

annotated
version

Machine Learning Course - CS-433

Matrix Factorizations

Nov 21, 2017

©Martin Jaggi and Mohammad Emtiyaz Khan 2016

minor changes by Martin Jaggi 2017

Last updated: November 21, 2017



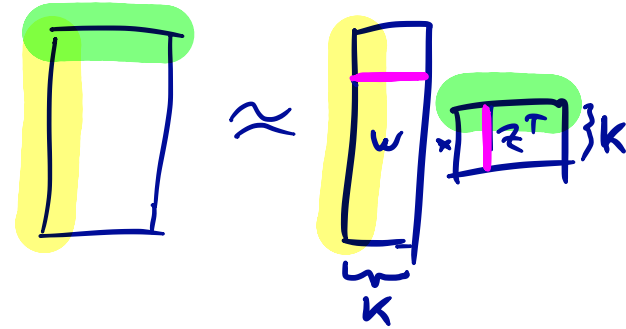
ÉCOLE POLYTECHNIQUE
FÉDÉRALE DE LAUSANNE

Note:
K-Means Notation was $X^T \approx WZ^T$

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).

$$X \approx (WZ^T)$$



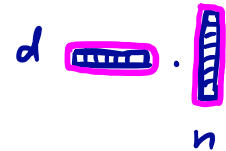
↓ user n

movie d →

	★	★★ ★		
		★★ ★★		
	★			
	★★		★★ ★	
★★ ★★				★★ ★
		★★		
	★★		★	★★ ★

3 stars

$$x_{dn} \approx (WZ^T)_{dn}$$



The Movie Ratings Data

Given **movies** $d = 1, 2, \dots, D$ and **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.

$$D = 20k$$

$$N = 500k$$

$$\begin{aligned} \# \text{observed entries} \\ = 100M \end{aligned}$$

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

Special case: $\Omega = \text{all} : \mathcal{L} = \|X - WZ^T\|_{\text{Frob}}^2$ see SVD

Prediction Using a Matrix Factorization

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

$$\mathbf{X} \approx \mathbf{W}\mathbf{Z}^T$$

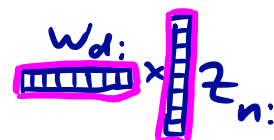
So we hope to 'explain' each rating x_{dn} by a numerical representation of the corresponding movie and user

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

Movie Matrix



User Matrix



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

$= f(\mathbf{w}\mathbf{z}^T)$

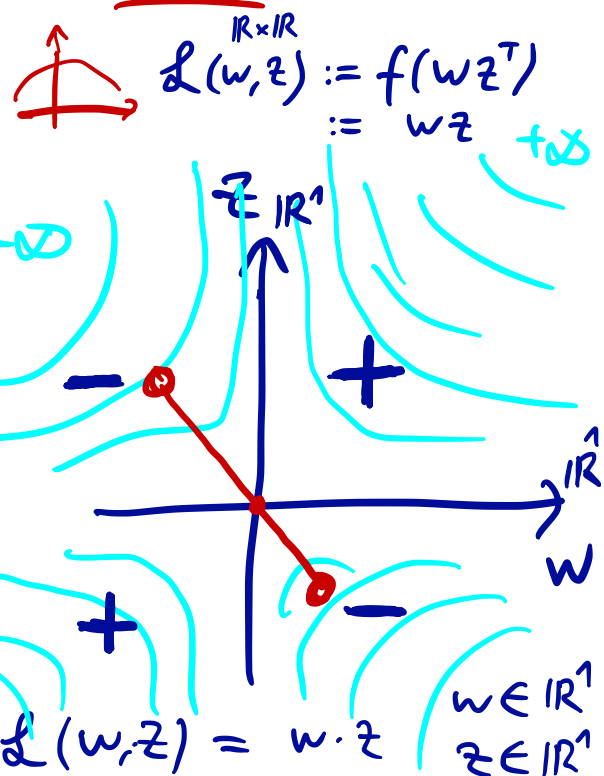
where $\mathbf{W} \in \mathbb{R}^{D \times K}$ and $\mathbf{Z} \in \mathbb{R}^{N \times K}$ are tall matrices, having only $K \ll D, N$ columns.

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.

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

① convex?



② optimal $\mathbf{w}^*, \mathbf{z}^*$
 $\Rightarrow 2 \cdot \mathbf{w}^*, \frac{1}{2} \mathbf{z}^*$ also optimal

② no!

Choosing K

K is the number of *latent* features.

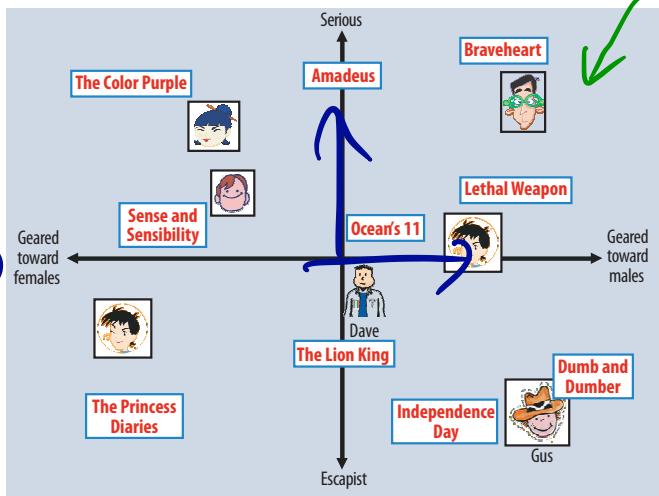
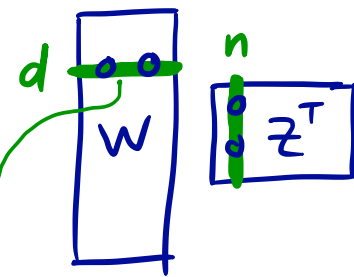


Figure 2. A simplified illustration of the latent factor approach, which characterizes both users and movies using two axes—male versus female and serious versus escapist.

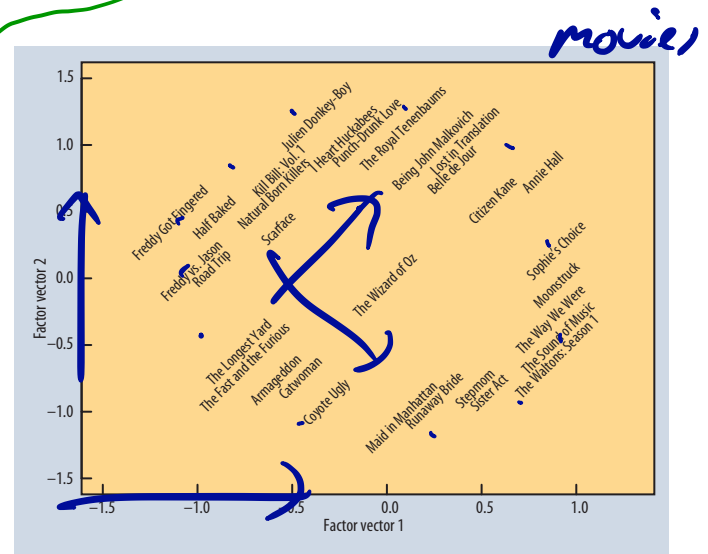


Figure 3. The first two vectors from a matrix decomposition of the Netflix Prize data. Selected movies are placed at the appropriate spot based on their factor vectors in two dimensions. The plot reveals distinct genres, including clusters of movies with strong female leads, fraternity humor, and quirky independent films.

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

if $K \geq \max\{D, N\}$:
Trivial Solution

$$X = \begin{bmatrix} \mathbf{1}_D & X \end{bmatrix}$$

$$X = \begin{bmatrix} X & \mathbf{1}_N \end{bmatrix}$$

Large K facilitates overfitting.

Regularization

We can add a regularizer and minimize the following cost:

$$\min_{\mathbf{W}, \mathbf{Z}} \frac{1}{2} \sum_{(d,n) \in \Omega} [x_{dn} - (\mathbf{W}\mathbf{Z}^T)_{dn}]^2 + \frac{\lambda_w}{2} \|\mathbf{W}\|_{\text{Frob}}^2 + \frac{\lambda_z}{2} \|\mathbf{Z}\|_{\text{Frob}}^2$$

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

quantifies the model complexity

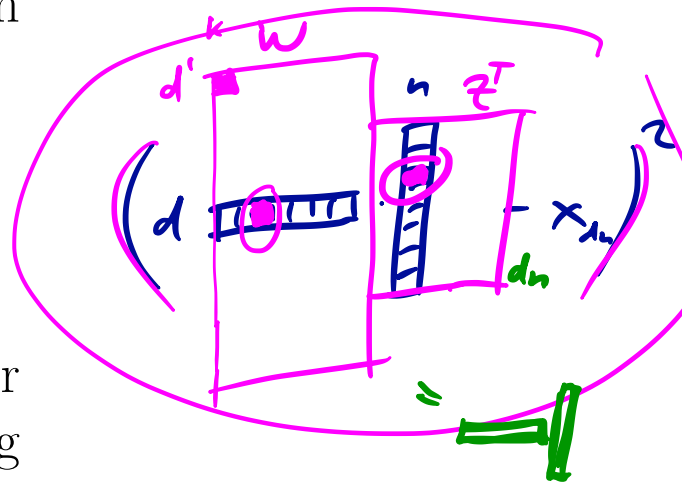


Invention: Blog Post "Netflix.. Try at home"
2006, Brandy Webb

Stochastic Gradient Descent (SGD)

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

$$\min_{\mathbf{W}, \mathbf{Z}} \frac{1}{2|\Omega|} \sum_{(d,n) \in \Omega} \frac{1}{2} \left[\overset{*}{x_{dn}} - \underbrace{(\mathbf{W}\mathbf{Z}^T)_{dn}}_{f_{d,n}(\mathbf{W}, \mathbf{Z})} \right]^2$$



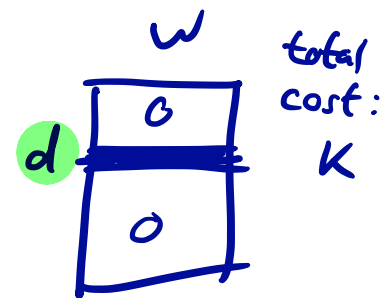
Derive the **stochastic gradient** for \mathbf{W}, \mathbf{Z} , given one observed rating $(d, n) \in \Omega$.

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:

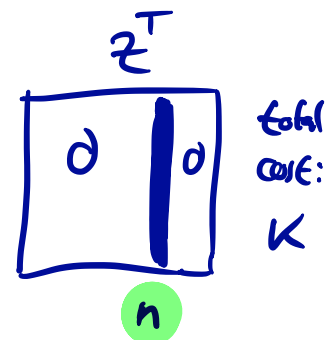
$$\begin{aligned} \nabla_{\mathbf{W}} f_{dn} &\in \mathbb{R}^{D \times K} \\ \nabla_{\mathbf{Z}} f_{dn} &\in \mathbb{R}^{N \times K} \\ \nabla_{\mathbf{W}, \mathbf{Z}} f_{dn}(\mathbf{W}, \mathbf{Z}) &\in \mathbb{R}^{(D+N) \times K} \end{aligned}$$

$$\frac{\partial}{\partial w_{d',k}} f_{d,n}(\mathbf{W}, \mathbf{Z}) = \begin{cases} -[x_{dn} - (\mathbf{W}\mathbf{Z}^T)_{dn}] z_{n,k} & \text{if } d' = d \\ 0 & \text{otherwise} \end{cases}$$

prediction error



$$\frac{\partial}{\partial z_{n',k}} f_{d,n}(\mathbf{W}, \mathbf{Z}) = \begin{cases} -[x_{dn} - (\mathbf{W}\mathbf{Z}^T)_{dn}] w_{d,k} & \text{if } n' = n \\ 0 & \text{otherwise} \end{cases}$$



updates:

$$\begin{aligned} \mathbf{W}^{(\epsilon+1)} &:= \mathbf{W}^{(\epsilon)} - \gamma \nabla f_{d,n}(\mathbf{W}^{(\epsilon)}, \mathbf{Z}^{(\epsilon)}) \\ \mathbf{Z}^{(\epsilon+1)} &:= \mathbf{Z}^{(\epsilon)} - \gamma \nabla f_{d,n}(\mathbf{W}^{(\epsilon)}, \mathbf{Z}^{(\epsilon)}) \end{aligned}$$

Alternating Minimization

Alternating Least-Squares (ALS)

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

Observation = all (n, d)

$$\mathcal{L}(w, z) = \frac{1}{2} \sum_{d=1}^D \sum_{n=1}^N [x_{dn} - (\mathbf{W}\mathbf{Z}^\top)_{dn}]^2 + \lambda_w \|\mathbf{w}\|_F^2 + \lambda_z \|\mathbf{z}\|_F^2$$

$\min_{w, z}$ $\mathcal{L}(w, z) = \frac{1}{2} \|\mathbf{X} - \mathbf{W}\mathbf{Z}^\top\|_{\text{Frob}}^2 + \dots + \dots$

Annotations: A green circle around \mathbf{W} and a blue circle around \mathbf{Z} in the matrix product $\mathbf{W}\mathbf{Z}^\top$. A green arrow points from the text "Ridge Regression target" to the \mathbf{Z} term in the equation.

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} .

$$\mathbf{Z}_* := (\mathbf{W}^\top \mathbf{W} + \lambda_z \mathbf{I}_K)^{-1} \mathbf{W}^\top \mathbf{X}$$

$$\mathbf{W}_*^\top := (\mathbf{Z}^\top \mathbf{Z} + \lambda_w \mathbf{I}_K)^{-1} \mathbf{Z}^\top \mathbf{X}^\top$$

Annotations: A green box around \mathbf{Z}_ and a blue box around \mathbf{W}_*^\top . A green arrow points from the text "Ridge Regression target" to the \mathbf{Z} term in the first equation. A blue arrow points from the text "argmin $\mathcal{L}(z, w)$ " to the \mathbf{Z}_* term. A blue arrow points from the text "argmin $\mathcal{L}(z, w)$ " to the \mathbf{W}_*^\top term. A green arrow points from the text "fixed" to the \mathbf{W} term in the first equation. A blue arrow points from the text "fixed" to the \mathbf{Z} term in the second equation.*

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

$$\nabla_w \mathcal{L} \stackrel{!}{=} 0$$

$$\nabla_z \mathcal{L} \stackrel{!}{=} 0$$

same as in
Least Squares

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.

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

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

↳ Normal equations

$$\nabla_{\mathbf{w}} \mathcal{L} \stackrel{!}{=} 0 \quad \Rightarrow \mathbf{w} = \dots$$

$$\nabla_{\mathbf{z}} \mathcal{L} \stackrel{!}{=} 0 \quad \Rightarrow \mathbf{z} = \dots$$