$$

Derive the stochastic gradient for \mathbf{W}, \mathbf{Z} , given one observed rating $(d, n) \in \Omega$.

For one fixed element (d, n) of the sum, we derive the gradient entry (d', k) for \mathbf{W} , that is $\frac{\partial}{\partial w_{d',k}} f_{d,n}(\mathbf{W}, \mathbf{Z})$, and analogously entry (n', k) of the \mathbf{Z} part:

$$\frac{\partial}{\partial w_{d',k}} f_{d,n}(\mathbf{W}, \mathbf{Z}) = \begin{cases} -[x_{dn} - (\mathbf{WZ}^\top)_{dn}] z_{n,k} & \text{if } d' = d \\ 0 & \text{otherwise} \end{cases}$$

= prediction

= prediction error

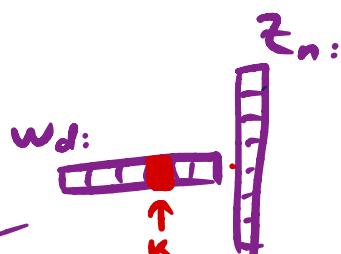
many "loss" functions

$$f_{dn}(w, z) \rightarrow \mathbb{R}$$

$$L = \sum_{(d,n)} f_{dn}$$

$$\nabla_w f_{dn}(w, z) \in \mathbb{R}^{D \times K}$$

$$\nabla_z f_{dn}(w, z) \in \mathbb{R}^{N \times K}$$



Cost: $O(K)$

for entire
 $\nabla_w f_{dn}$

$$\frac{\partial}{\partial z_{n',k}} f_{d,n}(\mathbf{W}, \mathbf{Z}) = \begin{cases} -[x_{dn} - (\mathbf{WZ}^\top)_{dn}] w_{d,k} & \text{if } n' = n \\ 0 & \text{otherwise} \end{cases}$$

$\frac{\partial}{\partial z_{n',k}}$

updates:

$$w^{(t+1)} := w^{(t)} - \delta \nabla_w f_{dn}(w, z)$$

$$z^{(t+1)} := z^{(t)} - \delta \nabla_z f_{dn}(w, z)$$

cost
 $O(K)$

Alternate ①, ②

Alternating Least-Squares (ALS)

For simplicity, let us first assume that there are **no missing** ratings, that is $\Omega := [D] \times [N]$. Then

$$\ell = \frac{1}{2} \sum_{d=1}^D \sum_{n=1}^N [x_{dn} - (\mathbf{WZ}^\top)_{dn}]^2 + \frac{\lambda_w}{2} \|\mathbf{W}\|_{\text{Frob}}^2 + \frac{\lambda_z}{2} \|\mathbf{Z}\|_{\text{Frob}}^2$$

add regularization

$$= \frac{1}{2} \|\mathbf{X} - \mathbf{WZ}^\top\|_{\text{Frob}}^2$$

①

②

$$\min_{\mathbf{w}, \mathbf{z}} \ell(\mathbf{w}, \mathbf{z})$$

$$\frac{1}{2} \|\mathbf{X} - \mathbf{wz}^\top\|_{\text{Frob}}^2$$

We can use coordinate descent to minimize the cost plus regularizer:

We first minimize w.r.t. \mathbf{Z} for fixed \mathbf{W} and then minimize \mathbf{W} given \mathbf{Z} .

$$\begin{aligned} ① \quad \mathbf{Z}^* &:= (\mathbf{W}^\top \mathbf{W} + \lambda_z \mathbf{I}_K)^{-1} \mathbf{W}^\top \mathbf{X} \\ ② \quad \mathbf{W}^* &:= (\mathbf{Z}^\top \mathbf{Z} + \lambda_w \mathbf{I}_K)^{-1} \mathbf{Z}^\top \mathbf{X}^\top \end{aligned}$$

What is the computational complexity? How can you decrease the cost when N and D are large?

$$\min_{\mathbf{w}, \mathbf{z}} \ell(\mathbf{w}, \mathbf{z})$$

②

$$\frac{1}{2} \|\mathbf{X} - \mathbf{wz}^\top\|_{\text{Frob}}^2$$

Least squares

Derivation:

$$\begin{aligned} \nabla_{\mathbf{w}} \ell(\mathbf{w}, \mathbf{z}) &\stackrel{!}{=} 0 & ① \\ \text{or} \quad \nabla_{\mathbf{z}} \ell(\mathbf{w}, \mathbf{z}) &\stackrel{!}{=} 0 & ② \end{aligned}$$

Cost: as in Ridge regression

per iteration

ALS with Missing Entries

Can you derive the ALS updates for the more general setting, when only the ratings $(d, n) \in \Omega$ contribute to the cost, i.e.

$$\mathcal{L} = \frac{1}{2} \sum_{(d,n) \in \Omega} [x_{dn} - (\mathbf{WZ}^\top)_{dn}]^2$$

*partially observed
(as in a recommender system)*

Hint: Compute the gradient with respect to each group of variables, and set to zero.

$$\nabla_w \mathcal{L}(w, z) \stackrel{!}{=} 0$$
$$\nabla_z \mathcal{L}(w, z) \stackrel{!}{=} 0$$