

$$\textcircled{1} \min_{x,y} f(x,y) = \min_x (\min_y f(x,y)) = \min_y (\min_x f(x,y))$$

Suppose  $x^*, y^*$  is optimal  $\textcircled{2}$  attained by  $\textcircled{2}$ , Suppose  $x^{**}, y^{**}$  is attained by  $\textcircled{1}$ , then we have

$$f(x^{**}, y^{**}) \leq f(x^*, y^*) = \min_x (\min_y f(x,y)) \leq \min_x (f(x, y^{**})) \leq f(x^{**}, y^{**})$$

$\textcircled{2}$  Latent Feature Lasso Question:

$$\min_{k \in N, Z \in \{0,1\}^{N \times k}, W \in \mathbb{R}^{k \times D}} \frac{1}{2N} \|X - ZW\|_F^2 + \frac{\lambda}{2} \|W\|_F^2$$

Consider the situation fitting  $k$  and  $Z$

$\Rightarrow$  New problem:

$$\min_{W \in \mathbb{R}^{k \times D}} \left\{ \frac{1}{2N} \|X - ZW\|_F^2 + \frac{\lambda}{2} \|W\|_F^2 \right\}$$

$\hookrightarrow$  Say it as  $L(W)$

Introducing slack variable  $E = ZW$

$$\Rightarrow \min_{W \in \mathbb{R}^{k \times D}, E \in \mathbb{R}^{N \times D}} \frac{1}{2N} \|X - E\|_F^2 + \frac{\lambda}{2} \|W\|_F^2 \text{ subject to } E = ZW \quad \textcircled{A}$$

before using Lagrange duality  $\Rightarrow$  we check strong duality first

$$\text{consider } \Omega = \begin{bmatrix} W \\ E \end{bmatrix} \begin{matrix} k \times D \\ N \times D \end{matrix}, \quad P_W \Omega = W \begin{matrix} k \times D \\ \end{matrix}$$

$$P_E \Omega = E \begin{matrix} N \times D \\ \end{matrix}$$

the problem becomes

$$\min_{\Omega} \frac{1}{2N} \|X - P_E \Omega\|_F^2 + \frac{\lambda}{2} \|P_W \Omega\|_F^2 \text{ subject to}$$

$P_E \Omega = Z P_W \Omega$ , which is convex problem of  $\Omega$  and is feasible, since simply choosing  $W$  and let  $E = ZW$  (Slater's condition)

Consider  $\textcircled{A}$ 's lagrange duality =

$$L(W, E, A) = \frac{1}{2N} \|X - E\|_F^2 + \frac{\lambda}{2} \|W\|_F^2 + \langle A, E - ZW \rangle$$

$$\textcircled{A} = \min_{W, E} \max_A L(W, E, A) = \max_A \min_{W, E} L(W, E, A) \quad (\text{by strong duality})$$

$$\max_A \min_E \underbrace{\left( \min_W \left( \frac{1}{2N} \|X - E\|_F^2 + \frac{\lambda}{2} \text{Tr}(W^T W) + \langle A, E - ZW \rangle \right) \right)}_{g(W)}$$

simply  
 $\nabla g(W) = 0 = \lambda W + \nabla \text{Tr}(-A^T Z W) = \lambda W + \nabla \text{Tr}(-W A^T Z) = \lambda W + Z^T A$

$W^*$  is given by  $-\frac{1}{\lambda} Z^T A$ , by substituting  $W^*$ , the question becomes

$$\max_A \min_E \left( \frac{1}{2N} \|X - E\|_F^2 + \frac{1}{2\lambda} \text{Tr}(A^T Z Z^T A) + \langle A, E \rangle + \langle A, -Z W^* \rangle \right)$$

$$= \max_A \min_E \left( L(E) + \frac{1}{2\lambda} \text{Tr}(A^T Z Z^T A) + \langle A, E \rangle - \frac{1}{\lambda} \text{Tr}(A^T Z Z^T A) \right)$$

$$= \max_A \min_E \left( L(E) + \frac{1}{2\lambda} \text{Tr}(A^T Z Z^T A) + \langle A, E \rangle \right) \quad \text{Recall convex conjugate } f^*(y) = \max_x y^T x - f(x)$$

$$= \max_A \left( \max_E (L(E) - \langle A, E Z \rangle) + \frac{1}{2\lambda} \text{Tr}(A^T Z Z^T A) \right)$$

$$= \max_A \left( -\frac{1}{2\lambda} \text{Tr}(A^T Z Z^T A) - L^*(A) \right) = \tilde{g}(M, A), \quad M = Z Z^T$$

when  $M$  is fixed,  $\tilde{g}(M, A)$  is a convex function of  $A$ ,

$g(M) \stackrel{\text{def}}{=} \max_A \tilde{g}(M, A)$  is still a convex function.  
 (pointwise maximum of convex)

The whole problem reduced to  
 $\min_Z g(M)$  while  $M = Z Z^T$

We introducing a trick called atomic (nuclear) norm regularization to force  $M$  has a structure of  $ZZ^T$ , Our problem is how

$\min_M \{g(M) + \mu \|M\|_S\}$  where  $\|M\|_S = \min_{c \geq 0, a \in S} \sum_{a \in S}^k c a$  s.t  $M = \sum_{a \in S} c a a$   
 $z \in \{0,1\}^N$  matrix  
 $S$  is a collection of all possible atoms,  $S = \{zz^T | z \in \{0,1\}^N\}$   
 (It can be shown that all possible combination of  $M$  covers  $ZZ^T$ , so it is actually a relaxation here).

Now we use greedy coordinate descent to solve this problem

$$\min_M \{g(M) + \lambda \|M\|_S\} = \min_{C \in \mathbb{R}_+^{\bar{K}}} \underbrace{\left\{ g\left(\sum_{j=1}^{\bar{K}} c_j z_j z_j^T\right) + \lambda \|C\|_1 \right\}}_{f(c)}$$

$\bar{K} = 2^{N-1}$

we first choose a coordinate of  $C$  to do coordinate descent, finding correspond atom  $z_j z_j^T$ , and then do proximal gradient update.

$$j^* = \operatorname{argmax}_j -\nabla_j f(c) = \operatorname{argmax}_j \langle \nabla g(M), z_j z_j^T \rangle$$

we can find this  $j$  using MAX-cut like problem.

Now consider the algorithm step, we have active set  $A$ , and new atom  $z_j z_j^T$  by maxcut, we want to update corresponding  $|A|$  number of  $C_i$  of our all atoms.

$$C^{r+1} \leftarrow \left[ C^r - \frac{\nabla f(C^r) + \lambda}{r |A|} \right]_+ \quad r=1,2,\dots,T_2 \quad (C > 0 \text{ differentiable})$$

$r$  is Lipschitz continuous of  $\nabla g_j f(c)$ .

By Danskin's theorem

$$\nabla c_j f(c) = z_j A^* A^{*T} z_j / (2N^2 T)$$

---

Danskin's theorem statement =

$f(x) = \max_{z \in Z} \phi(x, z)$  is a continuous function if  $\phi: \mathbb{R}^n \times Z \rightarrow \mathbb{R}$  is a continuous function,  $Z \subset \mathbb{R}^m$  is a compact set, further assume  $\phi(x, z)$  is convex in  $x$  for every  $z \in Z$  and  $Z$  is

$$f(x) = \max_{z \in Z} \phi(x, z)$$

$$Z_0(x) = \{ \bar{z} = \phi(x, \bar{z}) = \max_{z \in Z} \phi(x, z) \}$$

$\Rightarrow f(x)$  is convex

$$\nabla_y f(x) = \max_{z \in Z_0(x)} \phi'(x, z; y) \quad [\text{derivate of direction of } y \text{ of } x]$$

If  $\phi(x, z)$  is differentiable with respect to  $x$  for all  $z \in Z$ ,  $\frac{\partial \phi}{\partial x}$  is

$f(x)$  is differentiable at  $x$  if  $Z_0(x)$  consists of a single element

$\bar{z}$ , in this case, the derivative of  $f(x)$

$$\frac{\partial f(x)}{\partial x} = \frac{\partial \phi(x, \bar{z})}{\partial x}$$

---

$$g = \mathbb{R}^{n \times k} \times \mathcal{U} \rightarrow \mathbb{R}$$

$$\downarrow \quad \downarrow$$

$$z \quad A \in \mathbb{R}^{n \times d}$$

continuous ✓

convex if fix  $z$  ✓

also convex in  $A$ .

(compact if closed and bounded)  
I believe it because numerical  
error ?

$$g(z) = \max_{A \in \mathcal{U}} \tilde{g}(z, A)$$

$$\frac{\partial \text{Tr}(xx^T B)}{\partial x} = Bx + B^T x$$

convex: only one maximum  $\tilde{g}(z, A)$   $\left[ \begin{array}{l} \text{If is also} \\ \text{true when } A \text{ restricted on } \mathbb{R}^{+}_{z,z^T} \end{array} \right]$

$$\nabla_{C_j} g(C) = \nabla_{C_j} \left[ \max_A \left( -\frac{1}{2\lambda} \text{Tr} \left( A^T \sum_{j=1}^K C_j z_j z_j^T A \right) + L^*(A) \right) \right]$$

$$= \nabla_{C_j} \left[ -\frac{1}{2\lambda} \text{Tr} (A^T C_j z_j z_j^T A) + L^*(A) \right]$$

$$= -\frac{1}{\lambda} \text{Tr} (A^T z_j z_j^T A) = -\frac{1}{\lambda} (z_j^T A^T A z_j)$$

$$\begin{bmatrix} a_{11} & a_{12} \\ a_{21} & a_{22} \end{bmatrix}$$

$$\frac{\partial \text{Tr}(CA)}{\partial C} = \frac{\partial \text{Tr}(A)}{\partial C}$$

Maxcut Problem:

$$y = 2z - 1$$

$$z = \frac{y+1}{2}$$

$$\max_{z \in \{0,1\}^N} \langle C, zz^T \rangle = \max_{z \in \{0,1\}^N} \langle C, \underbrace{\frac{y+1}{2} \left( \frac{y+1}{2} \right)^T}_{\frac{yy^T + 2y^T + 11^T}{4}} \rangle$$

$$= \max_{z \in \{0,1\}^N} \frac{1}{4} (\langle C, yy^T \rangle + 2\langle C, 1y^T \rangle + \langle C, 11^T \rangle)$$

$$= \max_{\substack{(y_0, y) \\ y \in \{-1,1\}^{N+1}}} \frac{1}{4} \begin{bmatrix} y_0^T \\ y \end{bmatrix} \begin{bmatrix} 1^T C 1 & 1^T C \\ C 1 & C \end{bmatrix} \begin{bmatrix} y_0 \\ y \end{bmatrix}$$

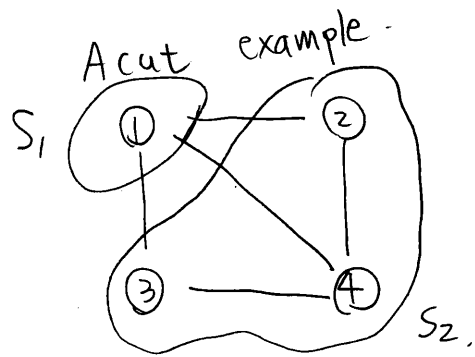
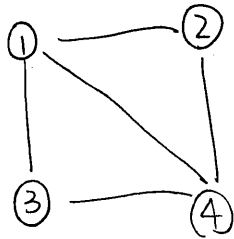
$$\begin{bmatrix} 1^T C 1 & 1^T C \\ C 1 & C \end{bmatrix} \begin{bmatrix} y_0 \\ y \end{bmatrix} = \begin{bmatrix} 1^T C 1 y_0 + 1^T C y \\ C 1 y_0 + C y \end{bmatrix}$$

$$\begin{bmatrix} y_0^T & y^T \end{bmatrix} \begin{bmatrix} 1^T C 1 y_0 + 1^T C y \\ C 1 y_0 + C y \end{bmatrix} = \begin{bmatrix} y_0^T 1^T C 1 y_0 + y_0^T 1^T C y \\ y_0^T C 1 y_0 + y^T C y \end{bmatrix}$$

$$= \langle C, yy^T \rangle + 2\langle C, 1y^T \rangle + \langle C, 11^T \rangle.$$

# Maxcut Problem:

$$A = \begin{bmatrix} 0 & 1 & 1 & 1 \\ 1 & 0 & 0 & 1 \\ 1 & 0 & 0 & 1 \\ 1 & 1 & 1 & 0 \end{bmatrix}$$



Suppose there are  $n$ -node,  $x_i$  characterized the cut by  $x_i = 1 \text{ if } i \in S_1$  and  $x_i = -1 \text{ if } i \in S_2$ .

$$\Rightarrow x = (1, -1, -1, -1)$$

each edge value of cut edge =  $\frac{1}{2} a_{ij} (1 - x_i x_j)$ .

$$W = \frac{1}{4} \sum_{i=1}^n \sum_{j=1}^n a_{ij} (1 - x_i x_j)$$

$$= \frac{1}{4} \left( \sum_{i=1}^n \sum_{j=1}^n x_i \overbrace{a_{ij}}^{=1} x_j - \sum_{i=1}^n \sum_{j=1}^n x_i a_{ij} x_j \right)$$

$$\Rightarrow W = \frac{1}{4} (x^T \text{Diag}(A \mathbf{1}) x - x^T A x)$$

transform (L)

$$L = \text{Diag}(A \mathbf{1}) - A$$

So the maximum cut problem can be formulated as

$$\max_{x \in \{-1, 1\}^n} \frac{1}{4} x^T L x$$

(~~similar~~ similar to spin-glass model VLSI.)

$\Rightarrow$  NP-hard problem,  $2^n$  possibilities

$\Rightarrow$  SDP-relaxation

$\max_{x \in \{-1,1\}^n} \frac{1}{4} x^T L x$  is a integer programming

$$\text{Trace}(x^T L x) = \text{Trace}(L x x^T) = L \bullet X, \quad X \text{ is real symmetric}$$

$$X \geq 0, \quad \text{diag}(X) \Leftrightarrow x_i^2 = 1$$

$$\begin{aligned} & (v^T x x^T v) \\ &= (x^T v)^T (x^T v) \geq 0 \end{aligned}$$



reformulation =

$$\text{Maximize } L \bullet X$$

$$\text{subject to } \text{diag}(X) = e, \quad \text{rank}(X) = 1 \text{ and } X \geq 0$$

$\Downarrow$  SDP - relaxation

$$\text{Maximize } L \bullet X$$

$$\text{subject to } \text{diag}(X) = e \text{ and } X \geq 0$$

$\hookrightarrow$  many possible solver!

$$\left( \begin{array}{l} \text{standard form:} \\ \min = L \bullet X \\ \text{subject to } A_i \bullet X = b_i \\ X \geq 0 \end{array} \right) \quad \text{pick } A_i = \begin{cases} a_{ii} = 1 \\ 0 \end{cases}$$