# CIL

## MATRIX APPROXIMATION & RECONSTRUCTION

## Problem Setup

Given a matrix $A \in \mathbb{R}^{m \times m}$ with observed entries $\mathcal{I} \subseteq [m] \times [m]$

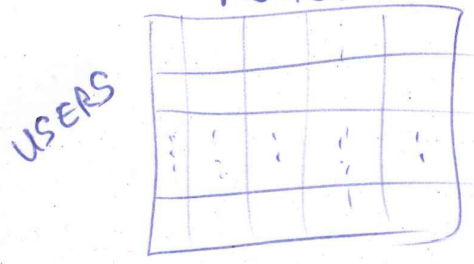Goal: Fill in the unobserved entries.

What makes it possible?

↳ Assume: there exist representations of the rows and the columns that require less than $m \times m$ parameters. "collaborative filtering"
(in other words: "there is something to learn")

## Examples

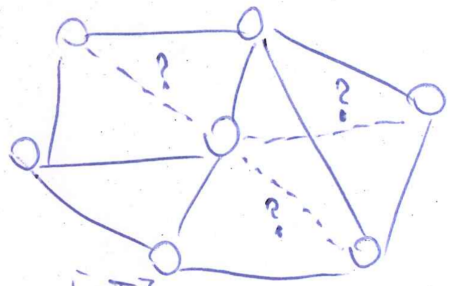① recommender systems

e.g. the Netflix contest

MOVIES

USERS



$\mathcal{I} = \{(i,j) \mid \text{user } i \text{ rated movie } j\}$

② sensor localisation

→ $m$ sensors: points $x_i \in \mathbb{R}^3$, $i=1,\dots,m$

→ observe partial info about pairwise distances

$$D_{ij} = \|x_i - x_j\|^2 =$$
$$= \|x_i\|^2 + \|x_j\|^2 - 2x_i^T x_j$$
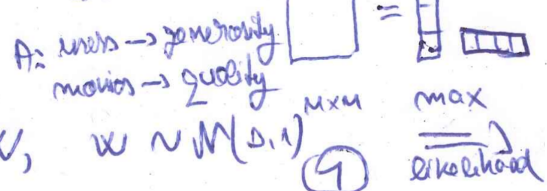
→ want to infer distances between all pairs of points.
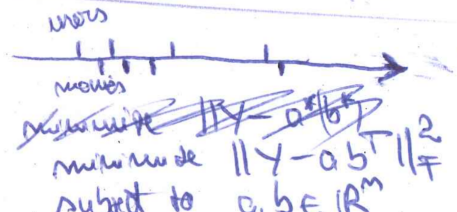


$X = \begin{bmatrix} x_1^T \\ \vdots \\ x_m^T \end{bmatrix} \in \mathbb{R}^{m \times 3}$, $\quad d = [\|x_1\|^2, \dots, \|x_m\|^2]^T$

$e = [1, \dots, 1]^T$

$\Rightarrow D = de^T + ed^T - 2XX^T$

low rank (!)

## Objectives

① minimize $\|A - X\|_{\mathcal{I}}^2$     (NP-)hard!
$X \in \mathbb{R}^{m \times m}$
subject to $\operatorname{rank}(X) \leq k$     "Is there a solution that achieves loss $C$?"

Notation: $\|X\|_{\mathcal{I}}^2 = \sum_{(i,j) \in \mathcal{I}} x_{ij}^2$     Low-rank matrix recovery

② minimize $\operatorname{rank}(X)$     Exact matrix recovery
$X$
subject to $\|A - X\|_{\mathcal{I}} = 0$

"simplest explanation"     Q: Interpretation for the case $k=1$?

## Statistical models?

Single-spike model     A: users → generosity
movies → quality

$Y = a^*(b^*)^T + W, \quad W \sim \mathcal{N}(0,1)$



$m \times m$     max
     likelihood

minimize $\|Y - ab^T\|_F^2$
subject to $a, b \in \mathbb{R}^m$

Recall

If $\mathcal{I} = [m] \times [m]$ (full matrix observed)

Thm problem ④ has a closed-form solution:

$$
\begin{aligned}
&\text{with } A = U\Delta V^T \text{ and } A_k = \sum_{i=1}^{k} \sigma_i u_i v_i^T \\
&\phantom{\text{with }} \underset{\min(m,n)}{} = \sum_{i=1}^{\min(m,n)} \sigma_i u_i v_i^T \quad (\text{for } k \leq \text{rank}(A)) \\
&\text{thm } \underset{\text{rank}(x)=k}{\min} \ \|A - x\|_F^2 = \|A - A_k\|_F^2 = \sum_{n=k+1}^{\text{rank}(A)} \sigma_n^2
\end{aligned}
$$

(Eckart-Young Theorem)

Q: How many degrees of freedom in a rank-$n$ matrix $\in \mathbb{R}^{m \times n}$?

$$nm + (m-n)n$$

---

Back to problem ①

$$
\underset{x:\ \text{rank}(x) \leq k}{\min} \ \|A - x\|_2^2
$$

~~$\{ \text{rank}(A), \text{rank}(B) \}$~~ $\text{rank}(A+B) \leq \text{rank}(A) + \text{rank}(B)$

Attempt I: Optimize approximately with gradient descent + projections.
  Can we project onto $\{x : \text{rank}(x) \leq k\}$ ? (Yes, see above)
  In practice the steps are too chaotic.  (Should round towclian: PCA)

Attempt II: Manifold optimization: take steps "on the manifold".

Attempt III: Reparametrise:

$$
\begin{aligned}
&\text{minimize } \|A - U^T V\|_{\mathcal{I}}^2 \\
&U \in \mathbb{R}^{k \times m} \\
&V \in \mathbb{R}^{k \times m}
\end{aligned}
$$

$$U = \begin{bmatrix} | & | & & | \\ u_1 & u_2 & \cdots & u_m \\ | & | & & | \end{bmatrix} \in \mathbb{R}^{k \times m}$$

$$V = \begin{bmatrix} | & | & & | \\ v_1 & v_2 & \cdots & v_m \\ | & | & & | \end{bmatrix} \in \mathbb{R}^{k \times m}$$
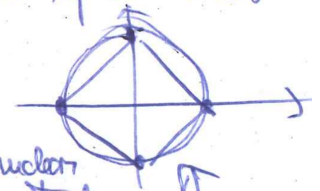
→ Regularised loss function:

$$
L(U,V) = \sum_{(i,j) \in \mathcal{I}} (a_{ij} - u_i^T v_j)^2 + \lambda \sum_{i=1}^{m} \|u_i\|^2 + \lambda \sum_{j=1}^{m} \|v_j\|^2
$$

why regularise?
  ↳ Avoid overfitting. For instance, all user and items representations
    are constrained to a ball:   $\|U\|_F \leq R$
    (similar to ridge regression)   $\|V\|_F \leq R$



Could try ~~sp~~ nuclear norm regularizer for sparsity!

→ "Solve" it optimistically with gradient descent?  Fix one, backprop for the other?
  ↳ Notice that we have closed form solution when one argument is fixed.

②

# Solution problem ②

① $L(U,V)$ is not convex!

Counter example: $M = M = 1$
$u \neq 0, v \neq 0$

$L(u,v) = (a - uv)^2 + \lambda u^2 + \lambda v^2$

~~Suffices~~ to Show $L'(u,v) = (a - uv)^2$
is not convex. (for simplicity) $\lambda = a$

$$\nabla^2 L'(u,v) = \begin{bmatrix} 2V^2 & 4uV \\ 4uv - 2a & 2u^2 \end{bmatrix}$$

(can, in higher dimensions,

$$U = \begin{bmatrix} u & 0 & \cdots & 0 \\ 0 & 0 & & 0 \\ \vdots & \vdots & & \vdots \\ 0 & 0 & \cdots & 0 \end{bmatrix} ; \quad V = \begin{bmatrix} v & 0 & \cdots & 0 \\ 0 & 0 & & 0 \\ \vdots & \vdots & & \\ 0 & 0 & \cdots & 0 \end{bmatrix})$$

__First__: check Twice - differentiability.

$(4uv - 2a)^2 = 4[-3(uv)^2 + 4a(uv) - a^2] < 0$ for

$|\nabla^2 L'(u,v)| = 4u^2 v^2 - \cancel{16u^2v^2} = \cancel{-12u^2v^2} < 0$

$uv \notin [\frac{a}{3}, a]$

Sylvester's criterion: a symmetric matrix is positive definite if and only if all leading principal minors have positive determinant.

In our case: $|2v| = 2v^2 > 0$ but $|\nabla^2 L'(u,v)| < 0$
(as a 1×1 matrix) (whole matrix)

⟹ $L(u,v)$ is not convex.

② $g(u) = L(U,V)$ is convex (think: ridge regression)

③ Update rule for $u_i$.

$$\frac{\partial L(U,V)}{\partial u_i} = -2 \sum_{j:(i,j) \in \mathcal{I}} (a_{ij} - u_i^T v_j) v_j + 2\lambda u_i \overset{!}{=} 0$$

$$\Rightarrow \sum_{j:(i,j) \in \mathcal{I}} a_{ij} v_j = \sum_{j:(i,j) \in \mathcal{I}} (u_i^T v_j) v_j + \lambda u_i =$$

$$= \sum v_j (v_j^T u_i) + \lambda u_i =$$

$$= (\sum v_j v_j^T + \lambda I_K) u_i$$

Can easily be used to embed new users and movies.

$$\Rightarrow u_i = \Big( \sum_{j:(i,j) \in \mathcal{I}} v_j v_j^T + \lambda I_K \Big)^{-1} \sum_{j:(i,j) \in \mathcal{I}} a_{ij} v_j$$

④ Complexity: $\underbrace{n_i K^2 + \boxed{K^3}}_{\text{inverse}} + \underbrace{n_i K + \boxed{K^2}}_{\text{③ product}} \Rightarrow O(n_i K^2 + K^3)$

(5) The update rule looks like

$$\mu_i = \left(V_{j \ni j}^T V_{j \ni j}\right)^{-1} V_{j \ni j}^T \, a_{i(j)}$$

$$V_{j \ni j} = \begin{bmatrix} - & V_{j(1)} & - \\ - & V_{j(2)} & - \\ & \vdots & \\ - & V_{j(m_i)} & - \end{bmatrix} \in \mathbb{R}^{m_i \times K}$$

if $\lambda = 0$, which is the solution to OLS.

Interpretation: • given the representation of the items $\{y_j\}$, compute the representation of the users independently, based on their ~~ratings~~ "ratings".

• the "ratings" $\{a_{ij}\}_j$ are orthogonally projected on the column space of $V_{j \ni j}$.

• the user representation $u_i$ is given by the "indices" in the orthogonal projection.

$$\left( \text{Recall } \hat{\beta}_{LS} = (X^T X)^{-1} X^T Y, \quad X \hat{\beta}(y) = Id_X Y, \quad \text{where } Id_X = X(X^T X)^{-1} X^T \text{ is} \right)$$
the orthogonal projection on the column space of $X$.

$\to$ How to use the representations?

$$a_{pq} = u_p^T v_q \quad \text{for } (p,q) \notin \mathcal{I}$$

Bonus: we have representations of both "movies" and "users" in the same space, so we can compute similarities between movies, distances between users and movies etc.

---

Next time

$$\text{rank}(X) = \dim(\text{span}(X)) = \#\{\sigma_i > 0\}$$
should remind of the "0-norm"

$$\min_X \text{rank}(X)$$
$$\text{s.t. } \|A - X\|_{\mathcal{I}} = 0$$

$\hookrightarrow$ gets relaxed to:

$$\min \|X\|_*$$
$$\text{s.t. } \|A - X\|_{\mathcal{I}} = 0$$

$\nearrow$ nuclear norm $\|X\|_* = \sum_i \sigma_i(X)$ $\nwarrow$ should remind of the L1-norm

these are proper norms



best convex approximation.

(4)

# Problem ①

$\begin{cases}\text{maximize} \quad \langle x, Ax \rangle = x^T A x =: f(x) \quad \text{with } A - \text{an } m \times m \text{ symmetric matrix} \\ \text{s.t.} \quad \|x\|^2 - 1 = 0 \\ \quad g(x) \end{cases}$

① $\quad v_1 = \arg\max\limits_{\|x\|^2-1=0} x^T A x \quad$ and $\quad \boxed{f(v_1) = \lambda_1}$

**Prove** $Av_1 = \lambda_1 v_1$.

Lagrange multiplier theory: $\exists \lambda \in \mathbb{R}$ such that $\nabla(f - \lambda g)(v_1) = 0$

We have: $\nabla g(x) = 2x$

$\nabla f(x) = 2Ax$

Thus: $2Av_1 = 2\lambda v_1 \implies Av_1 = \lambda v_1 \mid v_1^T(\cdot) \implies v_1^T A v_1 = \lambda v_1^T v_1 \implies$

$\implies f(v_1) = \lambda \quad$ (we also know $\boxed{f(v_1) = \lambda}$)

Hence: $Av_1 = \lambda_1 v_1$.

② $\quad v_2 = \arg\max\limits_{\begin{cases}g(x):=\|x\|^2-1=0 \\ h(x):=x^T v_1 = 0\end{cases}} f(x) \quad$ and $\quad \boxed{f(v_2)=\lambda_2}$

$S^{m-2}$ $\left(\begin{array}{l}\text{the } (m-2)\text{-dimensional} \\ \text{unit sphere}\end{array}\right)$

**Prove:** $Av_2 = \lambda_2 v_2$

As before: $\exists \lambda, \mu \in \mathbb{R}$ s.t. $\nabla(f - \lambda g - \mu h)(v_2) = 0$

$\nabla h(x) = v_1$

$\nabla g(x) = 2x$

$\nabla f(x) = 2Ax$

So we have: $2Av_2 = 2\lambda v_2 + \mu v_1 \mid v_1^T(\cdot) \implies$

$\implies 2v_1^T A v_2 = 2\lambda \underbrace{v_1^T v_2}_{?} + \mu \underbrace{v_1^T v_1}_{1}$

$\underset{x^Ty=y^Tx}{\iff} 2v_2^T A^T v_1 = \mu$

$\underset{A=A^T}{\iff} 2v_2^T A v_1 = \mu$

$\iff 2v_2^T v_1 = \mu$

$\iff 0 = \mu$

$\implies Av_2 = \lambda v_2 \mid \cdot v_2^T(\cdot)$

$\implies f(v_2) = \lambda \left(\begin{array}{l}\text{we also} \\ \text{have } f(v_2)=\lambda_2\end{array}\right)$

$\implies Av_2 = \lambda_2 v_2$

⑤

③  $v_3 = \arg\max\limits_{\substack{g(x)=\|x\|^2-1=0}} f(x)$   and   $f(v_3) = \lambda_3$

$h(x) = x^T v_1 = 0$      Prove: $\boxed{A v_3 = \lambda_3 v_3}$

$l(x) = x^T v_2 = 0$

Introduce a new Lagrange multiplier:

$$\nabla(f - \lambda g - \mu h - \vartheta l)(v_3) = 0$$

which yields

$$2 A v_3 = 2\lambda v_3 + \mu v_1 + \vartheta v_2$$

$\Rightarrow$ multiply by $v_1$ $\Rightarrow$ $\mu = 0$

$\Rightarrow$ multiply by $v_2$ $\Rightarrow$ $\vartheta = 0$

$\Rightarrow$ $f(v_3) = \lambda$ $\boxed{f(v_3) = \lambda_3}$ $\Rightarrow$ $A v_3 = \lambda_3 v_3$.

$\Rightarrow A v_3 = \lambda v_3 \;|\; v_3^T(\cdot)$

<u>NOTE</u>: By construction: $\lambda_1 \geq \lambda_2 \geq \lambda_3$.

④  Iterate the above procedure:

$$\langle v_j, A v_m \rangle = \langle A v_j, v_m \rangle = \lambda_j \langle v_j, v_m \rangle = 0, \qquad \forall j = 1, \dots, m-1$$

So   $A v_m = \lambda_m v_m$.

The vectors $\{v_1, \dots, v_m\}$ form an orthonormal basis

and $\lambda_1 \geq \lambda_2 \geq \dots \geq \lambda_m$.

⑤  In PCA, this is applied on the ~~empirical~~ empirical covariance

matrix $\dfrac{1}{m} \sum\limits_{i=1}^{m} x_i x_i^T = \Sigma$