

LEC6 DSA5103

Coordinate Descent Algorithm

Recap: { ① Proximal Operator for $\| \cdot \|_2$ }
 { ② Block Coordinate Descent Algo for SVM (SMO Algo) }

① Proximal Operator

$$P_f(x) := \underset{y \in \mathbb{R}^n}{\operatorname{argmin}} f(y) + \frac{1}{2} \|y - x\|_2^2$$

a) $f = \delta_C = \begin{cases} 0, & x \in C \\ +\infty, & x \notin C \end{cases}$

property:
 $u = P_f(x) \Leftrightarrow x - u \in \partial f(u)$

then $P_{\delta_C}(x) = \underset{y \in \mathbb{R}^n}{\operatorname{argmin}} \delta_C(y) + \frac{1}{2} \|y - x\|_2^2$

$$= \underset{y \in C}{\operatorname{argmin}} \frac{1}{2} \|y - x\|_2^2$$

$$= \Pi_C(x)$$

then $\Psi_{\delta_C}(x) = \frac{1}{2} \|\Pi_C(x) - x\|_2^2$

b) $f = \| \cdot \|_1$

then $P_{\|\cdot\|_1}(x) = \underset{y \in \mathbb{R}^n}{\operatorname{argmin}} \|y\|_1 + \frac{1}{2} \|y - x\|_2^2$

How to achieve this? { ① Through $u = P_f(x) \Leftrightarrow x - u \in \partial f(u)$
 ② Moreau-Yosida Decomposition }

approach 2: $x = P_f(x) + P_{f^*}(x)$ f is positive homogeneous
 $\Rightarrow P_{f^*}(x) = \delta_{\partial f(0)}(x)$

$$\Rightarrow P_{f^*}(x) = \Pi_{\partial f(0)}(x)$$

$$x \in \partial f(0) \Leftrightarrow f(y) \geq x^\top (y - 0) + f(0) \quad \forall y$$

$$\Leftrightarrow f(y) \geq \langle x, y \rangle \quad \forall y$$

$$\Leftrightarrow \|y\|_1 \geq \langle x, y \rangle \quad \forall y \text{ s.t. } \|y\|_1 = 1$$

$$\Leftrightarrow \langle x, y \rangle \leq 1 \quad \forall y \text{ s.t } \|y\|_2 = 1$$

$$\Leftrightarrow \sup_y \{ \langle x, y \rangle : \|y\|_2 = 1 \} \leq 1$$

$$\Leftrightarrow \|x\|_\infty \leq 1.$$

$$\Rightarrow P_f(x) = \Pi_{B_\infty^1}(x)$$

$$\Rightarrow P_f(x) = x - \Pi_{B_\infty^1}(x) = \begin{cases} x-1 & x \geq 1 \\ 0 & -1 \leq x \leq 1 \\ x+1 & x \leq -1 \end{cases}$$

Generalization . $P_{\lambda \Pi B_\infty^1}(x) = x - \Pi_{B_\infty^\lambda}(x) = \begin{cases} x-\lambda & x \geq \lambda \\ 0 & -\lambda \leq x \leq \lambda \\ x+\lambda & x \leq -\lambda \end{cases}$

② Block Coordinate Descent

Dual Formulation of SVM

$$\begin{array}{ll} \min_{\alpha} & \frac{1}{2} \sum_i \sum_j \alpha_i \alpha_j y_i y_j \langle x_i, x_j \rangle - \sum_i \alpha_i \\ \text{s.t.} & \sum_{i=1}^n \alpha_i y_i = 0 \\ & 0 \leq \alpha_i \leq C \quad i \in [n] \end{array}$$

① choose (i, j) pair

$$\begin{array}{ll} \text{② update } (\hat{\alpha}_i, \hat{\alpha}_j) = & \left\{ \begin{array}{l} \underset{\alpha_i, \alpha_j}{\operatorname{argmin}} \quad \frac{1}{2} K_{ii} \alpha_i^2 + \frac{1}{2} K_{jj} \alpha_j^2 + K_{ij} y_i y_j \alpha_i \alpha_j \\ - \alpha_i - \alpha_j \\ \text{s.t.} \quad y_i \alpha_i + y_j \alpha_j = C (\underline{\alpha_i - \alpha_j}) \\ 0 \leq \alpha_i \leq C \\ 0 \leq \alpha_j \leq C \end{array} \right. \end{array}$$

closed-form solution

$$\textcircled{3} \quad (\alpha_i, \alpha_j) \leftarrow (\hat{\alpha}_i, \hat{\alpha}_j)$$

→ It seems that, we are interested in solutions like

$$\tilde{x}_i = \underset{x_i}{\operatorname{argmin}} f(x_{-i}, x_i) \quad (\text{fix } x_{-i}) \quad \forall i$$

(coordinate-wise minimizer)

Today's content

1. Coordinate-wise Minimizer

→ { defn
in which condition, it can actually be global minimizer?

① Defn:

\bar{x} is coordinate-wise minimizer

$$\Leftrightarrow f(\bar{x} + d\epsilon_i) \geq f(\bar{x}) \quad \begin{cases} \forall d \in \mathbb{R} \\ \forall i \in [n] \end{cases}$$

Alternative defn:

$$\bar{x}_i \in \underset{x_i}{\operatorname{argmin}} f(x_i, \bar{x}_{-i}) \quad \forall i \in [n]$$

② Question:

Global Minimizer

↑ relation?

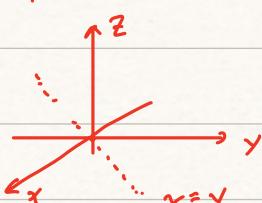
Coordinate-wise Minimizer

E.g. $f(x_1, x_2) = x_1^2 + x_2^2 + 20|x_1 - x_2|$

a) Coordinate-wise minimizer

$$g(x, y) = |x - y|$$

↓ convex!



$$\bar{x}_1 \in \underset{x_1}{\operatorname{argmin}} f(x_1, \bar{x}_2)$$

$$= \underset{x_1}{\operatorname{argmin}} x_1^2 + 20|x_1 - \bar{x}_2|$$

convex

$$\Leftrightarrow 0 \in 2\bar{x}_1 + 20 \partial g_{\bar{x}_2}(\bar{x}_1)$$

property:

$$f, g \text{ convex} \Rightarrow f+g \text{ convex}$$

Pf: $\forall x, y, \lambda$

$$\begin{aligned} & (f+g)(\lambda x + (1-\lambda)y) \\ &= f(\lambda x + (1-\lambda)y) + g(\lambda x + (1-\lambda)y) \\ &\leq \lambda(f+g)(x) + (1-\lambda)(f+g)(y) \end{aligned}$$

Here $g_{\bar{x}_2}(x) := |x - \bar{x}_2| \Rightarrow \partial g_{\bar{x}_2}(\bar{x}_1) = \begin{cases} [-1, 1] & \bar{x}_1 = \bar{x}_2 \\ -1 & \bar{x}_1 < \bar{x}_2 \\ 1 & \bar{x}_1 > \bar{x}_2 \end{cases}$

$$\Rightarrow 2\bar{x}_1 + 2_0 \partial g_{\bar{x}_2}(\bar{x}_1) = \begin{cases} [-2_0 + 2\bar{x}_2, 2_0 + 2\bar{x}_2] & \bar{x}_1 = \bar{x}_2 \\ -2_0 + 2\bar{x}_1 & \bar{x}_1 < \bar{x}_2 \\ 2_0 + 2\bar{x}_1 & \bar{x}_1 > \bar{x}_2 \end{cases}$$

\rightarrow Suppose $\bar{x}_1 > \bar{x}_2$ then $\bar{x}_1 = -1_0$ so that $0 \in 2\bar{x}_1 + 2_0 \partial g_{\bar{x}_2}(\bar{x}_1)$

$$\begin{aligned} \text{However, when } \bar{x}_2 < \bar{x}_1, \quad & 2\bar{x}_2 + 2_0 \partial g_{\bar{x}_2}(\bar{x}_2) \\ &= 2_0 + 2\bar{x}_2 \\ &< 2_0 + 2\bar{x}_1 \\ &< 0 \end{aligned}$$

which means (\bar{x}_1, \bar{x}_2) cannot contribute to "Coordinate-wise minimizer"

\rightarrow Suppose $\bar{x}_1 = \bar{x}_2$, then we require:

$$\underline{-1_0 \leq \bar{x}_1 = \bar{x}_2 \leq 1_0}$$

$\Rightarrow \{(\bar{x}_1, \bar{x}_2) : -1_0 \leq \bar{x}_1 = \bar{x}_2 \leq 1_0\}$ is "Coordinate-wise Minimizer"

b) Global Minimizer

obviously, $(0, 0)$ is global minimizer

(from above example) \rightarrow Personal understanding

Analysis: why Coordinate-wise minimizer cannot be Global minimizer?

$$f(x_1, x_2) = x_1^2 + x_2^2 + 2_0 |x_1 - x_2|$$

$$g(x) = h(x_1 - x_2)$$

$$\begin{aligned} \text{consider } \partial f(x) &= \nabla_x (x_1^2 + x_2^2) + 2_0 \partial g(x) \leftarrow h(z) = |z| \\ &= \begin{pmatrix} 2x_1 \\ 2x_2 \end{pmatrix} + 2_0 \begin{pmatrix} 1 \\ -1 \end{pmatrix} \cdot \partial h(x_1 - x_2) \end{aligned}$$

$$= \begin{pmatrix} 2x_1 \\ 2x_2 \end{pmatrix} + \begin{pmatrix} 2_0 \\ -2_0 \end{pmatrix} \cdot \begin{cases} [-1, 1] & x_1 = x_2 \\ -1, & x_1 < x_2 \\ 1, & x_1 > x_2 \end{cases}$$

if $-10 \leq x_1 = x_2 \leq 10$, \rightarrow coordinate-wise minimizer

then $\partial f(x) = \{(2x_1 + 20\lambda, 2x_2 - 20\lambda) : \lambda \in [-1, 1]\}$

$$= \{(2a + 20\lambda, 2a - 20\lambda) : \lambda \in [-1, 1], a = x_1\}$$

Let's check whether $\underline{0} = \begin{pmatrix} 0 \\ 0 \end{pmatrix} \in \partial f(x)$

if so, then there exists $\underline{\lambda} \in [-1, 1]$, s.t

$$\begin{cases} 0 = 2a + 20\lambda \\ 0 = 2a - 20\lambda \end{cases} \Rightarrow \begin{cases} \lambda = 0 \\ a = 0 = \underline{x}_1 = \underline{x}_2 \end{cases}$$

★
only global minimizer

Moreover, for a general function $f(\cdot)$. the main issue is:

if \bar{x} is a coordinate-wise minimizer

\rightarrow then suppose some regularity condition \rightarrow differentiable

$$\text{we can achieve } \nabla_i f(\bar{x}) = 0 \quad \forall i \in [n]$$

\Updownarrow

(intuition) $\nabla f(\bar{x}) = 0 \Rightarrow \bar{x}$ global minimizer

\rightarrow However, if $f(x)$ is not smooth

then we cannot control the change of

function value in all direction even if function
value is well-shaped in $x(y)$ direction

→ Main Issue comes from the "discontinuity" of "gradient"

\rightarrow Back to the relation between "Coordinate-wise minimizer"
and "global minimizer"

we can make the Formal Description as follows

a) if f is $\begin{cases} \text{convex} \\ \text{differentiable} \end{cases}$, then $\frac{\text{Coordinate-wise Minimizer}}{\Downarrow}$

Global Minimizer

(Easy to prove)

b)

Generally speaking

Coordinate-wise Minimizer



Global Minimizer

even if $f(\cdot)$ is convex

(example is one case)

c)

Target Problem → (when Coordinate-wise minimizer



is indeed global minimizer)

$$\min_{\mathbf{x}} f(\mathbf{x}) + \sum_{i=1}^n f_i(x_i)$$

Here: $\begin{cases} f(\cdot) \rightarrow \text{convex \& differentiable} \\ f_i(\cdot) \rightarrow \text{only convex} \end{cases}$

We have: Coordinate Minimizer \Leftrightarrow Global Minimizer

Pf: $\bar{x}_i \in \underset{x_i}{\operatorname{argmin}} f(x_i, \bar{x}_{-i}) + f_i(x_i) \quad \forall i \in [n]$

convex

$$\Leftrightarrow \underline{0} \in \nabla_i f(\bar{x}) + \underline{\partial f_i(\bar{x}_i)} \quad \forall i \in [n]$$

Scalar

1-dim

Question: $F(x) = \sum_{i=1}^n f_i(x)$

$$\text{Then } \nabla F(x) = \sum_{i=1}^n \underline{\partial f_i(x)} \rightarrow d\text{-dim}$$

Here, $f_i(x) := f_i(x_i) \rightarrow \text{only related to } x_i$

$$\text{Then } \underline{\zeta} \in \partial f_i(x) \Leftrightarrow \underline{\forall y}, f_i(y) \geq \zeta^\top (y - x) + f_i(x)$$

$$\rightarrow \text{choose } y_j = \begin{pmatrix} \zeta_j \\ x_i \end{pmatrix}, \text{ we have } 0 \geq \varepsilon_j \cdot \zeta_j \quad \forall j$$

forcing $\zeta_j = 0$ for $\forall j \neq i$

\rightarrow if all $\beta_j = 0$ for $\forall j \neq i$

then $f_i(y) \geq \beta^T(y - x) + f_i(x)$ by

$$\Leftrightarrow f_i(y_i) \geq \beta_i(y_i - x_i) + f_i(x_i)$$

$$\Leftrightarrow \beta_i \in \partial f_i(x_i)$$

$$\Rightarrow \partial f_i(x) = \{0\} \times \dots \times \partial f_i(x_i) \times \dots \times \{0\}$$

$$\Rightarrow \partial F(x) = \sum_{i=1}^n \partial f_i(x) = \prod_{i=1}^n \partial f_i(x_i)$$

i-th dimension cartesian product

$$\Leftrightarrow \underline{0} \in \nabla f(\bar{x}) + \partial F(\bar{x}) \Leftrightarrow \bar{x} \text{ is global minimizer}$$

Here $\nabla f(\bar{x}) := \begin{bmatrix} \nabla_1 f(\bar{x}) \\ \vdots \\ \nabla_n f(\bar{x}) \end{bmatrix}$

$$\partial F(\bar{x}) := \begin{bmatrix} \partial f_1(\bar{x}_1) \\ \vdots \\ \partial f_n(\bar{x}_n) \end{bmatrix}$$

* this part is highly nontrivial!

* and this is why some other general function cannot achieve this good property!

2. How to achieve Coordinate-wise minimizer? (Algo)

(Brute Force)

a) $\bar{x} \rightarrow$ Coordinate-wise Minimizer

Actually like
Nash Equilibrium!

\Rightarrow (Necessary Condition)

$$\bar{x}_i \in \underset{x_i}{\operatorname{argmin}} f(x_i, \bar{x}_{-i})$$

Consider $A_i := \{(x_i, x_{-i}) : x_i \in \underset{x_i}{\operatorname{argmin}} f(\bar{x}_i ; x_{-i})\}$

$\Rightarrow \{\text{Coordinate Minimizer}\} = \bigcap_{i=1}^n A_i$

(intractable)

b) Coordinate Descent Algorithm

$$\left\{ \begin{array}{l} \rightarrow \text{given } x^{(k)} \\ \rightarrow x_i^{(k+1)} = \underset{x_i}{\operatorname{argmin}} f(x_{<i}^{(k+1)}, x_i, x_{>i}^{(k)}) \\ \rightarrow k \rightarrow k+1 \end{array} \right.$$

3. Application

a) $\boxed{\min_x (x_1 - x_2)^2 + |x_1| + |x_2|}$

\rightarrow Algo:

$$\Rightarrow \text{given } x_1^{(k)}, x_2^{(k)}$$

$$\Rightarrow x_1^{(k+1)} = \underset{x_1}{\operatorname{argmin}} (x_1 - x_2^{(k)})^2 + |x_1|$$

$$= P_{\frac{1}{2}f} (x_2^{(k)})$$

$$= P_{\frac{1}{2}\| \cdot \|_1} (x_2^{(k)})$$

$$= \begin{cases} x_2^{(k)} - \frac{1}{2}, & x_2^{(k)} > \frac{1}{2} \\ 0, & 0 \leq \dots \\ x_2^{(k)} + \frac{1}{2}, & x_2^{(k)} < -\frac{1}{2} \end{cases}$$

$$\Rightarrow x_2^{(k+1)} = P_{\frac{1}{2}\| \cdot \|_1} (x_1^{(k+1)})$$

$$= \begin{cases} x_1^{(k+1)} - \frac{1}{2}, & x_1^{(k+1)} > \frac{1}{2} \\ 0, & 0 \leq \dots \\ x_1^{(k+1)} + \frac{1}{2}, & x_1^{(k+1)} < -\frac{1}{2} \end{cases}$$

b) Linear Reg with Coordinate Descent

$$\min_{\beta} \frac{1}{2} \|X\beta - y\|_2^2$$

differentiable
+ convex

$x_{\cdot i} \in \mathbb{R}^p$ $\beta_i \in \mathbb{R}$

if want to apply CD, then we should solve:

$$\hat{\beta}_i \in \arg \min_{\beta_i} \frac{1}{2} \|x_{\cdot i} \beta_i + x_{\cdot -i} \hat{\beta}_{-i} - y\|_2^2$$

$\frac{1}{2} z^T z$ $z = \leftarrow$

$$\Leftrightarrow \nabla_{\beta_i} L(\beta) = 0$$

$$\Leftrightarrow \nabla_{\beta_i} z \nabla_z L(z) = 0$$

$$\Leftrightarrow x_{\cdot i}^T \cdot (x_{\cdot i} \hat{\beta}_i + x_{\cdot -i} \hat{\beta}_{-i} - y) = 0$$

$$\Leftrightarrow \hat{\beta}_i = \frac{1}{\|x_{\cdot i}\|_2^2} [x_{\cdot i}^T (y - x_{\cdot -i} \hat{\beta}_{-i})]$$

↔

"Adaptive weight GD"

$$\begin{aligned} \hat{\beta}_i &\leftarrow \hat{\beta}_i - \frac{1}{\|x_{\cdot i}\|_2^2} (x_{\cdot i}^T (y - x_{\cdot -i} \hat{\beta}_{-i} - x_{\cdot i} \hat{\beta}_i)) \\ &= \hat{\beta}_i - \frac{1}{\|x_{\cdot i}\|_2^2} (x_{\cdot i}^T (y - X \hat{\beta})) \end{aligned}$$

GD update: $\hat{\beta}_i \leftarrow \hat{\beta}_i - x_{\cdot i}^T (y - X \hat{\beta})$

$$\Leftrightarrow \hat{\beta} \leftarrow \hat{\beta} - X^T (y - X \hat{\beta})$$

STILL cannot guarantee the sparsity of solution!



if we want, then we should re-formulate the optimization problem

ℓ_1 -norm LASSO

c) LASSO

$$\min_{\beta} \frac{1}{2} \|X\beta - y\|_2^2 + \lambda \|\beta\|_1$$

convex & differentiable separable & convex

(LASSO)

For this task, if we want to apply CP Framework,
then we should solve:

$$\begin{aligned}
 \hat{\beta}_i &\in \arg \min_{\beta_i} \frac{1}{2} \|x_{\cdot i} \beta_i + x_{\cdot \cdot i} \hat{\beta}_{-i} - y\|_2^2 + \lambda |\beta_i| \\
 &= \arg \min_{\beta_i} \frac{1}{2} \|x_{\cdot i}\|_2^2 \beta_i^2 + (x_{\cdot i}^\top (x_{\cdot \cdot i} \hat{\beta}_{-i} - y)) \beta_i + \lambda |\beta_i| \\
 &= \arg \min_{\beta_i} \frac{1}{2} \left(\beta_i^2 + \frac{x_{\cdot i}^\top \Delta}{\|x_{\cdot i}\|_2^2} \beta_i \right) + \frac{\lambda}{\|x_{\cdot i}\|_2^2} |\beta_i| \\
 &= \arg \min_{\beta_i} \frac{1}{2} \left(\beta_i + \frac{x_{\cdot i}^\top \Delta}{\|x_{\cdot i}\|_2^2} \right)^2 + \frac{\lambda}{\|x_{\cdot i}\|_2^2} |\beta_i| \\
 &= \arg \min_{\beta_i} \frac{\lambda}{\|x_{\cdot i}\|_2^2} |\beta_i| + \frac{1}{2} \left(\beta_i - \left(-\frac{x_{\cdot i}^\top \Delta}{\|x_{\cdot i}\|_2^2} \right) \right)^2 \\
 &= P_{\frac{\lambda}{\|x_{\cdot i}\|_2^2} \mathbb{I}} \left(-\frac{x_{\cdot i}^\top \Delta}{\|x_{\cdot i}\|_2^2} \right) \\
 &= P_{\frac{\lambda}{\|x_{\cdot i}\|_2^2} \mathbb{I}} \left(\boxed{\frac{x_{\cdot i}^\top (y - x_{\cdot \cdot i}^\top \hat{\beta}_{-i})}{\|x_{\cdot i}\|_2^2}} \right)
 \end{aligned}$$

the update for LR without regularization

through Coordinate Descent Method

$$\text{LR : } \hat{\beta}_i \leftarrow \frac{x_{\cdot i}^\top (y - x_{\cdot \cdot i}^\top \hat{\beta}_{-i})}{\|x_{\cdot i}\|_2^2}$$

$$\text{LASSO : } \hat{\beta}_i = \text{soft-threshold}(\hat{\beta}_i)$$

d) Box-constraint Regression

$$\begin{cases} \min_{\beta} \frac{1}{2} \|X\beta - y\|_2^2 \\ \text{s.t. } l \leq \beta \leq u \end{cases}$$

unconstrained optimization

$$\Leftrightarrow \min_{\beta \in \mathbb{R}^p} \frac{1}{2} \|X\beta - y\|_2^2 + \delta_C(\beta)$$

$C = \{\beta : l \leq \beta \leq u\}$
 $= C_1 \times C_2 \times \dots \times C_p$

$$\Leftrightarrow \min_{\beta \in \mathbb{R}^p} \frac{1}{2} \|X\beta - y\|_2^2 + \sum_{i=1}^p \delta_{C_i}(\beta_i)$$

$C_i = \{\beta_i : l_i \leq \beta_i \leq u_i\}$

if want to Apply Coordinate Descent Algo ,

then we need to solve :

$$\begin{aligned} \hat{\beta}_i &= \arg \min_{\beta_i} \frac{1}{2} \|x_{\cdot i} \beta_i + x_{-i} \hat{\beta}_{-i} - y\|_2^2 + \delta_{C_i}(\beta_i) \\ &= \arg \min_{\beta_i} \frac{1}{2} \|x_{\cdot i}\|_2^2 \beta_i^2 + x_{\cdot i}^\top (\underbrace{x_{-i}^\top \hat{\beta}_{-i} - y}_{\textcircled{4}}) \beta_i + \delta_{C_i}(\beta_i) \\ &= \arg \min_{\beta_i} \frac{1}{2} \left[\beta_i - \left(\frac{x_{\cdot i}^\top (y - x_{-i}^\top \hat{\beta}_{-i})}{\|x_{\cdot i}\|_2^2} \right) \right]^2 + \frac{1}{\|x_{\cdot i}\|_2^2} \delta_{C_i}(\beta_i) \\ &= P \frac{1}{\|x_{\cdot i}\|_2^2} \delta_{C_i}(\cdot) \left(\frac{x_{\cdot i}^\top (y - x_{-i}^\top \hat{\beta}_{-i})}{\|x_{\cdot i}\|_2^2} \right) \\ &= P \delta_{C_i}(\cdot) \left(\frac{x_{\cdot i}^\top (y - x_{-i}^\top \hat{\beta}_{-i})}{\|x_{\cdot i}\|_2^2} \right) \\ &= \Pi_{C_i} \left(\frac{x_{\cdot i}^\top (y - x_{-i}^\top \hat{\beta}_{-i})}{\|x_{\cdot i}\|_2^2} \right) \end{aligned}$$

4. Block CD & Application

$$\left\{ \begin{array}{l} \min_{x \in X} F(x) = f(x) + \sum_{i=1}^m r_i(x_i) \\ \text{s.t. } x_i \in X_i \quad \bigcup_{i=1}^m X_i = X \\ r_i: X_i \rightarrow (-\infty, +\infty] \end{array} \right.$$

function in block of
coordinates

CD (differentiable)
guarantee
 $\nabla f(\bar{x}) = 0$ (stationary point)

→ Algorithm is quite similar to CD Algo

Application (Matrix Factorization)

$$e) \left\{ \begin{array}{l} \min_{W \in \mathbb{R}^{m \times r}} \frac{1}{2} \|V - WH\|_F^2 \\ \text{s.t. } W \geq 0 \quad H \geq 0 \end{array} \right.$$

non-convex
CD method can only guarantee
Stationary point (local minimizer)
(element-wise, not PSD)

separable

$V \in \mathbb{R}^{m \times n}, \quad W \in \mathbb{R}^{m \times r} \quad H \in \mathbb{R}^{r \times n}$

Calculation:

$$\begin{aligned} ① \quad & \frac{1}{2} \|V - WH\|_F^2 \\ &= \frac{1}{2} \|V - \sum_{i=1}^r W_{\cdot i} H_{i \cdot}\|_F^2 \end{aligned}$$

$$= \frac{1}{2} \|V - W_{\cdot i} H_{i \cdot} - W_{\cdot -i} H_{-i \cdot}\|_F^2$$

$$:= \frac{1}{2} \|\Delta - W_{\cdot i} H_{i \cdot}\|_F^2 \quad \left(\text{where } \Delta := V - W_{\cdot -i} H_{-i \cdot} \in \mathbb{R}^{m \times n} \right)$$

$$= \frac{1}{2} \langle \Delta - W_{\cdot i} H_{i \cdot}, \Delta - W_{\cdot i} H_{i \cdot} \rangle_F$$

$$(\text{where } \langle A, B \rangle_F := \text{Tr}(A^T B))$$

$$= \frac{1}{2} \|\Delta\|_F^2 + \frac{1}{2} \|W_{\cdot i} H_{i \cdot}\|_F^2 - \langle \Delta, W_{\cdot i} H_{i \cdot} \rangle_F$$

$$= \frac{1}{2} \|\Delta\|_F^2 + \frac{1}{2} \text{Tr} (W_i H_i \cdot H_i^\top W_i^\top) - \dots$$

$$= \frac{1}{2} \|\Delta\|_F^2 + \frac{1}{2} \text{Tr} (H_i^\top W_i^\top W_i H_i) - \dots$$

$\downarrow \text{Tr}(H_i \Delta^\top W_i)$

$$= \frac{1}{2} \|\Delta\|_F^2 + \frac{1}{2} \|H_i^\top\|_2^2 \cdot \|W_i\|_2^2 - \text{Tr} (\Delta^\top W_i H_i)$$

$$= \frac{1}{2} \|\Delta\|_F^2 + \frac{1}{2} \|H_i^\top\|_2^2 \|W_i\|_2^2 - \begin{cases} < \Delta^\top W_i, H_i > \\ < \Delta H_i^\top, W_i > \end{cases}$$

$\Delta \in \mathbb{R}^{m \times n}$

② a) fix H & W_{-i} , update W_i by:

$$\rightarrow \tilde{W}_i = \underset{w_i}{\operatorname{argmin}} \quad \frac{1}{2} \|w_i\|_2^2 - \frac{\langle w_i, \tilde{\Delta} \tilde{H}_i^\top \rangle}{\|\tilde{H}_i^\top\|_2^2} + \delta_{\mathbb{R}_+^m}(w_i)$$

$$= \underset{w_i}{\operatorname{argmin}} \quad \frac{1}{2} \|w_i - \frac{\tilde{\Delta} \tilde{H}_i^\top}{\|\tilde{H}_i^\top\|_2^2}\|_2^2 + \delta_{\mathbb{R}_+^m}(w_i)$$

$$= P_{\delta_{\mathbb{R}_+^m}(\cdot)} \left(\frac{\tilde{\Delta} \tilde{H}_i^\top}{\|\tilde{H}_i^\top\|_2^2} \right)$$

$$= \Pi_{\mathbb{R}_+^m} \left(\frac{(V - \tilde{W}_{(-i)} \tilde{H}_{(-i)}) \tilde{H}_i^\top}{\|\tilde{H}_i^\top\|_2^2} \right)$$

$$\rightarrow \tilde{H}_i^\top = \underset{H_i}{\operatorname{argmin}} \quad \frac{1}{2} \|H_i\|_2^2 - \frac{\langle \tilde{\Delta}^\top \tilde{W}_i, H_i \rangle}{\|\tilde{W}_i\|_2^2} + \delta_{\mathbb{R}_+^n}(H_i)$$

$$= \underset{H_i}{\operatorname{argmin}} \quad \frac{1}{2} \|H_i - \frac{\tilde{\Delta}^\top \tilde{W}_i}{\|\tilde{W}_i\|_2^2}\|_2^2 + \delta_{\mathbb{R}_+^n}(H_i)$$

$$= P_{\delta_{\mathbb{R}_+^n}(\cdot)} \left(\frac{\tilde{\Delta}^\top \tilde{W}_i}{\|\tilde{W}_i\|_2^2} \right)$$

$$= \Pi_{\mathbb{R}_+^n} \left(\frac{(V - \tilde{W}_{(-i)} \tilde{H}_{(-i)})^\top \tilde{W}_i}{\|\tilde{W}_i\|_2^2} \right)$$

$$\text{Rmk: } \underset{x}{\underbrace{\min : f(x)}}$$

if f is non-convex, but differentiable

then we still have $\nabla f(\bar{x}) = 0$ (necessary condition for optimality)

which means \bar{x} will be a stationary point,
not a global optimizer

$$\underset{x}{\min : f(x) + \sum_{i=1}^n r_i(x_i)} := F(x)$$

if f is non-convex, but differentiable

then we still have $0 \in \partial F(\bar{x})$

"stationary point"