

Sparse Matrix Decomposition

Until now, we've been focusing on "simple" vectors:

$$\min_x \|y - x\|^2 \text{ s.t. } x \in \mathcal{C} : \{\text{"set of simple vectors"}\}$$
$$\mathcal{C} : \{x : \|x\|_0 \leq k\}.$$

What if now we formulate this for matrices?

$$\min_M \|Z - M\|_F^2 \text{ s.t. } M \in \mathcal{C}$$

- option 1: sparse matrices $\Rightarrow \mathcal{C} : \{M : \|M\|_0 \leq k\}$.

$$\Rightarrow \min_M \|Z - M\|_F^2 \text{ s.t. } \|M\|_0 \leq k$$

or

$$(P_1): \min_M \|Z - M\|_F^2 + \lambda \|M\|_1$$

$$\hat{M} = S_\lambda(Z) \Rightarrow \text{sign}(Z_{ij}) (|Z_{ij}| - \lambda)_+ = M_{ij}.$$

Option 2: Just like vectors, they aren't sparse themselves, but rather in some other domain. We will use some "natural decomposition".

Consider $X \in \mathbb{R}^{D \times N}$, and we want

$$X \approx \sum_{i=1}^r A_i \cdot c_i ; A_i \in \mathbb{R}^{D \times N} ; c_i \in \mathbb{R}.$$

"Basic" because $A_i = (u_i \cdot v_i^T)$, $u_i \in \mathbb{R}^D$, $v_i \in \mathbb{R}^N$ "basic" matrices.

$\Rightarrow A_i$: rank 1 matrices.

Thus, let's say our "simplicity" pursuit problem becomes:

$$\min_{A_i, c_i} \left\| X - \sum_{i=1}^d c_i A_i \right\|_F^2 \quad \text{s.t.} \quad \text{rank}(A_i) = 1 \\ d \leq \min(N, D).$$

This is the same as

$$\min_A \left\| X - A \right\|_F^2 \quad \text{s.t.} \quad \text{rank}(A) \leq d.$$

Recall that $\forall X: X = U \Sigma V^T$; Σ : diagonal $D \times D$.

and $\text{rank}(X) = \|\Sigma\|_0$. Thus,

$$\boxed{\min_A \left\| X - A \right\|_F^2 \quad \text{s.t.} \quad \|\Sigma\|_0 \leq d.}$$

where $A = U \Sigma V^T$

Assume hereafter that X is "centered"
column mean = 0.

How to solve?

Note:

$$\min_A \left\| X - A \right\|_F^2 = \min_{U, \Sigma_A, V_A} \left\| \Sigma_X - U \Sigma_A V_A^T \right\|_F^2$$

$$\begin{cases} U = U_X^T U_A \\ V = V_X^T V_A \end{cases}$$

$$= \min_{U, \Sigma_A, V_A} \left\| \Sigma_X \right\|_F^2 - 2 \langle \Sigma_X, U \Sigma_A V^T \rangle + \left\| \Sigma_A \right\|_F^2$$

Assume Σ_A : fixed.

$$\Rightarrow \max_{U, V} \langle \Sigma_X, U \Sigma_A V^T \rangle.$$

$$\text{tr}(F^T G) = \langle F, G \rangle$$

From Van Neumann's Inequality, $\langle F, G \rangle \leq \sum_{i=1}^n \sigma_i(F) \sigma_i(G)$

and equality holds iff $\exists V_F, U_G: F = U \Sigma_F V^T, G = U \Sigma_G V^T$.

$$\Rightarrow \max_{U, V} \langle \Sigma_X, U \Sigma_A V^T \rangle \leq \sum_{i=1}^d \sigma(\Sigma_X) \sigma(A) = \sigma(A)$$

To maximize this, sig. vectors should be the same

$$U = U_X^T U_A = I \quad \wedge \quad V = V_X^T V_A = I. \quad (2)$$

$$\rightarrow U_A = U_X, \quad V_A = V_X.$$

\rightarrow from here also:

$$\begin{aligned} \text{Thus, } \min_A \|X - A\|_F^2 &\Rightarrow \min_{\Sigma_A} \|\Sigma_A\|_F^2 - 2\langle \Sigma_X, \Sigma_A \rangle \\ &\Rightarrow \min_{\Sigma_A} \|\Sigma_X - \Sigma_A\|_F^2 \quad \text{s.t. } \|\Sigma_A\|_0 \leq d. \\ &\Rightarrow \min_{\sigma_i(A)} \sum_{i=1}^d \sigma_i(A)^2 - \sigma_i(X) \sigma_i(A) \\ &\Rightarrow \sigma_i(A) = \sigma_i(X) \text{ for } 1 \leq i \leq d \end{aligned}$$

$$\text{Thus, } A_d = \arg \min_A \|X - A\|_F^2 \text{ s.t. } \|\Sigma_X\|_0 \leq d$$

$$A_d = U_d \Sigma_d V_d^T \quad U_d, V_d : \text{top } d \text{ left and right s.v.}$$

$$\text{Note that we can write } A = U_X H_d(\Sigma_X) V_{dX}^T.$$

(Also, look @ Mantas Mazeika, Sing. V. Decomp. and low rank approx.).

Matrix Completion

Say we now observe only a subset of entries in X , Ω :

$$\begin{aligned} P_\Omega(X) &= \min_A \|P_\Omega(X) - P_\Omega(A)\|_F^2 \Rightarrow \text{infinite solutions!} \\ (X)_{ij} &\begin{cases} X_{ij} & \text{if } (i,j) \in \Omega \\ 0 & \text{otherwise} \end{cases} \\ \Rightarrow \min_A \|P_\Omega(X) - P_\Omega(A)\|_F^2 \text{ s.t. } \text{rk}(A) \leq d. \end{aligned}$$

This is now NP-hard, however.

Some Heuristic: iterative projections:

$$\text{init: } A^{(0)} = X$$

$$\text{iter: } - Z \leftarrow \arg \min_Z \|A^{(k)} - Z\|_F^2 \text{ s.t. } \text{rk}(Z) \leq d. \quad (\text{low-rank app})$$

$$- A^{(k+1)} \leftarrow P_\Omega(X) + Z_{\Omega^\perp} \quad (\text{better estimate closer to measurements})$$

depend greatly on initialization and not very robust.

Just Relax (again).

Instead of doing $\min_A \|P_2(X) - P_2(A)\|_F^2$ s.t. $\|\Sigma_A\|_0 \leq d$

$$\text{do } \min_A \|P_2(X) - P_2(A)\|_F^2 + \lambda \underbrace{\|\Sigma_A\|_1}$$

$$\sum_i \sigma_i(A) = \|A\|_*.$$

"Nuclear Norm" or
"trace norm".

Let's start by looking at the simpler
problem:

$$\min_X \|Y - X\|_F^2 + \lambda \|X\|_* \quad ; (P_\lambda^*)$$

\rightarrow convex - why?

Recall that this is equivalent to

$$\min_X \|\Sigma_Y - U \Sigma_X V^T\|_F^2 + \lambda \|X\|_*$$

$$\min_X \|\Sigma_Y - \Sigma_X\|_F^2 + \lambda \sum_i \sigma_i(X) \quad \text{from Von Neumann's.}$$

define $S_i = \text{diag}(\Sigma_i)$, then

$$\min \|S_Y - S_X\|_2^2 + \lambda \|S_X\|_1$$

$$\Rightarrow \hat{S}_X = S_\lambda(S_Y) \quad \text{- soft thresholding the singular values.}$$

$$\text{Thus, } \hat{X} = U S_\lambda(\Sigma_Y) V^T = \arg \min_X \|Y - X\|_F^2 + \lambda \|X\|_*.$$

$$\hat{X} = D_\lambda(Y) = \underset{\lambda \| \cdot \|_*}{\text{prox}}(Y).$$

(3)

an alternative (interesting) proof) is in Cai 2005.

Given $h(y) = \min_x \|y - x\|_F^2 + \lambda \|x\|_*$, recall that

$$z \in \partial f(x_0) \text{ if } f(x) \geq f(x_0) + \langle z, x - x_0 \rangle \quad \forall x.$$

and $\hat{x} = S_\lambda(y)$ minimizes $h(x)$ iff $0 \in \partial h(\hat{x})$:

$$0 \in \hat{x} - y + \lambda \partial \|\hat{x}\|_* \quad (1)$$

Let $X = U \Sigma V^T$. It is known that ($X \in \mathbb{R}^{n \times m}$)

$$\partial \|X\|_* = \left\{ UV^T + W : W \in \mathbb{R}^{n \times m}, U^T W = 0, W V = 0, \|W\|_2 \leq 1 \right\}.$$

Exercise: Verify that indeed (1) holds.

for $\hat{x} = U S_\lambda(\Sigma) V^T$.

So, going back to $(P_{\lambda, \Omega}^*)$: $\min_{\tilde{x}} \|P_\Omega(x) - P_\Omega(\tilde{x})\|_2 + \lambda \|\tilde{x}\|_*$

This is convex and can be solved with standard (Semi-Def. Program) but its expensive. Better: proximal gradient Descent:

Given y , init $\hat{x} = 0$, then:

$$\tilde{x}^{k+1} = \arg \min_{\tilde{x}} \|P_\Omega(x) - P_\Omega(\tilde{x})\|_2 + \lambda \|\tilde{x}\|_*$$

$$\Rightarrow X^{k+1} = \underset{\lambda \| \cdot \|_*}{\text{prox}} (X^k - \nabla f(X^k)).$$

$$X^{k+1} = S_{\lambda}^s (X^k - (P_{\Omega}(X^k) - P_{\Omega}(Y)))$$

Noting that $P_{\Omega}^{\perp}(X^k) = X^k - P_{\Omega}(X^k)$

$$\Rightarrow \boxed{X^{k+1} = S_{\lambda}^s (P_{\Omega}(Y) + P_{\Omega}^{\perp}(X^k))}$$

This converges to global optimum with rate $O(1/k)$.

Even better, note that $P_{\Omega}(Y) + P_{\Omega}^{\perp}(X^k) = \underbrace{P_{\Omega}(Y) - P_{\Omega}(X^k)}_{\text{sparse}} + \underbrace{X^k}_{\text{low rank}}$
 so this is efficient.

Guarantees for matrix completion

- No noise.
- How many samples from a low rank matrix do we need to recover it exactly?
 - first observation: if ~~no~~ no observations of a row ^{or} ~~and~~ columns, then it's impossible to recover the matrix exactly, even if it is of rank 1.
 - So how many samples do we need to observe at least 1 entry from every row & column? (with high probability)

actually: $N > p \log p$ ($p \times p$ matrix).

$$P_r[N > p \log p] \leq p^{-p}$$

\Rightarrow "coupon collector problem".

second observation:

(4)

We need a "good" matrix, \Rightarrow incoherent with the canonical basis.

Consider $Z = e_i e_i^T : = \begin{pmatrix} 1 & 0 & 0 & \dots \\ 0 & 0 & & \\ 0 & & \ddots & \\ \vdots & & & \ddots \end{pmatrix} \Rightarrow$ if we observe $N \ll p^2$ we won't see anything.

Intuition: we need the singular vectors to be "spread out".

Def: Matrix incoherence:

$X \in \mathbb{R}^{p \times m}$ is \mathcal{D} -incoherent w.r.t. sparse matrices if $\text{GOF rank} = r$.

$$\max_i \|u_i^T\|_2 \leq \mathcal{D} \sqrt{\frac{r}{p}}, \quad \max_j \|v_j^T\|_2 \leq \mathcal{D} \sqrt{\frac{r}{m}},$$

$$\|UV^T\|_\infty \leq \frac{\mathcal{D} \sqrt{r}}{\sqrt{pm}}.$$

u_i^T, v_j^T are the rows of U, V .

\Rightarrow This measures "spikiness"; it's low if entries in U, V are not concentrated.

Note that $\|U\|_2 = 1$ (columns). If these are concentrated ($U = I$),

then $\max_i \|u_i^T\|_2 = 1 \Rightarrow \mathcal{D} \geq \sqrt{p/r} > 1. \quad (\gg 1)$

If u_i is "spread out" maximally $\Rightarrow u_i = \pm \sqrt{1/p} \Rightarrow \mathcal{D} \approx 1$.
 $\mathcal{D} = 1$

In fact, this property is satisfied for Gaussian matrices as long as p is large enough.

Low Rank Matrix Recovery through Candes Options.

Let $X \in \mathbb{R}^{p \times p}$, ν -incoherent, and M is the expected number of observed entries sampled at random (eg. $M = p^2 \cdot t$)

Then, \exists c : \downarrow constant such that if

prob. that a sample is observed.

$$M \geq c \nu^4 p (\log p)^2 = \Omega(dp \text{ polylog}(p)).$$

[Candes & Recht, '09
Candes & Tao '10]

then $X = \arg \min_A \|A\|_*$ s.t. $P_\Omega(X) = P_\Omega(A)$. with $\Pr \geq 1 - p^{-3}$

Note that this is pretty tight:

How many degrees of freedom in a low rank matrix?

$$\text{Dof} = 2rp - r^2$$

Say we choose the first r columns to be ~~the~~ ~~the~~ of dim p (so $p \cdot r$). Say these are l.i.

Then, for the remainder columns, they will be linear combinations of the previous r columns. (r coeff)

$$\text{Thus Dof: } rp + (p-r) \cdot r = 2rp - r^2.$$

So up to a poly log factor, this is the lower bound.

Robust PCA:

So convex relaxation of the rank can recover a low rank matrix - if there's no noise. If some entries in X are corrupted, then the ℓ_2 will be very sensitive.

Upgrading this model:

$$X = L_0 + E_0, \quad \text{where } \text{rank}(L_0) \leq r.$$

E_0 : "outliers"

Thus:

and sparse $\|E\|_0 \leq k$.

$$\min_{L, E} \text{rank}(L) + \lambda \|E\|_0 \quad \text{s.t. } X = L + E.$$

This looks impossible! :

- p^2 equations (linear) for $2p^2$ unknowns.
- both losses are ^{non-}convex
- Solution might not be unique:

$$\text{imagine } Z = e, e^T \Rightarrow (L, E) = (Z, 0) \text{ or}$$

$$(L, E) = (0, Z) ?$$

The wellposed-ness of this problem will naturally depend on similar "incoherence" properties for the L_0 , and the outliers cannot be "conspicuously located", say block and entire column/row.

\Rightarrow Double Relax!

Principal Component Pursuit
(PCP)

$$\min_{L, E} \|L\|_* + \lambda \|E\|_1 \quad \text{s.t. } X = L + E$$

Guarantees: Say $X: p \times p$; $L_0 = U \Sigma V^T$, $U, V: p \times r$,
and $\text{supp}(E_0)$ is uniformly distributed, if

$$\text{rank}(L_0) \leq \frac{a \cdot p}{\sqrt{\log^2 p}} \quad \text{and} \quad \|E\|_0 \leq b \cdot p^2$$

for some constants $a, b > 0$, $\Rightarrow \exists c > 0$: $\hat{L} = L_0$, $\hat{E} = E_0$ with

$$Pr \geq 1 - c p^{-10} \quad (\text{as long as } \lambda = \frac{1}{\sqrt{p}}).$$

[Candes '11].