# Lecture: Matrix Completion

# Recommendation systems



References:
http://bicmr.pku.edu.cn/~wenzw/bigdata/07-recsys1.pdf
http://bicmr.pku.edu.cn/~wenzw/bigdata/08-recsys2.pdf
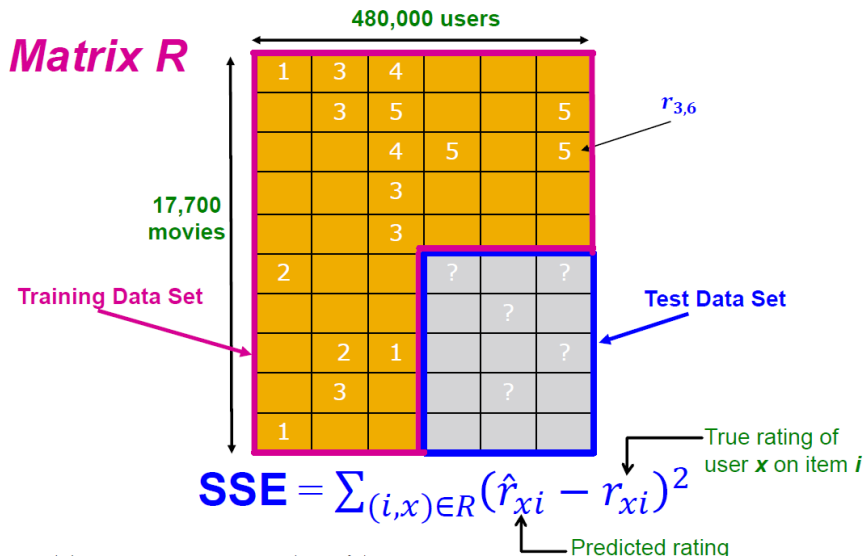
# The Netflix Prize

- Training data
  - 100 million ratings, 480,000 users, 17,770 movies
  - 6 years of data: 2000-2005

- Test data
  - Last few ratings of each user (2.8 million)
  - Evaluation criterion: root mean squared error (RMSE): $\sqrt{\sum_{xi}(r_{xi} - r_{xi}^*)^2}$: $r_{xi}$ and $r_{xi}^*$ are the predicted and true rating of $x$ on $i$
  - Netflix Cinematch RMSE: 0.9514

- Competition
  - 2700+ teams
  - $1 million prize for 10% improvement on Cinematch

**Matrix R**

480,000 users

17,700 movies

$r_{3,6}$

**Training Data Set**

**Test Data Set**

True rating of user $x$ on item $i$

$$\text{SSE} = \sum_{(i,x) \in R} (\hat{r}_{xi} - r_{xi})^2$$

Predicted rating

# Collaborative Filtering: weighted sum model

$$\hat{r}_{xi} = b_{xi} + \sum_{j \in N(i;x)} w_{ij}(r_{xj} - b_{xj})$$

- baseline estimate for $r_{xi}$: $b_{xi} = \mu + b_x + b_i$
  $\mu$: overall mean rating
  $b_x$: rating deviation of user x = (avg. rating of user x) - $\mu$
  $b_i$: (avg. rating of movie i) - $\mu$

- We sum over all movies j that are similar to i and were rated by x

- $w_{ij}$ is the interpolation weight (some real number). We allow:
  $\sum_{j \in N(i,x)} w_{ij} \neq 1$

- $w_{ij}$ models interaction between pairs of movies (it does not depend on user x)

- $N(i;x)$: set of movies rated by user x that are similar to movie i

# Finding weights $w_{ij}$?

Find $w_{ij}$ such that they work well on known (user, item) ratings:

$$\min_{w_{ij}} \quad F(w) := \sum_x \left( \left[ b_{xi} + \sum_{j \in N(i;x)} w_{ij}(r_{xj} - b_{xj}) \right] - r_{xi} \right)^2$$
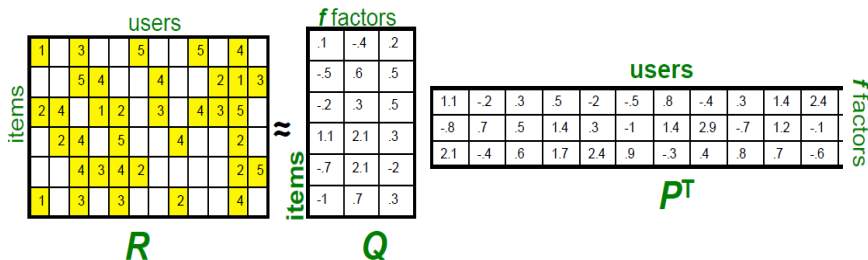
- Unconstrained optimization: quadratic function

$$\nabla_{w_{ij}} F(w) = 2 \sum_x \left( \left[ b_{xi} + \sum_{k \in N(i;x)} w_{ik}(r_{xk} - b_{xk}) \right] - r_{xi} \right) (r_{xj} - b_{xj}) = 0$$

for $j \in \{N(i,x), \forall i, x\}$

- Equivalent to solving a system of linear equations?

- Steepest gradient descent method: $w^{k+1} = w^k - \tau \nabla F(w)$

- Conjugate gradient method

# Latent factor models

- low rank factorization on Netflix data: $R \approx Q \cdot P^T$



**R**    **Q**    $P^T$

- For now let's assume we can approximate the rating matrix R as a product of "thin" $Q \cdot P^T$
  R has missing entries but let's ignore that for now! Basically, we will want the reconstruction error to be small on known ratings and we don't care about the values on the missing ones

# Ratings as products of factors

- How to estimate the missing rating of user x for item i?

$$\hat{r}_{xi} = q_i \cdot p_x^T = \sum_f q_{if} p_{xf},$$

where $q_i$ is row i of $Q$ and $p_x$ is column x of $P^T$

# Latent factor models

- Minimize SSE on training data!
- Use specialized methods to find P, Q such that $\hat{r}_{xi} = q_i \cdot p_x^T$

$$\min_{P,Q} \sum_{(i,x) \in \text{training}} (r_{xi} - q_i \cdot p_x^T)^2$$

  We don't require cols of P, Q to be orthogonal/unit length

- P, Q map users/movies to a latent space
- Add regularization:

$$\min_{P,Q} \sum_{(i,x) \in \text{training}} (r_{xi} - q_i \cdot p_x^T)^2 + \lambda \left[ \sum_x \|p_x\|_2^2 + \sum_i \|q_i\|_2^2 \right]$$

$\lambda$ is called regularization parameters

# Gradient descent method

$$\min_{P,Q} \quad F(P,Q) := \sum_{(i,x)\in \text{training}} (r_{xi} - q_i \cdot p_x^T)^2 + \lambda \left[ \sum_x \|p_x\|_2^2 + \sum_i \|q_i\|_2^2 \right]$$

Gradient decent:

- Initialize P and Q (using SVD, pretend missing ratings are 0)

- Do gradient descent:
  $P^{k+1} \leftarrow P^k - \tau \nabla_P F(P^k, Q^k)$,
  $Q^{k+1} \leftarrow Q^k - \tau \nabla_Q F(P^k, Q^k)$,
  where $(\nabla_Q F)_{if} = -2 \sum_{x,i} (r_{xi} - q_i p_x^T) p_{xf} + 2\lambda q_{if}$. Here $q_{if}$ is entry f of row $q_i$ of matrix Q

- Computing gradients is slow when the dimension is huge

# Stochastic gradient descent method

Observation: Let $q_{if}$ be entry f of row $q_i$ of matrix Q

$$
\begin{aligned}
(\nabla_Q F)_{if} &= \sum_{x,i} \left( -2(r_{xi} - q_i p_x^T) p_{xf} + 2\lambda q_{if} \right) = \sum_{x,i} \nabla_Q F(r_{xi}) \\
(\nabla_P F)_{xf} &= \sum_{x,i} \left( -2(r_{xi} - q_i p_x^T) q_{xf} + 2\lambda p_{if} \right) = \sum_{x,i} \nabla_P F(r_{xi})
\end{aligned}
$$

Stochastic gradient decent:

- Instead of evaluating gradient over all ratings, evaluate it for each individual rating and make a step

- $P \leftarrow P - \tau \nabla_P F(r_{xi})$
  $Q \leftarrow Q - \tau \nabla_Q F(r_{xi})$

- Need more steps but each step is computed much faster

# Latent factor models with biases

predicted models:

$$\hat{r}_{xi} = \mu + b_x + b_i + q_i \cdot p_x^T$$

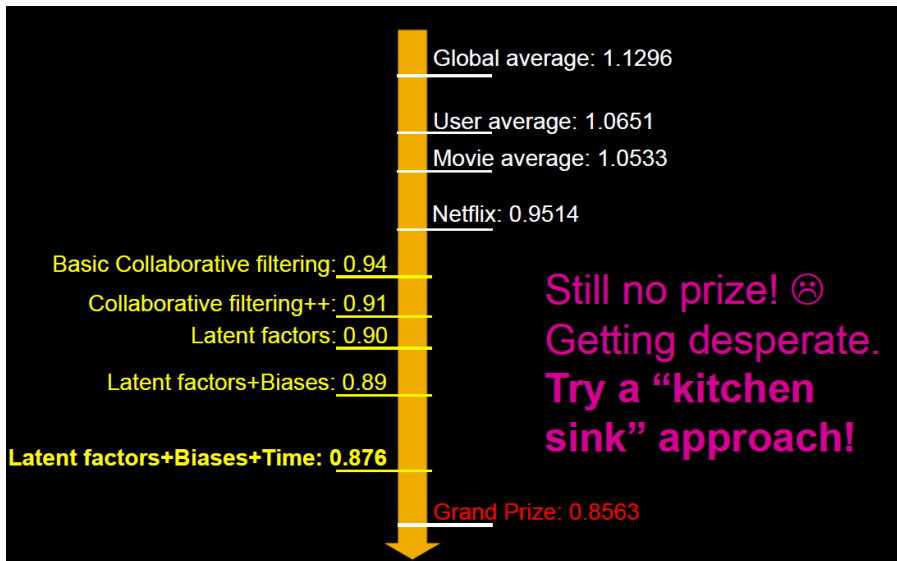$\mu$: overall mean rating, $b_x$: Bias for user x, $b_i$: Bias for movie i

New model:

$$\min_{P,Q,b_x,b_i} \sum_{(i,x)\in \text{training}} (r_{xi} - (\mu + b_x + b_i + q_i \cdot p_x^T))^2$$
$$+ \lambda \left[ \sum_x \|p_x\|_2^2 + \sum_i \|q_i\|_2^2 + \|b_x\|_2^2 + \|b_i\|_2^2 \right]$$

- Both biases $b_x$, $b_i$ as well as interactions $q_i$, $p_x$ are treated as parameters (we estimate them)

- Add time dependence to biases:

$$\hat{r}_{xi} = \mu + b_x(t) + b_i(t) + q_i \cdot p_x^T$$

# Netflix: performance



Global average: 1.1296

User average: 1.0651

Movie average: 1.0533

Netflix: 0.9514

Basic Collaborative filtering: 0.94

Collaborative filtering++: 0.91

Latent factors: 0.90

Latent factors+Biases: 0.89

**Latent factors+Biases+Time: 0.876**

Grand Prize: 0.8563

Still no prize! ☹
Getting desperate.
**Try a "kitchen
sink" approach!**

# Netflix: performance

General matrix completion

# Matrix completion

- Matrix $M \in \mathbb{R}^{n_1 \times n_2}$
- Observe subset of entries
- Can we guess the missing entries?

$$\begin{bmatrix} \times & ? & ? & ? & \times & ? \\ ? & ? & \times & \times & ? & ? \\ \times & ? & ? & \times & ? & ? \\ ? & ? & \times & ? & ? & \times \\ \times & ? & ? & ? & ? & ? \\ ? & ? & \times & \times & ? & ? \end{bmatrix}$$

# Which algorithm ?

Hope: only one low-rank matrix consistent with the sampled entries

Recovery by minimum complexity

$$\begin{aligned} \text{minimize} \quad &\text{rank}(X) \\ \text{subject to} \quad &X_{ij} = M_{ij}, \quad (i,j) \in \Omega \end{aligned}$$

## Problem

- This is NP-hard
- Doubly exponential in $n$ (?)

# SVD - Properties

## Theorem: SVD

If $A$ is a real $m$-by-$n$ matrix, then there exits

$$U = [u_1, \ldots, u_m] \in \mathbb{R}^{m \times m} \text{ and } V = [v_1, \ldots, v_n] \in \mathbb{R}^{n \times n}$$

such that $U^T U = I$, $V^T V = I$ and

$$U^T A V = \operatorname{diag}(\sigma_1, \ldots, \sigma_p) \in \mathbb{R}^{m \times n}, \quad p = \min(m, n),$$

where $\sigma_1 \geq \sigma_2 \geq \ldots \geq \sigma_p \geq 0$.

- Proof: Let $V_1 \in \mathbb{R}^{n \times r}$ has orthonormal columns, then exits $V_2 \in \mathbb{R}^{n \times (n-r)}$ such that $V = [V_1, V_2]$ is orthogonal.

- Let $x \in \mathbb{R}^n$ and $y \in \mathbb{R}^m$ be unit 2-norm vectors: $Ax = \sigma y$ with $\sigma = \|A\|_2$. Then exists $V_2 \in \mathbb{R}^{n \times (n-1)}$ and $U_2 \in \mathbb{R}^{m \times (m-1)}$ so $V = [x, V_2] \in \mathbb{R}^{n \times n}$ and $U = [y, U_2] \in \mathbb{R}^{m \times m}$ are orthogonal.

- Then it can be proved that $U^T A V$ has the following structure

$$U^T A V = \begin{pmatrix} \sigma & w^T \\ 0 & B \end{pmatrix} \equiv A_1.$$

Since

$$\left\| A_1 \begin{pmatrix} \sigma \\ w \end{pmatrix} \right\|_2^2 \geq (\sigma^2 + w^T w)^2,$$

we have $\|A_1\|_2^2 \geq (\sigma^2 + w^T w)$. But $\sigma^2 = \|A\|_2^2 = \|A_1\|_2^2$, and so we must have $w = 0$. An induction gives the proof.

Properties:

- $AV = U\Sigma$, $A^T U = V\Sigma^T$: $Av_i = \sigma u_i$, $A^T u_i = \sigma_i v_i$, $i = 1, \ldots, p$.
- $rank(A) = r$, $null(A) = span\{v_{r+1}, \ldots, v_n\}$, $ran(A) = span\{u_1, \ldots, u_r\}$
- $A = \sum_{i=1}^r \sigma_i u_i v_i^T$
- $\|A\|_F^2 = \sigma_1^2 + \ldots + \sigma_p^2$, $\|A\|_2 = \sigma_1$

# SVD - Best Low Rank Approximation

## Theorem

Let the SVD of $A \in \mathbb{R}^{m \times n}$ be given in Theorem: SVD. If $k < r = rank(A)$ and $A_k = \sum_{i=1}^{k} \sigma_i u_i v_i^T$, then

$$\min_{rank(B)=k} \|A - B\|_2 = \|A - A_k\|_2 = \sigma_{k+1}.$$

- Proof: Since $U^T A_k V = \operatorname{diag}(\sigma_1, \ldots, \sigma_k, 0, \ldots, 0)$ it follows that $rank(A_k) = k$ and $U^T(A - A_k)V = \operatorname{diag}(0, \ldots, 0, \sigma_{k+1}, \ldots, \sigma_p)$. Hence $\|A - A_k\|_2 = \sigma_{k+1}$.

- Suppose $rank(B) = k$ for some $B \in \mathbb{R}^{m \times n}$. We can find orthonormal vectors $x_1, \ldots, x_{n-k}$ so $null(B) = span\{x_1, \ldots, x_{n-k}\}$. A dimension argument shows:

$$span\{x_1, \ldots, x_{n-k}\} \cap span\{v_1, \ldots, v_{k+1}\} \neq \{0\}$$

- Let $z$ be a unit 2-norm vector in this intersection. Since $Bz = 0$ and

$$Az = \sum_{i=1}^{k+1} \sigma_i(v_i^T z)u_i,$$

we have

$$\|A - B\|_2^2 \geq \|(A - B)z\|_2^2 = \|Az\|_2^2 = \sum_{i=1}^{k+1} \sigma_i^2(v_i^T z)^2 \geq \sigma_{k+1}^2$$

Comments:

- So zeroing small $\sigma_i$ introduces less error

- How many $\sigma$s to keep? Rule of thumb: keep 80-90% of 'energy' ($= \sum \sigma_i^2$)

# SVD - Complexity

- To compute SVD: $O(nm^2)$ or $O(n^2m)$

- But:
    - Less work, if we just want singular values
    - or if we want first k singular vectors
    - or if the matrix is sparse

- Implemented in linear algebra packages like
    - Dense matrix: LAPACK
    - Sparse Matrix: ARPACK, PROPACK
    - High Level Software packages: Matlab, SPlus, Mathematica ...

# Relation to Eigen-decomposition

- SVD gives us $A = U\Sigma V^\top$

- Eigen-decomposition: $A = X\Lambda X^\top$
  - $A$ is symmetric
  - $U, V, X$ are orthonormal
  - $\Lambda, \Sigma$ are diagonal

- $AA^\top = U\Sigma\Sigma^\top U^\top$

- $A^\top A = V\Sigma\Sigma^\top V^\top$

- $\lambda_i(A^\top A) = \sigma_i^2(A)$

# Nuclear-norm minimization

Singular value decomposition

$$X = \sum_{k=1}^{r} \sigma_k u_k v_k^*$$

- $\{\sigma_k\}$: singular values, $\{u_k\}, \{v_k\}$: singular vectors

Nuclear norm ($\sigma_i(X)$ is $i$th largest singular value of $X$)

$$\|X\|_* = \sum_{i=1}^{n} \sigma_i(X)$$

Heuristic

$$\begin{aligned} \text{minimize} \quad & \|X\|_* \\ \text{subject to} \quad & X_{ij} = M_{ij}, \quad (i,j) \in \Omega \end{aligned}$$

- Convex relaxation of the rank minimization program

# Connections with compressed sensing

General setup

<table>
<tr><td>Rank minimization</td><td>Convex relaxation</td></tr>
<tr><td>minimize   $\text{rank}(X)$</td><td>minimize   $\|X\|_*$</td></tr>
<tr><td>subject to   $\mathcal{A}(X) = b$</td><td>subject to   $\mathcal{A}(X) = b$</td></tr>
</table>

Suppose $X = \text{diag}(x), x \in \mathbb{R}^n$

- $\text{rank}(X) = \sum_i 1_{(x_i \neq 0)} = \|x\|_{\ell_0}$
- $\|X\|_* = \sum_i |x_i| = \|x\|_{\ell_1}$

<table>
<tr><td>Rank minimization</td><td>Convex relaxation</td></tr>
<tr><td>minimize   $\|x\|_{\ell_0}$</td><td>minimize   $\|x\|_{\ell_1}$</td></tr>
<tr><td>subject to   $Ax = b$</td><td>subject to   $Ax = b$</td></tr>
</table>

This is compressed sensing!

# SOCP/SDP Duality

**(P)** $\min \quad c^\top x$
   s.t. $\quad Ax = b, x_{\mathcal{Q}} \succeq 0$

**(D)** $\max \quad b^\top y$
   s.t. $\quad A^\top y + s = c, s_{\mathcal{Q}} \succeq 0$

**(P)** $\min \quad \langle C, X \rangle$
   s.t. $\quad \langle A_1, X \rangle = b_1$
   $\quad \cdots$
   $\quad \langle A_m, X \rangle = b_m$
   $\quad X \succeq 0$

**(D)** $\max \quad b^\top y$
   s.t. $\quad \sum_i y_i A_i + S = C$
   $\quad S \succeq 0$

**Strong duality**

- If $p^* > -\infty$, (P) is **strictly** feasible, then (D) is feasible and $p^* = d^*$
- If $d^* < +\infty$, (D) is **strictly** feasible, then (P) is feasible and $p^* = d^*$
- If (P) and (D) has **strictly** feasible solutions, then both have optimal solutions.

# Semidefinite program

$$(D) \quad \begin{aligned} \min \quad & -b^\top y \\ \text{s.t.} \quad & y_1 A_1 + \ldots + y_m A_m \preceq C \end{aligned}$$

- $A_i, C \in \mathcal{S}^k$, multiplier is matrix $X \in \mathcal{S}^k$
- Lagrangian $\mathcal{L}(y, X) = -b^\top y + \langle X, y_1 A_1 + \ldots + y_m A_m - C \rangle$
- dual function

$$g(X) = \inf_y \quad \mathcal{L}(y, X) = \begin{cases} -\langle C, X \rangle, & \langle A_i, X \rangle = b_i \\ -\infty & \text{otherwise} \end{cases}$$

The dual of (D) is

$$\begin{aligned} \min \quad & \langle C, X \rangle \\ \text{s.t.} \quad & \langle A_i, X \rangle = b_i, X \succeq 0 \end{aligned}$$

$p^* = d^*$ if primal SDP is strictly feasible.

Suppose unknown matrix $X$ is positive semidefinite

$$
\begin{aligned}
\min \quad & \sum_{i=1}^{n} \sigma_i(X) \\
\text{s.t.} \quad & X_{ij} = M_{ij} \quad (i,j) \in \Omega \\
& X \succeq 0
\end{aligned}
\qquad \Leftrightarrow \qquad
\begin{aligned}
\min \quad & \text{trace}(X) \\
\text{s.t.} \quad & X_{ij} = M_{ij} \quad (i,j) \in \Omega \\
& X \succeq 0
\end{aligned}
$$

Trace heuristic: Mesbahi & Papavassilopoulos (1997), Beck & D'Andrea (1998)

# General SDP formulation

Let $X \in \mathbb{R}^{m \times n}$. For a given norm $\|\cdot\|$, the dual norm $\|\cdot\|_d$ is defined as

$$\|X\|_d := \sup\{\langle X, Y \rangle : Y \in \mathbb{R}^{m \times n}, \|Y\| \leq 1\}$$

Nuclear norm and spectral norms are dual:

$$\|X\| := \sigma_1(X), \quad \|X\|_* = \sum_i \sigma_i(X).$$

$$\text{(P)} \quad \begin{array}{c} \max\limits_{Y} \quad \langle X, Y \rangle \\ \text{s.t. } \|Y\|_2 \leq 1 \end{array} \quad \Leftrightarrow \quad \begin{array}{c} \max\limits_{Y} 2\langle X, Y \rangle \\ \text{s.t. } \begin{bmatrix} I_m & Y \\ Y^\top & I_n \end{bmatrix} \succcurlyeq 0 \end{array} \quad \Leftrightarrow \quad \begin{array}{c} \min\limits_{Z} - \left\langle Z, \begin{bmatrix} 0 & X \\ X^\top & 0 \end{bmatrix} \right\rangle \\ \text{s.t. } Z_1 = I_m \\ Z_2 = I_n \\ Z = \begin{bmatrix} Z_1 & Z_3 \\ Z_3^\top & Z_2 \end{bmatrix} \succeq 0 \end{array}$$

# General SDP formulation

The Lagrangian dual problem is:

$$\max_{W_1, W_2} \min_{Z \succeq 0} \quad -\left\langle Z, \begin{bmatrix} 0 & X \\ X^* & 0 \end{bmatrix} \right\rangle + \langle Z_1 - I_m, W_1 \rangle + \langle Z_2 - I_n, W_2 \rangle$$

*strong duality* after a scaling of $1/2$ and change of variables $X$ to $-X$

$$\text{minimize} \quad \frac{1}{2} \left( \text{trace}(W_1) + \text{trace}(W_2) \right)$$

(D)

$$\text{subject to} \quad \begin{bmatrix} W_1 & X \\ X^\top & W_2 \end{bmatrix} \succcurlyeq 0$$

Optimization variables: $W_1 \in \mathbb{R}^{n_1 \times n_1}, W_2 \in \mathbb{R}^{n_2 \times n_2}$.

Proposition 2.1 in "Guaranteed Minimum-Rank Solutions of Linear Matrix Equations via Nuclear Norm Minimization", Benjamin Recht, Maryam Fazel, Pablo A. Parrilo

# General SDP formulation

Nuclear norm minimization

$$\begin{array}{ll} \min \|X\|_* \\ \text{s.t. } \mathcal{A}(X) = b \end{array} \iff \begin{array}{ll} \max b^\top y \\ \text{s.t. } \|\mathcal{A}^*(y)\| \leq 1 \end{array}$$

SDP Reformulation

$$\begin{array}{ll} \min \dfrac{1}{2}\left(\text{trace}(W_1) + \text{trace}(W_2)\right) \\ \text{s.t. } \mathcal{A}(X) = b \\ \qquad \begin{bmatrix} W_1 & X \\ X^\top & W_2 \end{bmatrix} \succcurlyeq 0 \end{array} \iff \begin{array}{ll} \max b^\top y \\ \text{s.t. } \begin{bmatrix} I & \mathcal{A}^*(y) \\ (\mathcal{A}^*(y))^\top & I \end{bmatrix} \succcurlyeq 0 \end{array}$$

# Matrix recovery

$$M = \sum_{k=1}^{2} \sigma_k u_k u_k^*, \quad \begin{aligned} u_1 &= (e_1 + e_2)/\sqrt{2}, \\ u_2 &= (e_1 - e_2)/\sqrt{2} \end{aligned}$$

$$M = \begin{bmatrix} * & * & 0 & \dots & 0 & 0 \\ * & * & 0 & \dots & 0 & 0 \\ 0 & 0 & 0 & \dots & 0 & 0 \\ \vdots & \vdots & \vdots & \vdots & \vdots & \vdots \\ 0 & 0 & 0 & \dots & 0 & 0 \end{bmatrix}$$

Cannot be recovered from a small set of entries

Rank-1 matrix $M = xy^*$

$$M_{ij} = x_i y_j$$

$$\begin{bmatrix} \times & \times & \times & \times & \times & \times \\ & & & & & \\ \times & \times & \times & \times & \times & \times \\ \times & \times & \times & \times & \times & \times \\ \times & \times & \times & \times & \times & \times \\ \times & \times & \times & \times & \times & \times \end{bmatrix}$$

If single row (or column) is not sampled $\rightarrow$ recovery is not possible

*What happens for almost all sampling sets?*

$\Omega$ subset of $m$ entries selected uniformly at random

# References

- Jianfeng Cai, Emmanuel Candes, Zuowei Shen, *Singular value thresholding algorithm for matrix completion*
- Shiqian Ma, Donald Goldfarb, Lifeng Chen, *Fixed point and Bregman iterative methods for matrix rank minimization*
- Zaiwen Wen, Wotao Yin, Yin Zhang, *Solving a low-rank factorization model for matrix completion by a nonlinear successive over-relaxation algorithm*
- Onureena Banerjee, Laurent El Ghaoui, Alexandre d'Aspremont, *Model Selection Through Sparse Maximum Likelihood Estimation for Multivariate Gaussian or Binary Data*
- Zhaosong Lu, *Smooth optimization approach for sparse covariance selection*

# Matrix Rank Minimization

Given $X \in \mathbb{R}^{m \times n}$, $\mathcal{A} : \mathbb{R}^{m \times n} \to \mathbb{R}^p$, $b \in \mathbb{R}^p$, we consider

- the matrix rank minimization problem:

$$\min \ \text{rank}(X), \ \text{s.t. } \mathcal{A}(X) = b$$

- matrix completion problem:

$$\min \ \text{rank}(X), \ \text{s.t. } X_{ij} = M_{ij}, (i,j) \in \Omega$$

- nuclear norm minimization:

$$\min \ \|X\|_* \ \text{s.t. } \mathcal{A}(X) = b$$

where $\|X\|_* = \sum_i \sigma_i$ and $\sigma_i = i$th singular value of matrix $X$.

# Quadratic penalty framework

- Unconstrained Nuclear Norm Minimization:

$$\min \ F(X) := \mu \|X\|_* + \frac{1}{2} \|\mathcal{A}(X) - b\|_2^2.$$

- Optimality condition:

$$\mathbf{0} \in \mu \partial \|X^*\|_* + \mathcal{A}^*(\mathcal{A}(X^*) - b),$$

where $\partial \|X\|_* = \{UV^\top + W : U^\top W = 0, WV = 0, \|W\|_2 \le 1\}$.

- Linearization approach ($g$ is the gradient of $\frac{1}{2}\|\mathcal{A}(X) - b\|_2^2$):

$$
\begin{aligned}
X^{k+1} &:= \ \arg\min_X \ \mu\|X\|_* + \left\langle g^k, X - X^k \right\rangle + \frac{1}{2\tau}\|X - X^k\|_F^2 \\
&= \ \arg\min_X \ \mu\|X\|_* + \frac{1}{2\tau}\|X - (X^k - \tau g^k)\|_F^2
\end{aligned}
$$

# Matrix Shrinkage Operator

For a matrix $Y \in \mathbb{R}^{m \times n}$, consider:

$$\min_{X \in \mathbb{R}^{m \times n}} \nu \|X\|_* + \frac{1}{2}\|X - Y\|_F^2.$$

The optimal solution is:

$$X := S(Y, \nu) = U\mathrm{Diag}(s(\sigma, \nu))V^\top,$$

- SVD: $Y = U\mathrm{Diag}(\sigma)V^\top$
- Thresholding operator:

$$s(x, \nu) := \bar{x}, \text{ with } \bar{x}_i = \begin{cases} x_i - \nu, & \text{if } x_i - \nu > 0 \\ 0, & \text{o.w.} \end{cases}$$

# Fixed Point Method (Proximal gradient method)

Fixed Point Iterative Scheme

$$\begin{cases} Y^k = X^k - \tau \mathcal{A}^*(\mathcal{A}(X^k) - b) \\ X^{k+1} = S(Y^k, \tau \mu). \end{cases}$$

**Lemma:** Matrix shrinkage operator is non-expansive. i.e.,

$$\|S(Y_1, \nu) - S(Y_2, \nu)\|_F \leq \|Y_1 - Y_2\|_F.$$

Complexity of the fixed point method:

$$F(X^k) - F(X^*) \leq \frac{L_f \|X^0 - X^*\|^2}{2k}$$

# Accelerated proximal gradient (APG) method

APG algorithm ($t^{-1} = t^0 = 1$):

$$
\begin{aligned}
Y^k &= X^k + \frac{t^{k-1} - 1}{t^k}(X^k - X^{k-1}) \\
G^k &= Y^k - (\tau^k)^{-1}\mathcal{A}^*(\mathcal{A}(Y^k) - b) \\
X^{k+1} &= S_{\tau^k}(G^k), \quad t^{k+1} = \frac{1 + \sqrt{1 + 4(t^k)^2}}{2}
\end{aligned}
$$

Complexity:

$$
F(X^k) - F(X^*) \leq \frac{2L_f\|X^0 - X^*\|^2}{(k+1)^2}
$$

# SVT

Linearized Bregman method:

$$V^{k+1} := V^k - \tau \mathcal{A}^*(\mathcal{A}(X^k) - b)$$
$$X^{k+1} := S_{\tau\mu}(V^{k+1})$$

Convergence to

$$\min \ \tau\|X\|_* + \frac{1}{2}\|X\|_F^2, \ \text{s.t.} \ \mathcal{A}(X) = b$$

# Review of Bregman method

Consider the problem:

$$\min \ \|X\|_*, \ \text{s.t.} \ \mathcal{A}(X) = b$$

Bregman method:

- $D^{P^k}(X, X^k) := \|X\|_* - \|X^k\|_* - \langle P^k, X - X^k \rangle$
- $X^{k+1} := \arg\min_X \ \mu D^{P^k}(X, X^k) + \frac{1}{2}\|\mathcal{A}(X) - b\|_2^2$
- $P^{k+1} = P^k + \frac{1}{\mu}\mathcal{A}^\top(b - \mathcal{A}(X^{k+1}))$

Augmented Lagrangian (updating multiplier or $b$):

- $X^{k+1} := \arg\min_X \ \mu\|X\|_* + \frac{1}{2}\|\mathcal{A}(X) - b^k\|_2^2$
- $b^{k+1} = b + (b^k - \mathcal{A}(X^{k+1}))$

They are equivalent, see Yin-Osher-Goldfarb-Darbon

# Linearized approaches

Linearized Bregman method:

$$
\begin{aligned}
X^{k+1} &:= \arg\min \ \mu D^{p^k}(X, X^k) + \left\langle \mathcal{A}^\top(\mathcal{A}(X^k) - b), X - X^k \right\rangle + \frac{1}{2\delta}\|X - X^k\|_F^2, \\
P^{k+1} &:= P^k - \frac{1}{\mu\delta}(X^{k+1} - X^k) - \frac{1}{\mu}\mathcal{A}^\top(\mathcal{A}(X^k) - b),
\end{aligned}
$$

which is equivalent to

$$
\begin{aligned}
X^{k+1} &:= \arg\min \ \mu\|X\|_* + \frac{1}{2\delta}\|X - V^k\|_F^2 \\
V^{k+1} &:= V^k - \delta\mathcal{A}^\top(\mathcal{A}(X^{k+1}) - b)
\end{aligned}
$$

Bregmanized operator splitting:

$$
\begin{aligned}
X^{k+1} &:= \arg\min \ \mu\|X\|_* + \left\langle \mathcal{A}^\top(\mathcal{A}(X^k) - b^k), X - X^k \right\rangle + \frac{1}{2\delta}\|X - X^k\|_F^2 \\
b^{k+1} &= b + (b^k - \mathcal{A}(X^{k+1}))
\end{aligned}
$$

Are they equivalent?

# Linearized approaches

Linearized Bregman method:

$$X^{k+1} \quad := \quad \arg\min \ \mu D^{p^k}(X, X^k) + \left\langle \mathcal{A}^\top(\mathcal{A}(X^k) - b), X - X^k \right\rangle + \frac{1}{2\delta}\|X - X^k\|_F^2,$$

$$P^{k+1} \quad := \quad P^k - \frac{1}{\mu\delta}(X^{k+1} - X^k) - \frac{1}{\mu}\mathcal{A}^\top(\mathcal{A}(X^k) - b),$$

which is equivalent to

$$X^{k+1} := \mathcal{S}(V^k, \mu\delta)$$
$$V^{k+1} := V^k - \delta\mathcal{A}^\top(\mathcal{A}(X^{k+1}) - b)$$
or
$$X^{k+1} := \mathcal{S}(\delta A^\top(b^k), \mu\delta)$$
$$b^{k+1} := b + (b^k - \mathcal{A}(X^{k+1}))$$

Bregmanized operator splitting:

$$X^{k+1} \quad := \quad \mathcal{S}(X^k - \delta(\mathcal{A}^\top(\mathcal{A}(X^k) - b^k)), \mu\delta) = \mathcal{S}(\delta\mathcal{A}^\top(b^k) + X^k - \delta\mathcal{A}^\top(\mathcal{A}(X^k)), \mu\delta)$$

$$b^{k+1} \quad = \quad b + (b^k - \mathcal{A}(X^{k+1}))$$

# Low-rank factorization model

- Finding a low-rank matrix $W$ so that $\|\mathcal{P}_\Omega(W - M)\|_F^2$ or the distance between $W$ and $\{Z \in \mathbb{R}^{m \times n}, Z_{ij} = M_{ij}, \forall(i,j) \in \Omega\}$ is minimized.
- Any matrix $W \in \mathbb{R}^{m \times n}$ with $\operatorname{rank}(W) \leq K$ can be expressed as $W = XY$ where $X \in \mathbb{R}^{m \times K}$ and $Y \in \mathbb{R}^{K \times n}$.

## New model

$$\min_{X,Y,Z} \frac{1}{2} \|XY - Z\|_F^2 \ \text{ s.t. } \ Z_{ij} = M_{ij}, \forall(i,j) \in \Omega$$

- Advantage: SVD is no longer needed!
- Related work: the solver `OptSpace` based on optimization on manifold

# Nonlinear Gauss-Seidel scheme

First variant of alternating minimization:

$$
\begin{aligned}
X_+ &\leftarrow ZY^\dagger \equiv ZY^\top (YY^\top)^\dagger, \\
Y_+ &\leftarrow (X_+)^\dagger Z \equiv (X_+^\top X_+)^\dagger (X_+^\top Z), \\
Z_+ &\leftarrow X_+ Y_+ + \mathcal{P}_\Omega(M - X_+ Y_+).
\end{aligned}
$$

Let $\mathcal{P}_A$ be the orthogonal projection onto the range space $\mathcal{R}(A)$

- $X_+ Y_+ = \left( X_+ (X_+^\top X_+)^\dagger X_+^\top \right) Z = \mathcal{P}_{X_+} Z$
- One can verify that $\mathcal{R}(X_+) = \mathcal{R}(ZY^\top)$ .
- $X_+ Y_+ = \mathcal{P}_{ZY^\top} Z = ZY^\top (YZ^\top ZY^\top)^\dagger (YZ^\top) Z.$
- idea: modify $X_+$ or $Y_+$ to obtain the same product $X_+ Y_+$

# Nonlinear Gauss-Seidel scheme

Second variant of alternating minimization:

$$\begin{aligned}
X_+ &\leftarrow ZY^\top, \\
Y_+ &\leftarrow (X_+)^\dagger Z \equiv (X_+^\top X_+)^\dagger (X_+^\top Z), \\
Z_+ &\leftarrow X_+Y_+ + \mathcal{P}_\Omega(M - X_+Y_+).
\end{aligned}$$

Third variant of alternating minimization: $V = \mathtt{orth}(ZY^\top)$

$$\begin{aligned}
X_+ &\leftarrow V, \\
Y_+ &\leftarrow V^\top Z, \\
Z_+ &\leftarrow X_+Y_+ + \mathcal{P}_\Omega(M - X_+Y_+).
\end{aligned}$$
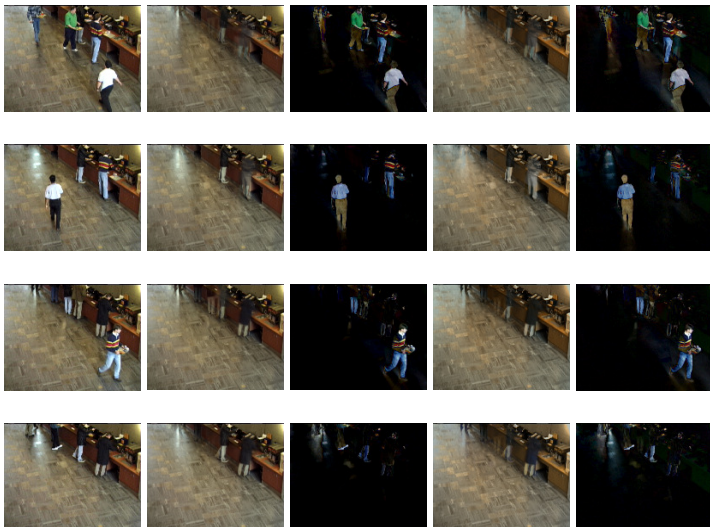
# Sparse and low-rank matrix separation

- Given a matrix $M$, we want to find a low rank matrix $W$ and a sparse matrix $E$, so that $W + E = M$.
- Convex approximation:

$$\min_{W,E} \ \|W\|_* + \mu\|E\|_1, \ \text{s.t.} \ W + E = M$$

- Robust PCA

# Video separation

- Partition the video into moving and static parts

# ADMM

Convex approximation:

$$\min_{W,E} \|W\|_* + \mu\|E\|_1, \text{ s.t. } W + E = M$$

Augmented Lagrangian function:

$$L(W,E,\Lambda) := \|W\|_* + \mu\|E\|_1 + \langle \Lambda, W + E - M \rangle + \frac{1}{2\beta}\|W + E - M\|_F^2$$

Alternating direction Augmented Lagrangian method

$$
\begin{aligned}
W^{j+1} &:= \arg\min_W L(W,\ E^j,\ \Lambda^j), \\
E^{j+1} &:= \arg\min_E L(W^{j+1},\ E,\ \Lambda^j), \\
\Lambda^{j+1} &:= \Lambda^j + \frac{\gamma}{\beta}(W^{j+1} + E^{j+1} - M).
\end{aligned}
$$

# W-subproblem

Convex approximation:

$$
\begin{aligned}
W^{j+1} &:= \arg\min_W \, L(W, \, E^j, \, \Lambda^j) \\
&= \arg\min_W \|W\|_* + \frac{1}{2\beta} \left\| W - \left( M - E^j - \beta\Lambda^j \right) \right\|_F^2 \\
&= S_\beta(M - E^j - \beta\Lambda^j) := U\mathrm{Diag}(s_\beta(\sigma))V^\top
\end{aligned}
$$

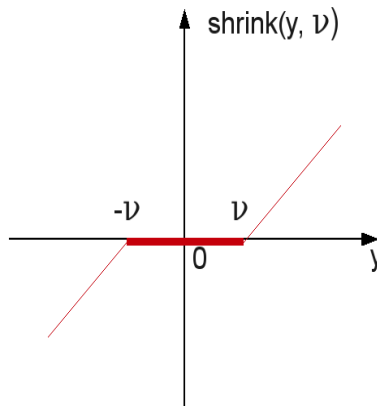- SVD: $M - E^j - \beta\Lambda^j = U\mathrm{Diag}(\sigma)V^\top$
- Thresholding operator:

$$
s_\nu(x) := \bar{x}, \text{ with } \bar{x}_i = \left\{ \begin{array}{ll} x_i - \nu, & \text{if } x_i - \nu > 0 \\ 0, & \text{o.w.} \end{array} \right.
$$

# E-subproblem

Convex approximation:

$$
\begin{aligned}
W^{j+1} &:= \arg\min_E \ L(W^{j+1}, \ E, \ \Lambda^j) \\
&= \arg\min_E \|E\|_1 + \frac{1}{2\beta\mu} \left\| E - \left(M - W^{j+1} - \beta\Lambda^j\right) \right\|_F^2 \\
&= s_{\beta\mu}(M - W^{j+1} - \beta\Lambda^j)
\end{aligned}
$$

$$
\begin{aligned}
s_\nu(y) : &= \arg\min_{x\in\mathbb{R}} \ \nu\|x\|_1 + \frac{1}{2}\|x - y\|_2^2 \\
&= \begin{cases} y - \nu\,\mathsf{sgn}(y), & \text{if } |y| > \nu \\ 0, & \text{otherwise} \end{cases}
\end{aligned}
$$



shrink(y, $\nu$)

# Low-rank factorization model for matrix separation

- Consider the model

$$\min_{Z,S} \ \|S\|_1 \ \text{s.t.} \ Z + S = D, \ \mathrm{rank}(Z) \leq K$$

- Low-rank factorization: $Z = UV$

$$\min_{U,V,Z} \ \|Z - D\|_1 \ \text{s.t.} \ UV - Z = 0$$

- Only the entries $D_{ij}$, $(i,j) \in \Omega$, are given. $\mathcal{P}_\Omega(D)$ is the projection of $D$ onto $\Omega$.

## New model

$$\min_{U,V,Z} \ \ \|\mathcal{P}_\Omega(Z - D)\|_1 \ \ \text{s.t.} \ \ UV - Z = 0$$

- Advantage: SVD is no longer needed!

## ADMM

Consider:

$$\min_{U,V,Z} \quad \|\mathcal{P}_\Omega(Z - D)\|_1 \quad \text{s.t.} \quad UV - Z = 0$$

Introduce the augmented Lagrangian function

$$\mathcal{L}_\beta(U, V, Z, \Lambda) = \|\mathcal{P}_\Omega(Z - D)\|_1 + \langle \Lambda, UV - Z \rangle + \frac{\beta}{2}\|UV - Z\|_F^2,$$

Alternating direction augmented Lagrangian framework (Bregman):

$$
\begin{aligned}
U^{j+1} &:= \arg\min_{U \in \mathbb{R}^{m \times k}} \mathcal{L}_\beta(U, \ V^j, \ Z^j, \ \Lambda^j), \\
V^{j+1} &:= \arg\min_{V \in \mathbb{R}^{k \times n}} \mathcal{L}_\beta(U^{j+1}, \ V, \ Z^j, \ \Lambda^j), \\
Z^{j+1} &:= \arg\min_{Z \in \mathbb{R}^{m \times n}} \mathcal{L}_\beta(U^{j+1}, \ V^{j+1}, \ Z, \ \Lambda^j), \\
\Lambda^{j+1} &:= \Lambda^j + \gamma\beta(U^{j+1}V^{j+1} - Z^{j+1}).
\end{aligned}
$$

# ADMM subproblems

- Let $B = Z - \Lambda/\beta$, then

$$U_+ = BV^\top(VV^\top)^\dagger \text{ and } V_+ = (U_+^\top U_+)^\dagger U_+^\top B$$

Since $U_+ V_+ = U_+(U_+^\top U_+)^\dagger U_+^\top B = \mathcal{P}_{U_+} B$, then:

$$Q := \texttt{orth}(BV^\top), \quad U_+ = Q \text{ and } V_+ = Q^\top B$$

- Variable $Z$:

$$
\begin{aligned}
\mathcal{P}_\Omega(Z_+) &= \mathcal{P}_\Omega\left(\mathcal{S}\left(U_+V_+ - D + \frac{\Lambda}{\beta}, \frac{1}{\beta}\right) + D\right) \\
\mathcal{P}_{\Omega^c}(Z_+) &= \mathcal{P}_{\Omega^c}\left(U_+V_+ + \frac{\Lambda}{\beta}\right)
\end{aligned}
$$