

Matrix Factorization

Mohammad Emtiyaz Khan
EPFL

Nov 19, 2015



ÉCOLE POLYTECHNIQUE
FÉDÉRALE DE LAUSANNE

©Mohammad Emtiyaz Khan 2015

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-rating data

We index D movies with $d = 1, 2, \dots, D$ and N users by $n = 1, 2, \dots, N$. We denote the rating of n 'th user for d 'th movie by x_{dn} . Define \mathbf{X} to be a $D \times N$ rating matrix with columns $\mathbf{x}_n \in \mathbb{R}^D$ containing ratings of the n 'th user.

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

What is the generalization error in this case?

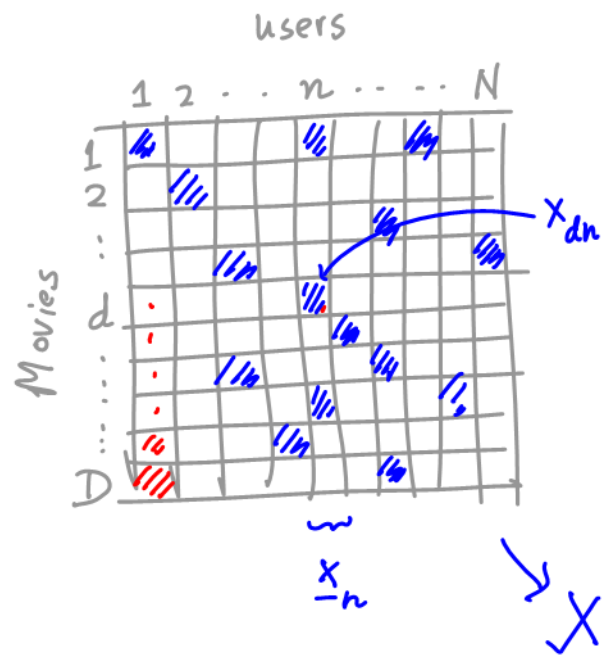
Let $O_n = \{ \text{all movies rated by the } n\text{'th user} \}$

Train error:
$$\sum_{n=1}^N \sum_{d \in O_n} (x_{dn} - \hat{x}_{dn})^2$$

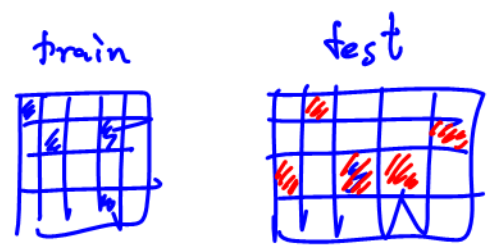
Generalization
Test error:
$$\sum_{n=1}^N \sum_{d \notin O_n} (x_{dn} - \hat{x}_{dn})^2$$

people

	↓	[Person Icon]	[Person Icon]	[Person Icon]
Star Wars		5	□	5
Star Trek		4	5	□
Notebook		□	1	□
⋮				



* Simulate the reality: Create train & test pairs by "punching holes".



* Read about A/B testing in Wikipedia

Minimizing reconstruction error

To simplify the problem, let us assume that there are no missing ratings.

Project data vectors \mathbf{x}_n to a smaller dimension $\mathbf{z}_n \in \mathbb{R}^M$ with $M < D$, such that the *reconstruction error* is minimized (assume M is known):

$$\min_{\mathbf{W}, \mathbf{Z}} \frac{1}{2} \sum_{n=1}^N \sum_{d=1}^D (x_{dn} - \mathbf{w}_d^T \mathbf{z}_n)^2$$

where \mathbf{W} is a $D \times M$ matrix and \mathbf{Z} is an $N \times M$ matrix, and \mathbf{w}_d^T and \mathbf{z}_n^T are rows of these matrices.

Is this cost jointly convex w.r.t. \mathbf{W} and \mathbf{Z} ? Is the model identifiable? **No**

Given \mathbf{W} , convex wrt \mathbf{Z} , and vice-versa.

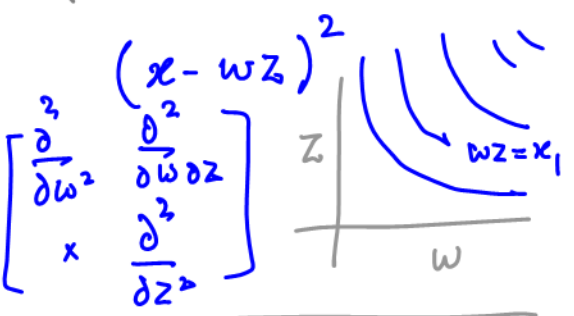
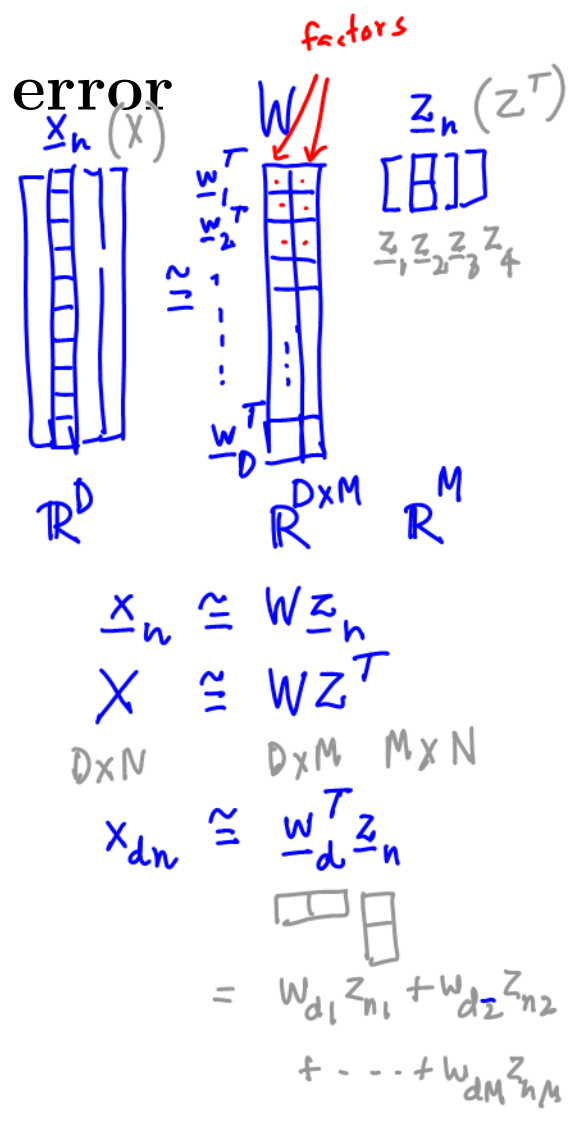
Regularization

We can add a regularizer and minimize the following cost:

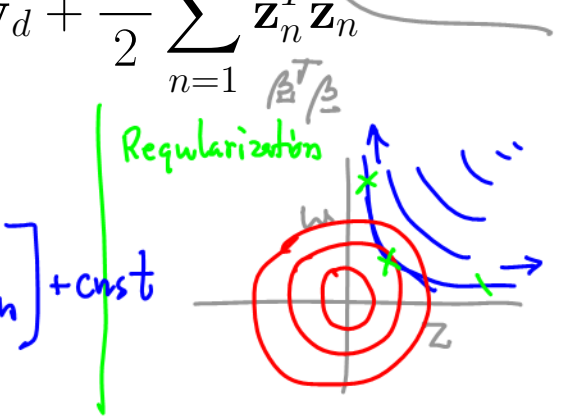
$$\mathcal{L}(\mathbf{W}, \mathbf{Z}) = \frac{1}{2} \sum_{n=1}^N \sum_{d=1}^D (x_{dn} - \mathbf{w}_d^T \mathbf{z}_n)^2 + \frac{\lambda_w}{2} \sum_{d=1}^D \mathbf{w}_d^T \mathbf{w}_d + \frac{\lambda_z}{2} \sum_{n=1}^N \mathbf{z}_n^T \mathbf{z}_n$$

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

$$\mathcal{L}(\mathbf{Z} | \mathbf{W}) = \frac{1}{2} \sum_{n=1}^N \left[(\mathbf{x}_n - \mathbf{W} \mathbf{z}_n)^T (\mathbf{x}_n - \mathbf{W} \mathbf{z}_n) + \lambda_z \mathbf{z}_n^T \mathbf{z}_n \right] + \text{const}$$



$\theta \rightarrow P_\theta(x)$ one-to-one
No since $\mathbf{W} \mathbf{Q} \mathbf{Q}^T \mathbf{Z}$ (orthogonal)



$$\underline{\beta}^* = (X^T X + \lambda I)^{-1} X^T \underline{y} \Rightarrow \underline{z}_n^* = (W^T W + \lambda_z I)^{-1} W^T \underline{x}_n$$

$M \times M$ $M \times D$ $D \times 1$
 $\underline{D M^2 + M^3}$

Alternating least-squares (ALS)

We can use coordinate descent algorithm, by first minimizing w.r.t. \mathbf{Z} given \mathbf{W} and then minimizing \mathbf{W} given \mathbf{Z} .

$$\mathbf{Z} \leftarrow (\mathbf{W}^T \mathbf{W} + \lambda_z \mathbf{I}_M)^{-1} \mathbf{W}^T \mathbf{X}$$

$$\mathbf{W} \leftarrow (\mathbf{Z}^T \mathbf{Z} + \lambda_w \mathbf{I}_M)^{-1} \mathbf{Z}^T \mathbf{X}^T$$

$$\underline{w}_d^* = (\mathbf{Z}^T \mathbf{Z} + \lambda_w \mathbf{I}_M)^{-1} \mathbf{Z}^T \underline{x}_d$$

\uparrow
 rows of \mathbf{X}

X	$D \times N$
W	$D \times M$
Z	$N \times M$

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

$(D M^2 + M^3) N$
 but since $D \gg M$, $\approx O(D N M^2)$
 Almost

Extensions

Dealing with missing entries

Denote the set of movies ~~not~~ rated by n 'th user by \mathbb{O}_n , then we can modify the cost as follows:

$$\frac{1}{2} \sum_{n=1}^N \sum_{d \in \mathbb{O}_n}^D (x_{dn} - \mathbf{w}_d^T \mathbf{z}_n)^2$$

$$\rightarrow \underline{z}_n^{(k+1)} = (W_{\mathbb{O}_n}^{(k)T} W_{\mathbb{O}_n}^{(k)} + \lambda_z I_M)^{-1} W_{\mathbb{O}_n}^{(k)T} \underline{x}_{\mathbb{O}_n}$$

Q: Write the update for \underline{z}_n given W & $\underline{x}_{\mathbb{O}_n}$

\underline{x}_n	$\underline{x}_{\mathbb{O}_n}$	$\underline{x}_{\mathbb{O}_n} \approx W_{\mathbb{O}_n} \underline{z}_n$
2×1	$2 \times M$	$M \times 1$

Algorithm

- * start with some $W^{(0)}$
- * Iterate

$\underline{z}^{(k+1)}$	\leftarrow	$W^{(k)}$
$W^{(k+1)}$	\leftarrow	$\underline{z}^{(k+1)}$

Derive the modified ALS algorithm.

Adding an offset

Since many ratings are missing, we cannot normalize the data. A solution is to add an offset term in the cost.

$$\frac{1}{2} \sum_{n=1}^N \sum_{d \in \mathcal{O}_n} (x_{dn} - \mathbf{w}_d^T \mathbf{z}_n - w_{0d} - z_{0n} - \mu)^2$$

where w_{0d} and z_{0n} are the offset terms for d 'th movie and n 'th user while μ is the global offset.

$$\left(\underbrace{x_{dn} - w_{0d} - \mu}_{\tilde{x}_{dn}} - \underbrace{\tilde{\mathbf{w}}_d^T}_{\tilde{\mathbf{w}}_d^T} \underbrace{\begin{bmatrix} z_{0n} \\ \mathbf{z}_n \end{bmatrix}}_{\tilde{\mathbf{z}}_n} \right)^2$$

Derive the modified ALS algorithm.

A Probabilistic Model

$$\prod_{n=1}^N \prod_{d \in \mathcal{O}_n} \mathcal{N}(x_{dn} / \frac{\mathbf{w}_d^T \mathbf{z}_n + 1}{-d - n})$$

$$\times \prod_n \mathcal{N}(\mathbf{z}_n / \sigma, \frac{1}{\lambda_z} \mathbf{I})$$

Write the probabilistic model associated with the regularized cost.

Generalization to Exponential Family

$$\times \prod_d \mathcal{N}(\frac{\mathbf{w}_d}{-d} / \sigma_0, \frac{1}{\lambda_w} \mathbf{I})$$

Using generalized linear model, we can extend matrix factorization to different types of datasets. See the following paper for details: "A Generalization of Principal Component Analysis to the Exponential Family" by Collins et. al., 2001.

Choosing number of factors M, λ_z, λ_w