

08/25/22 Thursday.

Least-Squares.

$$\vec{x}^* = \underset{\vec{x}}{\operatorname{arg\,min}} \|A\vec{x} - \vec{b}\|^2 = (A^T A)^{-1} A^T \vec{b}.$$

$$A^T(A\vec{x} - \vec{b}) = 0, A^T A \text{ invertible}$$

08/30/22 Tuesday.

Vector.  $\vec{x} \in \mathbb{R}^n$ .

2-norm | Euclidean Norm  $\|\vec{x}\|_2 = \sqrt{\sum_{i=1}^n x_i^2}$  distance to the origin.

Vector Norm.

Vector Space  $X$ .  $X \rightarrow \mathbb{R}$  is a norm if  $\|\vec{x}\| \geq 0 \quad \forall \vec{x} \in X$ , and  $\|\vec{x}\| = 0$  if  $\vec{x} = 0$ .

• triangular inequality.  $\|\vec{x} + \vec{y}\| \leq \|\vec{x}\| + \|\vec{y}\| \quad \forall \vec{x}, \vec{y} \in X$ .

•  $\|\lambda \vec{x}\| = |\lambda| \|\vec{x}\| \quad \forall \lambda \in \mathbb{R}, \vec{x} \in X$ .

Example:  $\ell_p$ -norm.  $\|\vec{x}\|_p = \left( \sum_{i=1}^n |x_i|^p \right)^{\frac{1}{p}}, 1 \leq p < \infty$ .

$$p=1, \|\vec{x}\|_1 = \sum_{i=1}^n |x_i|$$

$$p=\infty, \|\vec{x}\|_\infty = \max_{i=1, \dots, n} |x_i|.$$

Cauchy-Schwarz Inequality:  $\langle \vec{x}, \vec{y} \rangle = \vec{x}^T \vec{y} = \|\vec{x}\|_2 \|\vec{y}\|_2 \cos \theta$ .

$$\langle \vec{x}, \vec{y} \rangle \leq \|\vec{x}\|_2 \|\vec{y}\|_2. \quad |\vec{x}^T \vec{y}| \leq \|\vec{x}\|_2 \|\vec{y}\|_2.$$

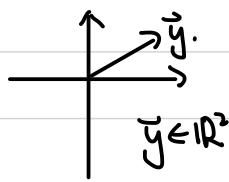
Hölder's Inequality:  $p, q \leq 1$  s.t.  $\frac{1}{p} + \frac{1}{q} = 1$ .

$$|\vec{x}^T \vec{y}| \leq \sum_{i=1}^n |x_i y_i| \leq \|\vec{x}\|_p \|\vec{y}\|_q.$$

$$\hookrightarrow \|\vec{x}\|_2 \|\vec{y}\|_2 \cos \theta.$$

optimization:  $\max \vec{x}^T \vec{y} \quad \|\vec{x}\|_p \leq 1$  given  $\vec{y} \in \mathbb{R}^n$  fixed.

$p=2$ .  $\|\vec{x}\|_2 \leq 1 \Rightarrow$  unit circle.



unit vector in the direction of  $\vec{y}$ .

$$\vec{x} = \vec{y} / \|\vec{y}\|_2.$$

$$p=\infty \cdot \begin{array}{c} \uparrow \\ \square \\ \downarrow \end{array} \rightarrow \|\vec{x}\|_\infty = \max_{i=1, \dots, n} |x_i|$$

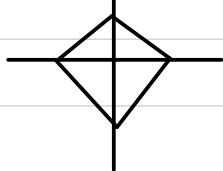
$$\vec{x}^T \vec{y} = \vec{x}_1 \vec{y}_1 + \vec{x}_2 \vec{y}_2 + \dots + \vec{x}_n \vec{y}_n \cdot \text{ independent: max each term.} \Rightarrow \text{max the sum.}$$

$$x_i = \text{sgn}(y_i) \quad \vec{x} = \text{sgn}(\vec{y}).$$

$$x_i y_i = \text{sgn}(y_i) \cdot y_i = |y_i|$$

$$\text{so } \max_{i=1}^n |x_i y_i| = \|\vec{x}^T \vec{y}\|_1 = \sum_{i=1}^n |y_i| = \|\vec{y}\|_1.$$

$$p=1 \cdot$$



$$|x_1| + |x_2| \leq 1. \quad \sum_{i=1}^n |x_i| \leq 1.$$

triangular inequality.

$$\vec{x}^T \vec{y} \leq \|\vec{x}^T \vec{y}\|_1 = \left\| \sum_{i=1}^n x_i y_i \right\|_1 \leq \sum_{i=1}^n |x_i y_i| = \sum_{i=1}^n |x_i| |y_i| \leq \sum_{i=1}^n |x_i| |y_{\max}| \leq |y_{\max}| = \|\vec{y}\|_1$$

$$y_1 \ y_2 \ \dots \ y_{\max} \ \dots \ y_n$$

$$\vec{x} = 0 \dots 0 \dots \text{sgn}(y_{\max}) \dots 0 \dots$$

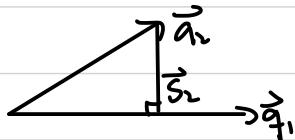
Gram-Schmidt: QR-decomposition.

$$X = \{\vec{a}_1, \vec{a}_2, \dots, \vec{a}_m\}, \text{ Basis.}$$

find orthonormal basis for the vector space.

$$\vec{q}_1 = \vec{a}_1 / \|\vec{a}_1\|_2. \quad \|\vec{q}_1\|_2 = 1.$$

$$\vec{q}_2 = \vec{s}_2 / \|\vec{s}_2\|_2.$$



must be normalized.

$$\text{project } \vec{a}_2 \text{ along } \vec{q}_1 = \vec{q}_1 \langle \vec{a}_2, \vec{q}_1 \rangle.$$

$$\vec{s}_2 = \vec{a}_2 - \vec{q}_1 \langle \vec{a}_2, \vec{q}_1 \rangle.$$

$$\vec{q}_3 = \vec{s}_3 / \|\vec{s}_3\|_2. \quad \vec{s}_3 = \vec{a}_3 - \langle \vec{a}_3, \vec{q}_1 \rangle \vec{q}_1 - \langle \vec{a}_3, \vec{q}_2 \rangle \vec{q}_2.$$

$$A = Q R \cdot \begin{matrix} \text{orthonormal} \\ \text{upper triangular.} \end{matrix} \quad [\vec{a}_1 \dots \vec{a}_m] = [\vec{q}_1 \dots \vec{q}_m] R.$$

09/01/22 Thursday.

## Fundamental Theorem of Linear Algebra.

$$A \in \mathbb{R}^{m \times n} \quad N(A) \oplus R(A^T) = \mathbb{R}^n.$$

$$R(A) \oplus N(A^T) = \mathbb{R}^m.$$

orthogonal decomposition thm.

let  $X$  be a vector space,  $S$  be a subspace

$$\forall \vec{x} \in X \quad \vec{x} = \vec{s} + \vec{r}. \quad \vec{s} \in S. \quad \vec{r} \in S^\perp \quad X = S \oplus S^\perp \quad S^\perp = \{\vec{r} \mid \langle \vec{r}, \vec{s} \rangle = 0, \forall \vec{s} \in S\}.$$

prove  $N(A) \oplus R(A^T) = \mathbb{R}^n$ :

$$N(A) \subseteq R(A^T)^\perp: (\text{let } \vec{x} \in N(A), A\vec{x} = 0. \text{ Want to show } \vec{x} \in R(A^T)^\perp \iff \exists \vec{w} \text{ s.t. } \vec{w} = A^T \vec{y}. \langle \vec{x}, \vec{w} \rangle = 0.)$$

$$\langle \vec{x}, \vec{w} \rangle = \langle \vec{x}, A^T \vec{y} \rangle = \vec{x}^T A^T \vec{y} \stackrel{A^T A \vec{x} = 0}{=} \vec{y}^T A \vec{x} = 0.$$

$$R(A^T)^\perp \subseteq N(A): (\text{let } \vec{x} \in R(A^T)^\perp \exists \vec{w} \text{ s.t. } \vec{w} = A^T \vec{y}. \langle \vec{x}, \vec{w} \rangle = 0. \text{ Want to show } A\vec{x} = 0.)$$

$$\langle \vec{x}, \vec{w} \rangle = \langle \vec{x}, A^T \vec{y} \rangle = \vec{x}^T A^T \vec{y} = \vec{y}^T A \vec{x} = 0. \text{ since it holds for all } \vec{y}, A\vec{x} = 0$$

## Minimum Norm Problem.

$$A\vec{x} = \vec{b}. \quad \boxed{A} \boxed{\vec{x}} = \boxed{\vec{b}} \quad \text{overdetermined.} \quad \# \text{ solutions} = 0.$$

$$\boxed{A} \boxed{\vec{x}} = \boxed{\vec{b}} \quad \text{underdetermined.} \quad \# \text{ solutions} = \infty.$$

$$\min \|\vec{x}\|_2^2 \text{ s.t. } A\vec{x} = \vec{b}$$

by FTLA,  $\vec{x} = \vec{y} + \vec{z}$ . s.t.  $\vec{y} \in N(A)$ ,  $\vec{z} \in R(A^T)$ .

$$A(\vec{y} + \vec{z}) = 0 + A\vec{z} = \vec{b}.$$

$$\|\vec{x}\|_2^2 = \|\vec{y}\|_2^2 + \|\vec{z}\|_2^2. \text{ Since } \vec{y} \perp \vec{z}.$$

$$\vec{z} = A^T \vec{w}. \quad A\vec{z} = \vec{b} \Rightarrow A A^T \vec{w} = \vec{b} \Rightarrow \vec{w} = (A A^T)^{-1} \vec{b}. \quad \text{Hence } \vec{z} = A^T (A A^T)^{-1} \vec{b}.$$

## Symmetric Matrices.

$$A \in \mathbb{S}^n \quad A = A^T. \quad A_{ij} = A_{ji}.$$

e.g. Covariance matrices.  $A = B \cdot B^T. \quad A^T = (B \cdot B^T)^T = B \cdot B^T$ .

Always diagonalizable.  $A = U \Lambda V^{-1}$  algebraic multiplicity  $\downarrow$  geometric multiplicity.

orthogonal.  $\det(\lambda I - A) / \# \text{ times } \lambda \text{ is a root.}$  dimension of  $N(\lambda I - A)$ .

$$\text{e.g. } \begin{pmatrix} 1 & 1 \\ 0 & 1 \end{pmatrix} \Rightarrow \begin{vmatrix} \lambda - 1 & -1 \\ 0 & \lambda - 1 \end{vmatrix} = (\lambda - 1)^2 = 0 \Rightarrow \lambda = 1. \text{ w/ multiplicity = 2.}$$

$$\text{Null}(\lambda I - A) = \begin{pmatrix} 0 & -1 \\ 0 & 0 \end{pmatrix} \begin{pmatrix} x \\ y \end{pmatrix} = \begin{pmatrix} -y \\ 0 \end{pmatrix} \Rightarrow \begin{pmatrix} 1 \\ 0 \end{pmatrix} \text{ dimension }= 1. \text{ Hence Not diagonalizable.}$$

## Spectral Thm

$A \in S^n$ .  $\lambda_i \in \mathbb{R}$ . with algebraic multiplicity  $m_i$ .  $\mathcal{E}_i = N(\lambda_i I - A)$ .

①  $\lambda_i \neq \lambda_j$ .  $\mathcal{E}_i \cap \mathcal{E}_j = \emptyset$ .

②  $\dim(\mathcal{E}_i) = m_i$ .

Proof: lemma:  $(\lambda, \mu)$  be an eigenpair for  $A$ .  $\exists$  orthonormal  $V$  st.  $UTAU = \begin{pmatrix} \lambda & 0 & \dots & 0 \\ 0 & \lambda & \dots & 0 \\ \vdots & \vdots & \ddots & 0 \\ 0 & 0 & \dots & \lambda \end{pmatrix}$   $B \in S^{n-1}$ .

by construction,  $V = \begin{bmatrix} \vec{u}_1 & \vec{u}_2 & \dots & \vec{u}_n \end{bmatrix}$

$$UTAU = \begin{bmatrix} \vec{u}_1^T \\ \vec{u}_2^T \\ \vdots \\ \vec{u}_n^T \end{bmatrix} \begin{bmatrix} A \end{bmatrix} \begin{bmatrix} \vec{u}_1 & \vec{u}_2 & \dots & \vec{u}_n \end{bmatrix} = \begin{bmatrix} \vec{u}_1^T \\ \vec{u}_2^T \\ \vdots \\ \vec{u}_n^T \end{bmatrix} \begin{bmatrix} \vec{u}_1^T & \vec{u}_2^T & \dots & \vec{u}_n^T \end{bmatrix} = \begin{bmatrix} \lambda & 0 & \dots & 0 \\ 0 & \lambda & \dots & 0 \\ \vdots & \vdots & \ddots & 0 \\ 0 & 0 & \dots & \lambda \end{bmatrix} = \lambda I_n.$$

Let  $B = V^T \Lambda V \in \mathbb{R}^{(n-1) \times (n-1)}$   
orthonormal  
diagonal

$$V U T A U V^T = V V^T \Lambda V^T = \begin{bmatrix} \lambda & 0 \\ 0 & \lambda \end{bmatrix}$$

$$\begin{bmatrix} \lambda & 0 & \dots & 0 \\ 0 & \lambda & \dots & 0 \\ \vdots & \vdots & \ddots & 0 \\ 0 & 0 & \dots & \lambda \end{bmatrix} = B = U^T A U.$$

→ orthonormal vectors.

For matrix, orthogonal = orthonormal.

② SVD.

09/06/2022 Tuesday. Eigenvalue decomposition = diagonalization

Spectral theorem:  $A \in S^n$ .  $A = U \Lambda U^T$   
diagonal  
orthonormal (unitary)

$$A = \underbrace{\begin{bmatrix} \vec{u}_1 & \vec{u}_2 & \dots & \vec{u}_r & \dots & \vec{u}_n \end{bmatrix}}_{\text{range}} \begin{bmatrix} \lambda_1 & & & & & 0 \\ & \lambda_2 & & & & 0 \\ & & \ddots & & & 0 \\ & & & \lambda_r & & 0 \\ & & & & \ddots & 0 \\ & & & & & 0 \end{bmatrix} \underbrace{U^T}_{\text{null space}}$$

Variational Characteristics of eigenvalue for symmetric matrix.

$$R = \frac{\vec{x}^T A \vec{x}}{\vec{x}^T \vec{x}}. \quad A \in S \quad \lambda_{\min}(A) \leq \frac{\vec{x}^T A \vec{x}}{\|\vec{x}\|_2^2} \leq \lambda_{\max}(A).$$

↓  
 Rayleigh coeff.  $\lambda_{\max}(A) = \max_{\|\vec{x}\|_2=1} \vec{x}^T A \vec{x} = \vec{x}$  be the eigenvector corr. to  $\lambda_{\max}$ .

$\lambda_{\min}(A) = \min_{\|\vec{x}\|_2=1} \vec{x}^T A \vec{x} = \vec{x}$  be the eigenvector corr. to  $\lambda_{\min}$ .

$$A = U \Lambda U^T. \quad U \vec{x} = \vec{y}. \quad \vec{x}^T A \vec{x} = \vec{x}^T U \Lambda U^T \vec{x} = \vec{y}^T \Lambda \vec{y} = \sum_{i=1}^n \lambda_i y_i^2 \leq \lambda_{\max} \sum_{i=1}^n y_i^2 = \lambda_{\max} \|\vec{y}\|_2^2 = \lambda_{\max} \|\vec{x}\|_2^2.$$

$\|\vec{x}\|_2^2 = 1 \Rightarrow \|\vec{y}\|_2^2 = 1$

"=" when  $\vec{x}$  corr. to  $\lambda_{\max}$ .

Principal Component Analysis.

$$\vec{x}_1, \dots, \vec{x}_n \quad \vec{x}_i \in \mathbb{R}^P \quad \text{Data Matrix} \quad X = \begin{bmatrix} -\vec{x}_1^T \\ -\vec{x}_2^T \\ \vdots \\ -\vec{x}_n^T \end{bmatrix} \quad n \times p.$$

$$C = \frac{1}{n} X X^T = \frac{1}{n} \begin{bmatrix} \|\vec{x}_1\|^2 & \langle \vec{x}_1, \vec{x}_2 \rangle \\ \langle \vec{x}_2, \vec{x}_1 \rangle & \|\vec{x}_2\|^2 \\ \vdots & \vdots \\ \|\vec{x}_n\|^2 & \langle \vec{x}_n, \vec{x}_1 \rangle \end{bmatrix} \quad D = \frac{1}{n} X^T X. \quad \text{Symmetric.}$$

want to find direction of PC1  $\vec{w}$ .  $\|\vec{w}\|_2 = 1$ .

projection onto  $\vec{w}$ :  $\langle \vec{x}_i, \vec{w} \rangle \vec{w}$ . Error:  $\|\vec{x}_i - \langle \vec{w}, \vec{x}_i \rangle \vec{w}\|^2 = e_i^2$ . Avg Error:  $\frac{1}{n} \sum e_i^2$ .

$$\begin{aligned} \|\vec{x}_i - \langle \vec{w}, \vec{x}_i \rangle \vec{w}\|^2 &= (\vec{x}_i - \langle \vec{w}, \vec{x}_i \rangle \vec{w})^T (\vec{x}_i - \langle \vec{w}, \vec{x}_i \rangle \vec{w}) \\ &= \|\vec{x}_i\|_2^2 + \langle \vec{w}, \vec{x}_i \rangle^2 \underbrace{\|\vec{w}\|_2^2}_{=1} - 2 \langle \vec{w}, \vec{x}_i \rangle^2 \\ &= \|\vec{x}_i\|_2^2 - \langle \vec{w}, \vec{x}_i \rangle^2 \end{aligned}$$

$$\frac{1}{n} \vec{w}^T + \sum_{i=1}^n \|\vec{x}_i\|_2^2 - \frac{1}{n} \sum_{i=1}^n \langle \vec{w}, \vec{x}_i \rangle^2 \Rightarrow \frac{1}{n} \vec{w}^T + \sum_{i=1}^n \langle \vec{w}, \vec{x}_i \rangle^2 = \frac{1}{n} \|\vec{x} \cdot \vec{w}\|^2 = \frac{1}{n} (\vec{x} \cdot \vec{w})^T (\vec{x} \cdot \vec{w}) = \frac{1}{n} \vec{w}^T \vec{x}^T \vec{x} \vec{w} \text{ Symmetric.}$$

so  $\vec{w}$  is the eigenvector corresponding to  $\lambda_{\max}$  of  $\vec{x}^T \vec{x}$ .

Singular Value Decomposition.

$$A \in \mathbb{R}^{m \times n} \quad A = \underbrace{U \Sigma V^T}_{m \times m \quad m \times n \quad n \times n}$$

$$A = \underbrace{U \Sigma V^T}_{m \times m} \quad \Sigma = \begin{bmatrix} \sigma_1 & & \\ & \ddots & \\ & & \sigma_r & 0 \\ & & & \vdots \\ & & & 0 \end{bmatrix} \quad \text{Rank } A = r.$$

$U$ -eigenvectors of  $A^T A$ .

$V$ -eigenvectors of  $A A^T$ .

$$A = U \begin{bmatrix} \sigma_1 & & \\ & \ddots & \\ & & \sigma_r & 0 \\ & & & \vdots \\ & & & 0 \end{bmatrix} V^T$$

$$A = \sigma_1 \vec{u}_1 \vec{v}_1^T + \sigma_2 \vec{u}_2 \vec{v}_2^T + \dots + \sigma_r \vec{u}_r \vec{v}_r^T$$

$$A \in \mathbb{R}^{m \times n} \text{ consider } A^T A \quad \lambda_1 \geq \lambda_2 \geq \dots \geq \lambda_r > 0 \quad \lambda_{r+1} = \dots = \lambda_n = 0.$$

$\vec{v}_1, \vec{v}_2, \dots, \vec{v}_n$  eigenvectors.

$$A^T A \vec{v}_i = \lambda_i \vec{v}_i \quad V = \begin{bmatrix} \vec{v}_1 & \vec{v}_2 & \dots & \vec{v}_n \end{bmatrix} \quad \text{define } \sigma_i = \sqrt{\lambda_i}, \vec{u}_i \text{ s.t. } A \vec{v}_i = \sigma_i \vec{u}_i. \quad i \leq r.$$

want to prove  $\vec{u}_i$  and  $\vec{u}_j$  are orthonormal.

$$\vec{u}_i^T \vec{u}_j = \frac{(A \vec{v}_i)^T}{\sigma_i} \frac{(A \vec{v}_j)}{\sigma_j} = \frac{1}{\sigma_i \sigma_j} \vec{v}_i^T A^T A \vec{v}_j = \frac{1}{\sigma_i \sigma_j} \vec{v}_i^T \lambda_j \vec{v}_j = \frac{\lambda_j}{\sigma_i \sigma_j} \vec{v}_i^T \vec{v}_j = \begin{cases} 1 & i=j \\ 0 & i \neq j \end{cases} \quad \text{ATA} \in \mathbb{R}^{n \times n} \text{ orthogonal}$$

Givann-Schmidt to find remaining  $\vec{u}_i$ .  $U: m \times m$ .

$$A \cdot V_r = \begin{bmatrix} \vec{u}_1 & \dots & \vec{u}_r \end{bmatrix} \begin{bmatrix} \sigma_1 & & \\ & \ddots & \\ & & \sigma_r \end{bmatrix} \quad A \vec{v}_i = \sigma_i \vec{u}_i.$$

$$AV = US$$

SVD is not unique e.g. GS.

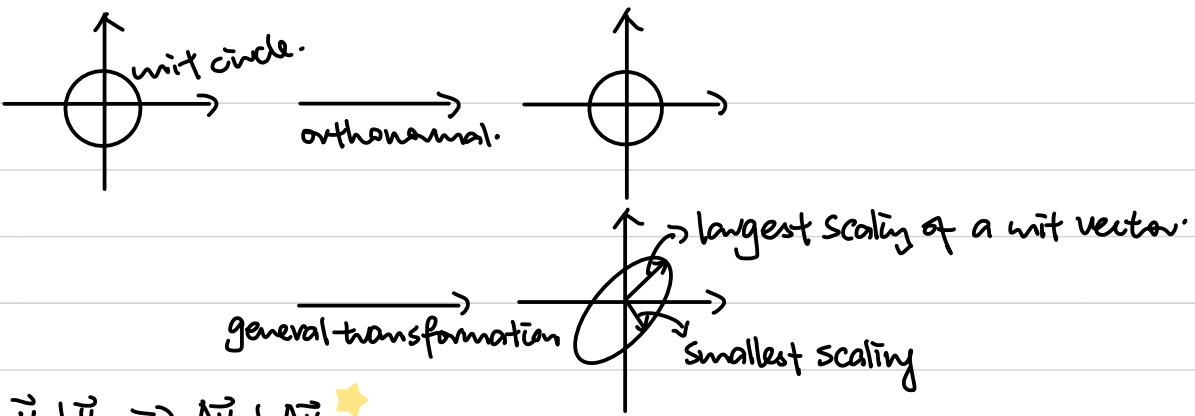
(not unique even if order the  $\lambda$ : repeated  $\lambda \Rightarrow$  multiple bases).

09/08/22 Thursday.

Geometry of SVD

$$A = U \Sigma V^T \quad A \vec{x} = U \Sigma V^T \vec{x} \quad \begin{array}{l} \rightarrow \text{scaling} \\ \downarrow \text{rotation/reflection} \\ \text{rotation/reflection} \end{array}$$

$\lambda_i = 0 \mid$  Not full rank  $\Rightarrow$  collapse dimension.  
e.g. line.



$$\vec{v}_1 \perp \vec{v}_2 \Rightarrow A\vec{v}_1 \perp A\vec{v}_2.$$

orthogonal after transformation  $\Rightarrow$  rotate/reflect/scale vectors.

$U^T \vec{x}$ : change  $\vec{x}$  into a basis s.t.  $A$  is just scaling.

$\Sigma V^T \vec{x}$ : do the scaling.

$U\Sigma V^T \vec{x}$ : change back to the original basis.

for arbitrary  $\vec{x}$ , think of  $\vec{x} = \lambda_1 \vec{v}_1 + \lambda_2 \vec{v}_2$ .

### Low Rank Approximation.

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

#### Matrix-Norm.

$\hookrightarrow$  block of data.  $\hookrightarrow$  Frobenius Norm:  $\|A\|_F = \sqrt{\sum_{i=1}^m \sum_{j=1}^n |a_{ij}|^2} = \sqrt{\text{trace}(A^T A)}$ .

$\hookrightarrow$  operator.  $\hookrightarrow$  invariant to orthonormal trans.  $\|UAU^T\|_F = \|AU\|_F = \|A\|_F$ .

$\hookrightarrow$  Spectral Norm /  $\ell_2$ -norm:  $\|A\|_2 = \frac{\max}{\|\vec{x}\|_2=1} \|A\vec{x}\|_2$ . max scaling of unit vector.  
 $= \frac{\max}{\|\vec{x}\|_2=1} \sqrt{\vec{x}^T A^T A \vec{x}} = \sqrt{\lambda_{\max}(A^T A)} = \sigma_{\max}(A)$ .

$$\|A^{-1}\|_2 = \frac{1}{\sigma_{\min}(A)}$$

### Eckart - Young - Mirsky Thm

$$A \in \mathbb{R}^{m \times n}. \quad \sigma_1 \geq \sigma_2 \geq \dots \geq \sigma_n. \quad A = U \Sigma V^T = \sum_{i=1}^n \sigma_i \vec{u}_i \vec{v}_i^T. \quad A_k = \sum_{i=1}^k \sigma_i \vec{u}_i \vec{v}_i^T.$$

$$\text{argmin}_{B \in \mathbb{R}^{m \times n}} \text{Rank}(B) = k. \quad \|A - B\|_F = A_k.$$

$$\text{argmin}_{B \in \mathbb{R}^{m \times n}} \text{Rank}(B) = k. \quad \|A - B\|_2 = A_k.$$

$$\text{Proof: } \|A - A_k\|_2 = \left\| \sum_{i=k+1}^n \sigma_i \vec{u}_i \vec{v}_i^T \right\|_2 = \sigma_{k+1}$$

$$\|A - B\|_2 \geq \|(A - B)\vec{w}\|_2 \quad \forall \vec{w} \in \text{Null}(B). \quad \|(A - B)\vec{w}\|_2 = \|A\vec{w}\|_2.$$

$$\text{Consider } V_{k+1} = [\vec{v}_1 \vec{v}_2 \dots \vec{v}_{k+1}] \quad \text{Rank}(V_{k+1}) = k+1. \quad \text{Dim}(\text{Null}(B)) = n - k.$$

$(n-k) + (k+1) = n+1 > n \Rightarrow \exists$  dimension overlap in  $\text{Range}(V_{k+1})$  and  $\text{Null}(B)$ .

$$\text{Suppose } \vec{w} \in \text{Range}(V_{k+1}). \quad \vec{w} = V \cdot \vec{a} = [V_{k+1} \quad V_{\text{rest}}] \begin{bmatrix} \vec{a}_1 \\ \vec{a}_2 \\ \vdots \\ \vec{a}_{k+1} \end{bmatrix} = \vec{a}_1 \vec{v}_1 + \vec{a}_2 \vec{v}_2 + \dots + \vec{a}_{k+1} \vec{v}_{k+1}.$$

$$\|\vec{w}\|^2 = 1 \Rightarrow \sum_{i=1}^{k+1} \vec{a}_i^2 = 1.$$

$$\|A\vec{w}\|_2^2 = \|U \Sigma V^T \vec{w}\|_2^2 = \|U \Sigma \vec{z}\|_2^2 = \|\Sigma \vec{z}\|_2^2 = \sum_{i=1}^{k+1} \sigma_i^2 \vec{a}_i^2 \geq \sigma_{k+1}^2 \left( \sum_{i=1}^{k+1} \vec{a}_i^2 \right) = \sigma_{k+1}^2.$$

09/13/22 Tuesday.

$$\|AUV\|_F = \|UAV\|_F = \|AU\|_F. \quad \|A\|_F = \sqrt{\text{trace}(A^T A)}$$

$$\text{Proof: } \|AUV\|_F = \sqrt{\text{trace}((AU)^T AU)} = \sqrt{\text{trace}(U^T A^T AU)} = \sqrt{\text{trace}(U U^T A^T A)} = \sqrt{\text{trace}(A^T A)} = \|A\|_F.$$

$$\|A\|_F = \|\Sigma\|_F = \sqrt{\sum_{i=1}^k \sigma_i^2}$$

Proof: want  $\|A-B\|_F \geq \|A-A_k\|_F$ .

$$\|A-A_k\|_F = \left\| \sum_{i=k+1}^n \sigma_i \bar{u}_i \bar{v}_i^T \right\|_F = \sqrt{\sum_{i=k+1}^n \sigma_i^2}$$

$$\text{Want } \sum_{i=k+1}^n \sigma_i^2 / \|A-B\| \geq \sum_{i=k+1}^n \sigma_i^2 / \|A\|$$

enough to prove:

$$\sigma_i(A-B) \geq \sigma_{k+i}(A) \cdot \frac{\|A\|}{\sum_{i=k+1}^n \sigma_i^2(A)} \Rightarrow \sum_{i=k+1}^n \sigma_i^2(A-B) \geq \sum_{i=k+1}^n \sigma_i^2(A)$$

$$\sigma_{k+i}(A) = k+i \text{th largest SV of } A = \|A - A_{k+i-1}\|_2. \quad \text{where } A_j = \sum_{i=1}^j \sigma_i \bar{u}_i \bar{v}_i^T$$

$$(\text{let } A-B=C. \quad \sigma_i(A-B) = \sigma_i(C) = \|C-C_{i-1}\|_2.)$$

$\text{rank}(B) = k$ .

$$\sigma_{k+1}(B) = 0. \quad \|B-B_k\|_2 = 0.$$

$$\sigma_i(A-B) = \|C-C_{i-1}\|_2 + \|B-B_k\|_2 \geq \|C+B-C_{i-1}-B_k\|_2 = \|A-C_{i-1}-B_k\|_2.$$

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

$$\text{rank}(B_k) = k. \quad \text{rank}(C_{i-1}) \leq i-1. \quad (\text{let } D = C_{i-1} + B_k. \quad \text{rank}(D) \leq k+i-1.)$$

$$\sigma_i(A-B) = \|A-D\|_2. \quad \sigma_{k+i}(A) = \|A-A_{k+i-1}\|_2. \Rightarrow \sigma_i(A-B) \geq \sigma_{k+i}(A)$$

False Proof:  $\min_{\text{rank}(B)=k} \|A-B\|_F = \min \|U\Sigma V^T - B\|_F = \min \|\Sigma - V^T B V\|_F \times \min_{\substack{\Sigma = \text{diagonal} \\ \Sigma = \text{rank}(k)}} \|\Sigma - \bar{\Sigma}\|_F.$

$$\begin{pmatrix} 1 & 0 & 0 \\ 0 & 1 & 0 \\ 0 & 0 & 0 \end{pmatrix} \begin{pmatrix} 1 & 0 & 0 \\ 0 & 0 & 0 \\ 0 & 0 & 0 \end{pmatrix}.$$

ingredients:  $\text{rank}(B) = k$ .

Frob. Norm  $\Leftrightarrow$  Singular Values.

Spectral norm = max singular value.

09/15/2022 Thursday.

Vector Calculus.

Taylor's Theorem:  $f(x): \mathbb{R} \rightarrow \mathbb{R}$ .  $x_0 \in \mathbb{R}$  · fixed point.

$$f(x_0 + \Delta x) = f(x_0) + \frac{df}{dx} \Big|_{x=x_0} (\Delta x) + \frac{1}{2!} \frac{d^2 f}{dx^2} \Big|_{x=x_0} (\Delta x)^2 + \dots$$

Taylor's Thm for Vectors:  $f(\vec{x}): \mathbb{R}^n \rightarrow \mathbb{R}$

first-order approximation.

$$f(\vec{x}_0 + \Delta \vec{x}) = f(\vec{x}_0) + \nabla f \Big|_{\vec{x}=\vec{x}_0}^T \Delta \vec{x} + \frac{1}{2!} (\Delta \vec{x})^T \nabla^2 f \Big|_{\vec{x}=\vec{x}_0} \Delta \vec{x} + \dots$$

gradient

$$\nabla \vec{x} = b. \quad \text{e.g. } b = 0.$$

always a column vector.  
same dimension as  $\vec{x}$ .

Gradient

Hessian: often symmetric.

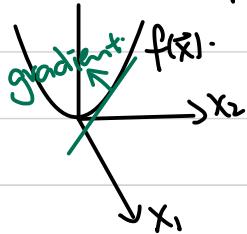
$$\nabla^2 f(\vec{x})_{ij} = \frac{\partial^2 f}{\partial x_i \partial x_j}$$

hyperplane.

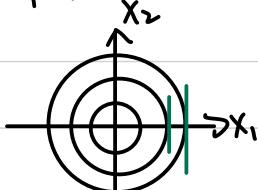
$\nabla f(\vec{x})$  captures the change wrt. to all components of  $\vec{x}$ .

$$\nabla f(\vec{x}) = \left[ \frac{\partial f}{\partial x_1}, \frac{\partial f}{\partial x_2}, \dots, \frac{\partial f}{\partial x_n} \right]^T$$

Example:  $f(\vec{x}) = \|\vec{x}\|_2^2$ .  $f: \mathbb{R}^n \rightarrow \mathbb{R}$ .



(level sets:  $f(\vec{x}) = \text{constant}$ ).



$$\nabla f(\vec{x}) = \begin{bmatrix} 2x_1 \\ 2x_2 \end{bmatrix} = 2\vec{x}.$$

$$\nabla^2 f(\vec{x}) = \begin{bmatrix} 2 & 0 \\ 0 & 2 \end{bmatrix} \xrightarrow{\frac{\partial(2x_1)}{\partial(x_2)}} = 0.$$

$$f(\vec{x} + \Delta\vec{x}) = (x_1^2 + x_2^2) + (2x_1, 2x_2) \begin{bmatrix} \Delta x_1 \\ \Delta x_2 \end{bmatrix} + \frac{1}{2} (\Delta x_1, \Delta x_2) \begin{bmatrix} 2 & 0 \\ 0 & 2 \end{bmatrix} \begin{bmatrix} \Delta x_1 \\ \Delta x_2 \end{bmatrix}$$

$$= x_1^2 + x_2^2 + 2x_1 \Delta x_1 + 2x_2 \Delta x_2 + \Delta x_1^2 + \Delta x_2^2$$

$\approx x_1^2 + x_2^2 + (x_1 + \Delta x_1)^2$  Error = 0 · quadratic approximation to quadratic function.

Example:  $f(\vec{x}) = \vec{x}^T \vec{a} = \sum_{i=1}^n x_i a_i$ .  $\nabla f(\vec{x}) = \begin{bmatrix} a_1 \\ a_2 \\ \vdots \\ a_n \end{bmatrix} = \vec{a}$ .

Example:  $f(\vec{x}) = \vec{x}^T A \vec{x} = \sum_i \sum_j x_i a_{ij} x_j$  terms including  $x_i: \sum_{j \neq i} x_i a_{ij} x_j + \sum_{j \neq i} x_j a_{ji} x_i + x_i^2 a_{ii}$ .  
 $\frac{\partial f}{\partial x_i} = \sum_{j \neq i} a_{ij} x_j + \sum_{j \neq i} x_j a_{ji} + 2x_i a_{ii} = \sum_j (a_{ij} + a_{ji}) x_j$   $\nabla f(\vec{x}) = (A + A^T) \begin{pmatrix} x_1 \\ \vdots \\ x_n \end{pmatrix}$ .  
 $\nabla^2 f = (A + A^T)$

The Main Thm: optimal is obtained at boundary OR derivative = 0.

$f: \mathbb{R}^n \rightarrow \mathbb{R}$ . differentiable everywhere.  $\Omega: \text{open set} \subseteq \mathbb{R}^n$ .  $\underset{\vec{x} \in \Omega}{\text{minimize}} f(\vec{x})$ .  $\frac{df}{d\vec{x}}(\vec{x}^*) = 0$ .

Proof:  $f(\vec{x}^* + \Delta\vec{x}) = f(\vec{x}^*) + \frac{df}{d\vec{x}}(\vec{x}^*) \Delta\vec{x} + \dots$

$f(\vec{x}^*)$  is the minimal  $\Rightarrow f(\vec{x}^*) \leq f(\vec{x}^* + \Delta\vec{x})$ .

open set  $\Rightarrow \exists \varepsilon > 0$  s.t.  $\forall \vec{x} \in \Omega$ ,  $|\vec{x} - \vec{x}^*| < \varepsilon$ .

(let  $0 < \Delta\vec{x} < \varepsilon$ .  $f(\vec{x}^*) \leq f(\vec{x}^*) + \frac{df}{d\vec{x}}(\vec{x}^*) \Delta\vec{x} + \dots \Rightarrow 0 \leq \frac{df}{d\vec{x}}(\vec{x}^*) \Delta\vec{x} + \dots$

$0 \leq \frac{df}{d\vec{x}}(\vec{x}^*) + \text{Const. } \frac{\Delta\vec{x}^2}{\Delta\vec{x}} + \text{Const. } \frac{\Delta\vec{x}^3}{\Delta\vec{x}}$ .  $\lim_{\Delta\vec{x} \rightarrow 0} \frac{\Delta\vec{x}^2}{\Delta\vec{x}} = 0$ ,  $0 \leq \frac{df}{d\vec{x}}(\vec{x}^*)$ .

wlog w/  $\vec{x}^* - \Delta\vec{x}$ ,  $\frac{df}{d\vec{x}}(\vec{x}^*) \leq 0$ .

09/20/2022 Tuesday.

Sensitivity / Perturbation Analysis.

$A\vec{x} = \vec{y}$ .  $A \in \mathbb{R}^{n \times n}$ , invertible.  $\vec{y} \rightarrow \vec{y} + \vec{\delta}y \Rightarrow \vec{x} \rightarrow \vec{x} + \vec{\delta}x$ . How big is  $\vec{\delta}x$ ?

$\frac{\|\vec{\delta}x\|_2}{\|\vec{x}\|_2} = \frac{\|\vec{y} + \vec{\delta}y - A\vec{x}\|_2}{\|\vec{x}\|_2} = \frac{\|\vec{\delta}y\|_2}{\|\vec{x}\|_2} \leq \|A^{-1}\|_2 \|\vec{\delta}y\|_2$ . by def of spectral Norm.

$\|\vec{\delta}y\|_2 = \|A\vec{\delta}x\|_2 \leq \|A\|_2 \|\vec{\delta}x\|_2 \Rightarrow \|\vec{\delta}x\|_2 \geq \|\vec{\delta}y\|_2 / \|A\|_2$ .

$A$  is invertible  $\Rightarrow A \neq 0$ .

$$\frac{\|\vec{x}\|_2}{\|\vec{x}\|_2} \leq \|A\|_2 \|A^{-1}\|_2 \left[ \frac{\|\vec{y}\|_2}{\|\vec{y}\|_2} \right] = \frac{1}{\sigma_{\min}} \frac{\|\vec{y}\|_2}{\|\vec{y}\|_2}$$

"Condition Number"

Least-Squares:  $\vec{x} = (\mathbf{A}\mathbf{A}^T)^{-1}\mathbf{A}\vec{b}$ .  $(\mathbf{A}\mathbf{A}^T)\vec{x} = \mathbf{A}\vec{b}$ .  $\in \mathbb{R}^{m \times n}$ .

(A doesn't need to be square, ATA should be invertible).

Condition Number ↑

if Add a constant (e.g. 200) to the dataset (A), the error will explode.

$(A^T A + \lambda I) \tilde{x} = A^T b$  works.  $\lambda$ -hyperparameter.

eigenvalue of  $A + \lambda I$ .

Shift property of eigenvalues :  $(A + \lambda I)\vec{v}_i = A\vec{v}_i + \lambda\vec{v}_i = \lambda_1\vec{v}_i + \lambda\vec{v}_i = (\lambda_1 + \lambda)\vec{v}_i$ .

$\Rightarrow$  increase  $\lambda_{\min}$  and  $\lambda_{\max} \Rightarrow$  decrease condition number.

## Ridge Regression

$$\min_{\tilde{X}} \|\tilde{A}\tilde{X} - \tilde{b}\|_2^2 + \lambda^2 \|\tilde{X}\|_2^2 \quad \lambda > 0.$$

$$f(\mathbf{x}) = (\mathbf{A}\mathbf{x} - \mathbf{b})^T(\mathbf{A}\mathbf{x} - \mathbf{b}) + \lambda^2 \mathbf{x}^T \mathbf{x} = \mathbf{x}^T \mathbf{A}^T \mathbf{A} \mathbf{x} - \mathbf{x}^T \mathbf{A}^T \mathbf{b} - \mathbf{b}^T \mathbf{A} \mathbf{x} + \lambda^2 \mathbf{x}^T \mathbf{x} + \mathbf{b}^T \mathbf{b}. \quad \text{PD} \Leftrightarrow \text{invertible.}$$

$$\nabla f(\bar{x}) = 2A^T A \bar{x} - 2(B^T A)^T + 2\lambda^2 \bar{x} = 0 \Rightarrow (A^T A + \lambda^2 I) \bar{x}^* = A^T B \Rightarrow \bar{x}^* = (A^T A + \lambda^2 I)^{-1} A^T B.$$

Another way of interpretation:

$$\begin{array}{c} A \\ \boxed{100000} \\ \boxed{100000} \end{array} = \begin{array}{c} 10 \\ 10 \end{array}$$

$$\vec{x}^* = \left( \begin{bmatrix} A^T & \lambda I \\ A & \lambda I \end{bmatrix} \right)^{-1} \begin{bmatrix} A^T & \lambda I \\ b \\ 0 \end{bmatrix} = (A^T A + \lambda^2 I)^{-1} (A^T b + 0)$$

We are confident by  $\lambda$  degree that  $x_2$  is (approx.) 0.

greater  $\lambda$ , greater confidence.

09/22/2022 Thursday.

## Tikhonov Regularization.

$$\frac{m}{\lambda} \|w_1(Ax - b)\|_2^2 + \|w_2(x - x_0)\|_2^2.$$

→ weight matrix → importance of each row.

$w_1 A \bar{x} = \bar{b}$  incorporate side info.

Incorporate Side Info.

$$y_i = g(\vec{x}_i) + \epsilon_i \quad (\vec{x}_i, y_i) \text{ data points} \quad \epsilon_i \sim N(0, \sigma^2) \quad f(\epsilon_i) = \frac{e^{-\epsilon_i^2/2\sigma^2}}{\sqrt{2\pi}\sigma} \quad i.i.d.$$

→ linear model.

$$g(\vec{x}_i) = \vec{x}_i^T \vec{w} \quad \text{Not LS since we know } z_i \sim N(0, \sigma_i^2).$$

### Maximum Likelihood Estimation

$\arg\max_{\vec{w}_0} f(Y_1=y_1, Y_2=y_2, \dots, Y_n=y_n | \vec{w}=\vec{w}_0)$   $\vec{w}$  that makes observed data most likely.

↑ density (integrate to get probability).

$= \arg\max_{\vec{w}_0} \prod_{i=1}^n f(y_i=y_i | \vec{w}=\vec{w}_0)$  since independent.

$$= f(\vec{x}_i^T \vec{w}_0 + z_i = y_i | \vec{w}=\vec{w}_0) = f(z_i = y_i - \vec{x}_i^T \vec{w}_0 | \vec{w}=\vec{w}_0) = \frac{e^{-\frac{(y_i - \vec{x}_i^T \vec{w}_0)^2}{2\sigma_i^2}}}{\sqrt{2\pi\sigma_i^2}}$$

$$= \arg\max_{\vec{w}_0} \frac{1}{(\sqrt{2\pi})^n \prod_{i=1}^n \sigma_i} e^{-\frac{\sum_{i=1}^n (y_i - \vec{x}_i^T \vec{w}_0)^2}{2\sigma_i^2}}$$

$$= \arg\min_{\vec{w}_0} \sum_{i=1}^n \frac{(y_i - \vec{x}_i^T \vec{w}_0)^2}{2\sigma_i^2} \quad \text{log is strictly increasing.}$$

$$= \arg\min_{\vec{w}_0} \|S(\vec{x}\vec{w}_0 - \vec{y})\|_2^2. \quad S = \begin{bmatrix} \frac{1}{\sigma_1} & 0 & \dots & 0 \\ 0 & \frac{1}{\sigma_2} & \dots & 0 \\ \vdots & \vdots & \ddots & \vdots \end{bmatrix} \quad \text{weighted LS.}$$

### Maximum a-posteriori (MAP) w/ prior on $\vec{w}$

$$\vec{y}_i = \vec{x}_i^T \vec{w} + z_i \quad z_i \sim N(0, \sigma_i^2) \quad \text{covariance = 0.}$$

$$\vec{w} \sim N(\vec{\mu}, \Sigma_w) \quad \vec{w} \sim N(\vec{\mu}, \Sigma_w) \quad \Sigma_w = \begin{bmatrix} \sigma_1^2 & 0 & \dots & 0 \\ 0 & \sigma_2^2 & \dots & 0 \\ \vdots & \vdots & \ddots & \vdots \end{bmatrix}.$$

What is most likely  $\vec{w}$ , given data  $y_1, y_2, \dots, y_n$ .

$$\arg\max_{\vec{w}} f(\vec{w} | \vec{y} = \vec{y}). \quad f(\vec{w} | \vec{Y} = \vec{y}) = \frac{f(\vec{Y} = \vec{y} | \vec{w}) f(\vec{w})}{\text{constant.}}$$

$$= \arg\max_{\vec{w}} f(\vec{Y} = \vec{y} | \vec{w}) \cdot f(\vec{w}) = \arg\max_{\vec{w}} \left( \prod_{i=1}^n f(Y_i = y_i | \vec{w}) \right) \cdot f(\vec{w}) = \arg\max_{\vec{w}} \prod_{i=1}^n \frac{\exp\{-\frac{(y_i - \vec{x}_i^T \vec{w})^2}{2\sigma_i^2}\}}{\sqrt{2\pi\sigma_i^2}} \cdot \frac{\exp\{-\frac{1}{2}(\vec{w} - \vec{\mu})^T \Sigma_w^{-1} (\vec{w} - \vec{\mu})\}}{(\sqrt{2\pi})^n |\Sigma_w|}.$$

$$= \arg\max_{\vec{w}} \exp\left\{ \sum_{i=1}^n \frac{(y_i - \vec{x}_i^T \vec{w})^2}{2\sigma_i^2} + -(\vec{w} - \vec{\mu})^T \Sigma_w^{-1} (\vec{w} - \vec{\mu}) \right\}. \quad \text{multivariate Gaussian Distribution. } w \sim \frac{e^{-\frac{\|\vec{w} - \vec{\mu}\|^2}{2\sigma^2}}}{\sqrt{2\pi\sigma^2}}$$

$$= \arg\min_{\vec{w}} \sum_{i=1}^n \frac{(y_i - \vec{x}_i^T \vec{w})^2}{2\sigma_i^2} + (\vec{w} - \vec{\mu})^T \Sigma_w^{-1} (\vec{w} - \vec{\mu})$$

$$= \arg\min_{\vec{w}} \|S(\vec{x}\vec{w} - \vec{y})\|_2^2 + \|\sqrt{\Sigma_w^{-1}}(\vec{w} - \vec{\mu})\|_2^2. \quad \Sigma_w^{-1} \in S \Rightarrow \Sigma_w^{-1} = \sqrt{\Sigma_w^{-1}} \cdot \sqrt{\Sigma_w^{-1}}$$

$\vec{\mu} = 0 \Rightarrow$  Ridge.  $\vec{\mu} \neq 0 \Rightarrow$  Tikhonov.

09/27/2022 Tuesday.

Convex Combinations:  $\sum_{i=1}^n \lambda_i \vec{x}_i$  is a convex combination of  $\vec{x}_i$  if  $\lambda_i \geq 0$  and  $\sum_{i=1}^n \lambda_i = 1$ .

Convex Set: A set  $C$  is convex if the line segment joining any two points in the set is contained in the set.

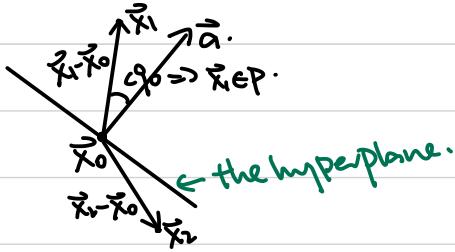
$\theta \vec{x}_1 + (1-\theta) \vec{x}_2 \in C$  if  $\vec{x}_1, \vec{x}_2 \in C$  and  $\theta \in [0, 1]$ . traverse along the line segment

Example:  $C = \{ \vec{x} | \vec{a}^T \vec{x} = b \}$ . A.k.A.  $\vec{a}^T (\vec{x} - \vec{x}_0) = 0$ . A.k.A. hyperplane.

Consider  $\vec{x}_1, \vec{x}_2 \in C$ ,  $\vec{x}_3 = \theta \vec{x}_1 + (1-\theta) \vec{x}_2$

$$\vec{a}^T \vec{x}_3 = \theta \vec{a}^T \vec{x}_1 + (1-\theta) \vec{a}^T \vec{x}_2 = b \Rightarrow \vec{x}_3 \in C \Rightarrow C \text{ is convex}$$

$$P: \{\vec{x} | \vec{a}^T (\vec{x} - \vec{x}_0) \geq 0\} \quad N: \{\vec{x} | \vec{a}^T (\vec{x} - \vec{x}_0) \leq 0\}.$$



Example:  $P = \{A | A \in S^n, A \text{ is PSD}\}$ .

Consider  $A_1, A_2 \in P$ .  $A_3 = \theta A_1 + (1-\theta) A_2$ .

$$\vec{x}^T A_3 \vec{x} = \theta (\vec{x}^T A_1 \vec{x}) + (1-\theta) (\vec{x}^T A_2 \vec{x}) \geq 0 \Rightarrow A_3 \in P \Rightarrow P \text{ is convex.}$$

Separating Hyperplane Theorem: let  $C, D$  be convex sets,  $C \cap D = \emptyset$ , then  $\exists$  hyperplane  $\vec{a}^T \vec{x} = b$  s.t.

$$\forall \vec{x} \in C, \vec{a}^T \vec{x} \geq b. \quad \forall \vec{x} \in D, \vec{a}^T \vec{x} \leq b. \quad \Rightarrow \text{only.}$$

*largest lower bound. → in case C,D are open sets.*

Proof:  $\text{dist}(C, D) = \inf \{ \| \vec{c} - \vec{d} \|_2 \mid \vec{c} \in C, \vec{d} \in D \}$ .  $C, D$  are closed, bounded sets. Assume  $\vec{c}, \vec{d}$  exist, are distinct.

Consider  $\vec{d} - \vec{c}$  as the normal vector passes through  $\frac{\vec{d} + \vec{c}}{2}$

$$f(\vec{x}) = (\vec{d} - \vec{c})^T (\vec{x} - \frac{\vec{d} + \vec{c}}{2}) \quad f(\vec{x}) = 0 \text{ is the hyperplane.}$$

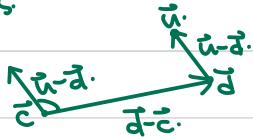
$$f(\vec{c}) = -\frac{1}{2} \|\vec{d} - \vec{c}\|^2 \quad f(\vec{x}) \leq 0 \quad \forall \vec{x} \in C. \quad f(\vec{d}) = \frac{1}{2} \|\vec{d} - \vec{c}\|^2 \quad f(\vec{x}) \geq 0 \quad \forall \vec{x} \in D.$$

If possible, let  $\vec{u} \in D$  s.t.  $f(\vec{u}) < 0$  (prove by contradiction).

$$f(\vec{u}) = (\vec{d} - \vec{c})^T (\vec{u} - \frac{1}{2}(\vec{d} + \vec{c}) + \vec{d} - \vec{c}) = (\vec{d} - \vec{c})^T ((\vec{u} - \vec{c}) + \frac{1}{2}(\vec{d} - \vec{c})) = \langle \vec{d} - \vec{c}, \vec{u} - \vec{c} \rangle + \frac{1}{2} \|\vec{d} - \vec{c}\|_2^2 < 0$$

*directions.*

$$\text{So } \langle \vec{d} - \vec{c}, \vec{u} - \vec{c} \rangle < 0$$



01/29/2022 Thursday.

tricky case: ①  $C = \{(x, y) | x \geq 0\}$ .  $D = \{(x, y) | xy \geq 1\}$ .

② \*open sets.



Consider  $\vec{p} = \vec{a} + t(\vec{u} - \vec{a}) = t\vec{u} + (1-t)\vec{a}$ .

$D$  is a convex set.  $\vec{u} \in D, \vec{a} \in D, \text{ if } t \in [0,1], \vec{p} \in D$ .

$$\|\vec{c} - \vec{p}\|_2^2 = \|\vec{c} - \vec{a} - t(\vec{u} - \vec{a})\|^2 = ((\vec{c} - \vec{a}) - t(\vec{u} - \vec{a}))^\top ((\vec{c} - \vec{a}) - t(\vec{u} - \vec{a})) = \|\vec{c} - \vec{a}\|^2 + t^2\|\vec{u} - \vec{a}\|^2 - 2t\langle \vec{c} - \vec{a}, \vec{u} - \vec{a} \rangle$$

$$\text{Want } \|\vec{c} - \vec{p}\|_2^2 < \|\vec{c} - \vec{a}\|_2^2 \Rightarrow t^2\|\vec{u} - \vec{a}\|^2 + 2t\langle \vec{c} - \vec{a}, \vec{u} - \vec{a} \rangle < 0$$

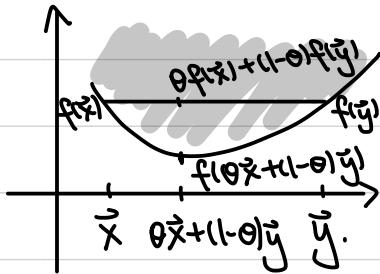
$$\Rightarrow t < \frac{-2\langle \vec{c} - \vec{a}, \vec{u} - \vec{a} \rangle}{\|\vec{u} - \vec{a}\|^2} \quad t \text{ can be arbitrary small.}$$

## Convex function.

upwards

Def: A function  $f: \mathbb{R}^n \rightarrow \mathbb{R}$  is convex if  $\text{dom}(f)$  is a convex set and  $\forall \vec{x}, \vec{y} \in \text{dom}(f)$  and  $\theta \in [0,1]$ , we have  $f(\theta\vec{x} + (1-\theta)\vec{y}) \leq \theta f(\vec{x}) + (1-\theta)f(\vec{y})$ . [Strictly convex].

Def: A function  $f: \mathbb{R}^n \rightarrow \mathbb{R}$  is concave if  $\text{dom}(f)$  is a convex set and  $\forall \vec{x}, \vec{y} \in \text{dom}(f)$  and  $\theta \in [0,1]$ , we have  $f(\theta\vec{x} + (1-\theta)\vec{y}) \geq \theta f(\vec{x}) + (1-\theta)f(\vec{y})$ . [Strictly concave].



Epigraph:  $\text{Epi } f = \{(x, t) \mid x \in \text{dom}(f), f(x) \leq t\}$ .

Prop:  $f$  is a convex function  $\Leftrightarrow \text{Epi } f$  is a convex set.

Jensen's inequality: if  $f$  is convex,  $\vec{x}_1, \vec{x}_2, \dots, \vec{x}_k \in \text{dom}(f)$ , and  $\theta_1, \theta_2, \dots, \theta_k \geq 0$  with  $\sum_{i=1}^k \theta_i = 1$ , then  $f(\theta_1\vec{x}_1 + \theta_2\vec{x}_2 + \dots + \theta_k\vec{x}_k) \leq \theta_1 f(\vec{x}_1) + \theta_2 f(\vec{x}_2) + \dots + \theta_k f(\vec{x}_k)$ .

first order condition: Suppose  $f$  is differentiable,  $f$  is convex iff  $\text{dom}(f)$  is convex and FOC  $\forall \vec{x}, \vec{y} \in \text{dom}(f), f(\vec{y}) \geq f(\vec{x}) + \nabla f(\vec{x})^\top (\vec{y} - \vec{x})$  tangent line below the fn.

Second order condition: Suppose  $f$  is twice differentiable.  $f$  is convex iff  $\text{dom}(f)$  is convex and  $\nabla^2 f(\vec{x})$  (Hessian) is PSD for all  $\vec{x} \in \text{dom}(f)$ .

Example of convex fn that is not differentiable.  $f(\vec{x}) = \|\vec{x}\|$ .

Implication: if  $\nabla f(\vec{x}^*) = 0$  and if  $f$  is convex,  $f(\vec{y}) \geq f(\vec{x}^*) + 0(\vec{y} - \vec{x}^*) \Rightarrow \vec{x}^*$  is the global minimum.

Proof of first order condition:

$$\Rightarrow f((1-t)\vec{x} + t\vec{y}) \leq (1-t)f(\vec{x}) + t f(\vec{y}) \Rightarrow t f(\vec{y}) \geq t f(\vec{x}) - f((1-t)\vec{x} + t\vec{y}).$$

$$\Rightarrow f(\vec{y}) \geq f(\vec{x}) + t[f(\vec{x} + t(\vec{y} - \vec{x})) - f(\vec{x})].$$

$$f'(\vec{x}) = \lim_{\Delta x \rightarrow 0} \frac{f(\vec{x} + \Delta x) - f(\vec{x})}{\Delta x} \stackrel{\Delta x = t(\vec{y} - \vec{x})}{=} \lim_{t \rightarrow 0} \frac{f(\vec{x} + t(\vec{y} - \vec{x})) - f(\vec{x})}{t(\vec{y} - \vec{x})}$$

$$\text{take limit } t \rightarrow 0, f(\vec{y}) \geq f(\vec{x}) + f'(\vec{x})(\vec{y} - \vec{x})$$

$\Leftarrow$ : let  $z = \theta x + (1-\theta)y$ .

$$f(x) \geq f(z) + f'(z)(x-z), \quad f(y) \geq f(z) + f'(z)(y-z)$$

$$\theta f(x) + (1-\theta)f(y) \geq \theta f(z) + \theta f'(z)(x-z) + (1-\theta)f(z) + (1-\theta)f'(z)(y-z).$$

$$= f(z) + f'(z)(\theta x - \theta z + y - z - \theta y + \theta z) = f(z) + f'(z)[\theta x + (1-\theta)y - z] = f(z).$$

10/04/2022 Tuesday. Strong convexity  $\Rightarrow$  strict convexity  $\Rightarrow$  convexity.

### Strict Convexity

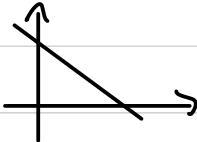
domain(f) is convex,  $\theta \in (0,1)$   $\forall x, y \in \text{domain}(f)$ .

Zeroth order condition:  $f(\theta \bar{x} + (1-\theta)\bar{y}) < \theta f(\bar{x}) + (1-\theta)f(\bar{y})$ .

First order condition:  $f'(\bar{y}) > f'(\bar{x}) + \nabla f(\bar{x})^T(\bar{y}-\bar{x})$  Taylor first order approximation.

Second order condition:  $\nabla^2 f(\bar{x}) \succ 0$  PD.  $A \succeq B \Rightarrow A - B \succeq 0$ .

Example:



both concave and convex. Not strictly.

### Strong Convexity

$f$  is differentiable, domain  $f$  is convex.  $\mu > 0$

$f$  is  $\mu$ -strongly convex if  $\forall \bar{x}, \bar{y} \in \text{domain}(f)$ ,  $\bar{x} \neq \bar{y}$   $f(\bar{y}) \geq f(\bar{x}) + \nabla f(\bar{x})^T(\bar{y}-\bar{x}) + \frac{\mu}{2} \|\bar{y}-\bar{x}\|^2$ .

Taylor second order approximation:  $f(\bar{y}) \approx f(\bar{x}) + \nabla f(\bar{x})^T(\bar{y}-\bar{x}) + \frac{1}{2}(\bar{y}-\bar{x})^T \nabla^2 f(\bar{x})(\bar{y}-\bar{x})$ .

quadratic fn under

$$= \frac{1}{2}(\bar{y}-\bar{x})^T \begin{bmatrix} \mu & \dots & \mu \\ \vdots & \ddots & \vdots \\ \mu & \dots & \mu \end{bmatrix} (\bar{y}-\bar{x})$$

$$\mu I \leq \nabla^2 f(\bar{x})$$

$$\mu I \leq \nabla^2 f(\bar{x})$$

the gradient is not changing too slow.

### Unconstrained optimization problem

#### Gradient Descent.

$$f(\bar{x} + \Delta \bar{x}) = f(\bar{x}) + \nabla f(\bar{x})^T \Delta \bar{x}$$

stepsize direction.

$$S>0 \quad f(\bar{x} + S\bar{v}) = f(\bar{x}) + S \langle \nabla f(\bar{x}), \bar{v} \rangle. \quad \text{maximize magnitude} \Rightarrow \bar{v} \text{ aligned to gradient.}$$

$$\text{minimize } f \Rightarrow \bar{v} = -\nabla f(\bar{x}).$$

+ if close to  $\bar{x}^*$ ,  $\nabla f(\bar{x})$  is small, step  $\langle \nabla f(\bar{x}), \bar{v} \rangle$  is small.

$$\bar{x}_{k+1} = \bar{x} - \eta \nabla f(\bar{x}_k). \quad \bar{x}_0: \text{initial point}. \quad \eta: \text{stepsize}. \quad \text{if } \bar{x}_k = \bar{x}^*, \bar{x}_{k+1} = \bar{x}^* - 0.$$

$$\text{Least squares: } f(x) = \|Ax - b\|^2. \quad \nabla f(x) = 2A^T(Ax - b).$$

$$\bar{x}^* = (A^T A)^{-1} A^T b. \quad \bar{x}_{k+1} = \underbrace{(-2\eta A^T A)}_{\text{eigenvalues.}} \bar{x}_k + 2\eta A^T b.$$

Why GD? - faster.

$$\begin{aligned}\vec{x}_{k+1} - \vec{x}^* &= \vec{x}_{k+1} - (\text{ATA})^{-1} A^T \vec{b} = (\mathbb{I} - 2\eta(\text{ATA})) \vec{x}_k + 2\eta(\text{ATA})(\text{ATA})^{-1} A^T \vec{b} - (\text{ATA})^{-1} A^T \vec{b} \\ &= (\mathbb{I} - 2\eta(\text{ATA})) \vec{x}_k - (\text{ATA})^{-1} A^T \vec{b}. \\ |\text{eigenvalues}| < 1 &\Rightarrow \text{Converge} \\ &= (\mathbb{I} - 2\eta(\text{ATA}))^k (\vec{x}_0 - (\text{ATA})^{-1} A^T \vec{b}).\end{aligned}$$

10/06/2022 Thursday.

L-smooth:  $\forall \vec{x}, \vec{y} \in \text{domain}, f(\vec{y}) \leq f(\vec{x}) + Df(\vec{x})^T(\vec{y} - \vec{x}) + \frac{L}{2} \|\vec{y} - \vec{x}\|_2^2$ .

↑  
the gradient is not changing too fast.

Thm:  $f: \mathbb{R}^n \rightarrow \mathbb{R}^n$   $\mu$ -strong convexity AND L-smooth.  $\min_{\vec{x} \in \mathbb{R}^n} f(\vec{x}) \quad \vec{x}_{t+1} = \vec{x}_t - \eta Df(\vec{x}_t)$ .  
then  $\|\vec{x}_{t+1} - \vec{x}^*\|_2^2 \leq C^{t+1} \|\vec{x}_0 - \vec{x}^*\|_2^2$  (choose  $\eta$  s.t.  $|C| < 1$ )  $\Rightarrow \lim_{t \rightarrow \infty} \vec{x}_t = \vec{x}^*$

Lemma: if  $f$  is L-smooth,  $\|Df(\vec{x})\|_2^2 \leq 2L(f(\vec{x}) - f(\vec{x}^*))$  gradient close to optimum is small.

Proof:  $f(\vec{x}^*) \leq f(\vec{x}) \quad f(\vec{x}^*) \leq f(\vec{x} - \frac{Df(\vec{x})}{L})$

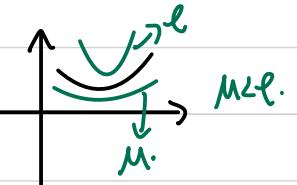
$$f(\vec{x} - \frac{Df(\vec{x})}{L}) \leq f(\vec{x}) + Df(\vec{x})^T(-\frac{Df(\vec{x})}{L}) + \frac{L}{2} \left\| -\frac{Df(\vec{x})}{L} \right\|_2^2 = f(\vec{x}) - \frac{1}{2L} \|Df(\vec{x})\|_2^2$$

$\mu$ -strong convexity: let  $\vec{y} = \vec{x}^*$   $Df(\vec{x})^T(\vec{x} - \vec{x}^*) \geq f(\vec{x}) - f(\vec{x}^*) + \frac{\mu}{2} \|\vec{x}^* - \vec{x}\|_2^2$ .

$$\begin{aligned}\text{Proof: } \|\vec{x}_{t+1} - \vec{x}^*\|_2^2 &= \|(\vec{x}_t - \vec{x}^*) - \eta Df(\vec{x}_t)\|_2^2 = \|\vec{x}_t - \vec{x}^*\|_2^2 + \eta^2 \|Df(\vec{x}_t)\|_2^2 - 2\eta Df(\vec{x})^T(\vec{x}_t - \vec{x}^*) \\ &\leq \|\vec{x}_t - \vec{x}^*\|_2^2 + \eta^2 \cdot 2L(f(\vec{x}_t) - f(\vec{x}^*)) - 2\eta(f(\vec{x}) - f(\vec{x}^*)) + \frac{\mu}{2} \|\vec{x}^* - \vec{x}\|_2^2.\end{aligned}$$

$$= (1 - \eta \cdot \mu) \|\vec{x}_t - \vec{x}^*\|_2^2 + (2\eta^2(1 - 2\eta)) [f(\vec{x}_t) - f(\vec{x}^*)].$$

choose  $\eta = \frac{1}{L}$   $\Rightarrow 2\eta^2(1 - 2\eta) = 0$ ; then  $\|\vec{x}_{t+1} - \vec{x}^*\|_2^2 \leq \underbrace{(1 - \frac{\mu}{L})^{t+1}}_{< 1} \|\vec{x}_0 - \vec{x}^*\|_2^2$ .



## Stochastic Gradient Descent.

$$f(\vec{x}) = \sum_{i=1}^m \frac{1}{m} f_i(\vec{x})$$

$$\vec{x}_{k+1} = \vec{x}_k - \eta_k Df_i(\vec{x}_k) \quad E[Df_i(\vec{x}_k)] = Df(\vec{x}_k)$$

$\curvearrowleft$   
pointwise maximize convex fns  $\rightarrow$  convex.

concave  $\rightarrow$  convex/concave.

$\curvearrowleft$   
pointwise minimize convex fns  $\rightarrow$  convex/concave. concave  $\rightarrow$  concave.



10/13/2022 Thursday.

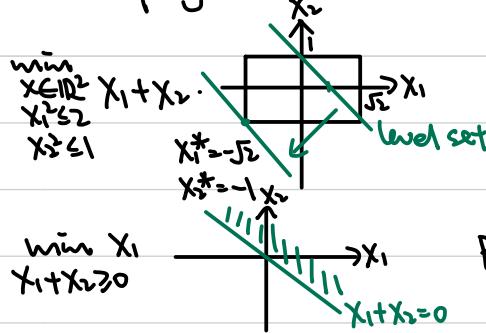
Active Constraint: optimum on the constraint / boundary.

Inactive Constraint:

Convex optimization problem: minimize convex/maximize concave.

infeasible optimization problem  $\rightarrow$  feasible set =  $\emptyset$ .  $p^* = \infty$

linear Program:



$\rightarrow c^T$  (orthogonal to level set).

$[c^T]x$  minimize  $\Rightarrow$  move in negative  $\nabla f(x) = c$  direction.

Then:  $\min_{\vec{x} \in X} \vec{c}^T \vec{x}$   $X$ : convex set, closed. If  $\vec{x}^*$  is an optimum, then  $\vec{x}^* \in \text{Boundary}(X)$ .

Proof: Assume  $\vec{x}^*$  is in the interior of the set  $X$ .

$$\forall \vec{\varepsilon} \quad \|\vec{x} - \vec{\varepsilon}\|_2 \leq r \Rightarrow \vec{\varepsilon} \in X. \quad (\text{Let } \vec{\varepsilon}_0 = -2 \cdot \vec{c} \quad 2 \cdot \|\vec{c}\|_2 > 0)$$

$$f_0(\vec{x}^* + \vec{\varepsilon}) = \vec{c}^T(\vec{x}^* + \vec{\varepsilon}) = \vec{c}^T \vec{x}^* - 2 \cdot \vec{c}^T \vec{c} < f_0(\vec{x}^*)$$

= -(the concave fn.)

Problem Transformation:

(maximize concave fn) = - (minimize convex fn).

① Addition of "slack" Variable.  $\min_{\vec{x} \in X} f_0(\vec{x}) \Leftrightarrow \min_{\vec{x} \in X} f_0(\vec{x}) + t$

$$\text{e.g. } \min_{\vec{x}} \|A\vec{x} - b\|_2^2 + \|\vec{x}\|_1. \text{ LASSO} \Leftrightarrow \|\vec{x}\|_1 + \|\vec{A}\vec{x} - \vec{b}\|_2^2 + \sum_i t_i.$$

②  $\phi(\vec{x})$  conti. + strictly increasing.

$$\min_{\vec{x}} f(\vec{x}) = \phi(f_0(\vec{x})) \Leftrightarrow \min_{\vec{x}} f(\vec{x}) = \phi(\phi(f_0(\vec{x}))) \quad \text{Should be the same.}$$

$$\text{e.g. } \min_{\substack{\vec{x}, \vec{t} \\ \vec{x} \geq 0}} \vec{t}^T \vec{t} \quad \text{Not convex.} \quad \Leftrightarrow \min_{\substack{\vec{x}, \vec{t} \\ \log(\vec{x}) \geq 0}} \vec{t}^T \vec{t} \quad \text{Convex.}$$

obj. strictly decreasing over domain  
f\_i nondecreasing over domain

$\Downarrow$   
equality constraint  $\Rightarrow$  inequality ( $\leq$ )

Duality:  $\rightarrow$  inequality constraint.

$$p^* = \min_{\vec{x}} \max_{\vec{\lambda}} f_0(\vec{x}) \quad \text{PRIMAL} \quad \text{Lagrangian: } L(\vec{x}, \vec{\lambda}, \vec{v}) = f_0(\vec{x}) + \sum_{i=1}^m \lambda_i f_i(\vec{x}) + \sum_{j=1}^p v_j h_j(\vec{x}) \quad \lambda_i \geq 0$$

$h_i(\vec{x}) = 0 \quad \forall i \in \{1, \dots, p\} \rightarrow$  equality constraint.

$$\min_{\vec{x}} L(\vec{x}, \vec{\lambda}, \vec{v}) \quad \text{s.t. } \vec{\lambda} \geq 0; \vec{v} := g(\vec{\lambda}, \vec{v})$$

L is affine fn  $\rightarrow$  both convex and concave. ✓

g is pointwise min of concave fn  $\Rightarrow$  concave.

prop:  $g(\vec{x}, \vec{\mu})$  is the lower bound of  $p^*$ .  $\forall \vec{x} \geq 0, \vec{\mu}, g(\vec{x}, \vec{\mu}) \leq p^*$ .

proof: consider  $\vec{x}$  some feasible set for the primal.  $f_i(\vec{x}) \leq 0 \Rightarrow g_i(\vec{x}) = 0$ .

$$L(\vec{x}, \vec{\lambda}, \vec{\mu}) = f_0(\vec{x}) + \sum_{i=0}^m \lambda_i f_i(\vec{x}) + \sum_{i=0}^m \mu_i h_i(\vec{x}) \leq f_0(\vec{x})$$

$$g(\vec{x}, \vec{\mu}) = \min_{\vec{\lambda}} L(\vec{x}, \vec{\lambda}, \vec{\mu}) \Rightarrow g(\vec{x}, \vec{\mu}) \leq p^*$$

$$M(\vec{x}) = f_0(\vec{x}) + \sum_{i=0}^m \mathbb{1}_{f_i(\vec{x}) \leq 0} f_i(\vec{x}) + \sum_{i=0}^m \mathbb{1}_{h_i(\vec{x}) \geq 0} h_i(\vec{x}). \quad \mathbb{1}_{f_i(\vec{x}) \leq 0} = \begin{cases} 0 & \text{if } f_i(\vec{x}) \leq 0 \\ \infty & \text{if } f_i(\vec{x}) > 0 \end{cases} \quad \mathbb{1}_{h_i(\vec{x}) \geq 0} = \begin{cases} 0 & \text{if } h_i(\vec{x}) = 0 \\ \infty & \text{if } h_i(\vec{x}) > 0 \end{cases}$$

*Hard threshold: encourage  $f_i(\vec{x}) \leq 0$ .*

$$\sum_{i=0}^m \mathbb{1}_{f_i(\vec{x}) \leq 0} \leftrightarrow \sum_{i=0}^m \lambda_i f_i(\vec{x}).$$

$$f_i(\vec{x}) \leq 0 \quad 0 \quad \geq \quad \leq 0$$

$$f_i(\vec{x}) > 0 \quad \infty \quad \geq \quad > 0$$

largest lower bound.

max  $\vec{x} \geq 0, \vec{\mu} \quad g(\vec{x}, \vec{\mu}) = d^*$ . DUAL problem to Primal.

\* Maximization of a concave fn, linear constraint.  $\Rightarrow$  Convex Program.

\* # Variable = # Constraints of primal.

\* Always Convex program regardless of primal  $\xrightarrow{f_0(\vec{x})}$ .

WEAK DUALITY:  $d^* \leq p^*$ . Duality Gap  $p^* - d^*$  STRONG DUALITY:  $d^* = p^*$   
\* Always hold regardless of convexity of  $f_0(\vec{x}), f_i(\vec{x})$ .

Example:  $\min_{\vec{x} \in \mathbb{R}^n} \vec{x}^T \vec{x} \quad A \in \mathbb{R}^{m \times n} \text{ mnn.}$

$$A\vec{x} = \vec{b} \quad \vec{x} \geq 0$$

$$L(\vec{x}, \vec{\mu}) = \vec{x}^T \vec{x} + \vec{\mu}^T (A\vec{x} - \vec{b}) \quad \nabla_{\vec{x}} L(\vec{x}, \vec{\mu}) = 2\vec{x} + A^T \vec{\mu} \Rightarrow \vec{x}^* = -\frac{1}{2} A^T \vec{\mu}$$

$$g(\vec{\mu}) = \min_{\vec{x}} L(\vec{x}, \vec{\mu}) = L(\vec{x}^*, \vec{\mu}) = \frac{1}{4} \vec{\mu}^T A A^T \vec{\mu} + \vec{\mu}^T (-\frac{1}{2} A A^T \vec{\mu} - \vec{b}) = -\frac{1}{4} \vec{\mu}^T A A^T \vec{\mu} - \vec{\mu}^T \vec{b}$$

$$\nabla_{\vec{\mu}} g(\vec{\mu}) = -\frac{1}{2} A A^T \vec{\mu} - \vec{b} = 0 \Rightarrow \vec{\mu}^* = -2(A A^T)^{-1} \vec{b} \Rightarrow \vec{x}^* = -\frac{1}{2} A^T (A A^T)^{-1} \vec{b} = A^T (A A^T)^{-1} \vec{b}$$

Strong duality by state since

Partitioning Problem s.t.  $x_i^j = 1 \quad i \in \{1, n\} \quad \vec{x}^T W \vec{x} \quad W \in \mathbb{R}^{n \times n} \quad$  Not Convex.  $\leftarrow$  discrete domains. no fi.

$$\vec{x}^T W \vec{x} = \sum_i \sum_j x_i w_{ij} x_j. w_{ij} \rightarrow \text{unhappiness of two people in the same group.}$$

$$L(\vec{x}, \vec{\mu}) = \vec{x}^T W \vec{x} + \sum_i \mu_i (x_i^2 - 1) = \vec{x}^T W \vec{x} + \vec{x}^T \text{diag}(\vec{\mu}) \vec{x} - \sum_i \mu_i = \vec{x}^T (W + \text{diag}(\vec{\mu})) \vec{x} - \sum_i \mu_i$$

$$g(\vec{\mu}) = \min_{\vec{x}} L(\vec{x}, \vec{\mu}) = \begin{cases} -\sum_i \mu_i & \text{if } W + \text{diag}(\vec{\mu}) \text{ PSD} \\ -\infty & \text{otherwise: have negative eigenvalues.} \end{cases}$$

DUAL:  $W + \text{diag}(\vec{\mu})$  is PSD  $\Rightarrow -\sum_i \mu_i \quad \vec{\mu}^* = \lambda_{\min}(W) \quad p^* \geq n \cdot \lambda_{\min}(W)$ .

10/20/2022 Thursday.

Minimax Inequality: Sets X, Y F is any function.  $\min_{x \in X} \max_{y \in Y} f(x, y) \geq \max_{y \in Y} \min_{x \in X} f(x, y)$  first

Proof:  $\forall x_0 \in X, y_0 \in Y$ . Define  $h(y_0) := \min_{x \in X} f(x, y_0)$ ,  $g(x_0) := \max_{y \in Y} f(x_0, y)$ .

$$\forall x_0, y_0 \quad h(y_0) \leq f(x_0, y_0) \leq g(x_0) \Rightarrow \max_{y \in Y} h(y_0) \leq \min_{x \in X} g(x_0)$$

$$d^* = \max_{\bar{x} \geq 0, \bar{\mu}} \min_{\bar{X}, \bar{\lambda}} L(\bar{X}, \bar{\lambda}, \bar{\mu})$$

$$\max_{\bar{x} \geq 0, \bar{\mu}} \left[ f_0(\bar{x}) + \sum_{i=1}^m \lambda_i f_i(\bar{x}) + \sum_{i=1}^r \bar{\mu}_i h_i(\bar{x}) \right] = \begin{cases} f_0(\bar{x}) & \text{if } \bar{x} \text{ is feasible.} \\ \infty & \text{otherwise.} \end{cases} \Rightarrow p^* = \min_{\bar{X}} \max_{\bar{x} \geq 0, \bar{\mu}} L(\bar{X}, \bar{\lambda}, \bar{\mu})$$

$$\text{Hence } p^* \geq d^*$$

weak duality holds with Non-Convex dimension things.

Slater's Condition Strong duality holds if  $\exists \bar{x}_0 \in \text{relative interior}(D)$  s.t.  $\nabla f(\bar{x}_0) \leq 0$   
CONVEX (constraints, obj funs, domains).

refined slater. CONVEX:  $f_1, f_2, \dots, f_k$  are affine.

$$\exists \bar{x}_0 \in \text{relative interior}(D) \text{ s.t. } f_i(\bar{x}_0) \leq 0 \text{ for } i \in \{1, k\} \text{ AND } f_i(\bar{x}_0) < 0 \text{ for } i \in \{k+1, n\}$$

→ strong duality always hold if feasible.

linear program & duality.

$$\min_{A\bar{x} \leq \bar{b}} \bar{c}^T \bar{x} \quad L(\bar{x}, \bar{\lambda}) = \bar{c}^T \bar{x} + \bar{\lambda}^T (A\bar{x} - \bar{b}) = (\bar{A}^T \bar{\lambda} + \bar{c})^T \bar{x} - \bar{b}^T \bar{\lambda}.$$

$$g(\bar{x}) = \min_{\bar{\lambda}} L(\bar{x}, \bar{\lambda}) = \begin{cases} -\infty & \text{if } \bar{A}^T \bar{\lambda} + \bar{c} \neq 0 \\ -\bar{b}^T \bar{\lambda} & \text{otherwise} \end{cases} \Rightarrow d^* = \max_{\substack{\bar{x} \geq 0 \\ \bar{A}^T \bar{\lambda} + \bar{c} = 0}} -\bar{b}^T \bar{\lambda}. \text{ Also Linear program.}$$

200kg merlot. 300kg shiraz. Blend 1: 4kgM + 1kgS \$20. Blend 2: 2kgM + 3kgs \$15.

$$\max_{q_1, q_2} 20q_1 + 15q_2 \text{ s.t. } 4q_1 + 2q_2 \leq 200, q_1 + 3q_2 \leq 300, q_1, q_2 \geq 0$$

+ sell grapes directly:  $\max_{q_1, q_2} 20q_1 + 15q_2 + \lambda_1(200 - 4q_1 - 2q_2) + \lambda_2(300 - q_1 - 3q_2)$  Lagrangian

$$= q_1, q_2 \geq 0 \quad \frac{(20 - 4\lambda_1 - \lambda_2)q_1 + (15 - 2\lambda_1 - 3\lambda_2)q_2 + 200\lambda_1 + 300\lambda_2}{\lambda_1 = 0 \Rightarrow \text{No Blend 1} \quad \lambda_2 = 0 \Rightarrow \text{No Blend 2.}}$$

At indifference point:  $\lambda_1, \lambda_2 \rightarrow \text{shadow prices.}$

$$\min_{\lambda_1, \lambda_2 \geq 0} 200\lambda_1 + 300\lambda_2 \text{ s.t. } 20 - 4\lambda_1 - \lambda_2 = 0 \text{ AND } 15 - 2\lambda_1 - 3\lambda_2 = 0 \text{ DUAL.}$$

10/25/2022 Tuesday.

Certificate Property.

$(\lambda_1, \mu_1)$  - feasible dual point  $x_1$  - feasible primal

$$p^* \geq g(\lambda_1, \mu_1) \Rightarrow f_0(x_1) - p^* \leq f_0(x_1) - g(\lambda_1, \mu_1) \iff p^* \in [g(\lambda_1, \mu_1), f_0(x_1)] \text{ w/ strong duality. } d^* \in [g(\lambda_1, \mu_1), f_0(x_1)].$$

## Complementary Slackness. \*regardless of the convexity.

Primal optimal ( $\vec{x}^*$ ) Dual Optimal ( $\vec{\lambda}^*, \vec{\mu}^*$ )

$$g(\vec{\lambda}^*, \vec{\mu}^*) = \min_{\vec{x}} \left( f_0(\vec{x}) + \sum_{i=1}^m \lambda_i^* f_i(\vec{x}) + \sum_{i=1}^p \mu_i^* h_i(\vec{x}) \right) \stackrel{(1)}{\leq} f_0(\vec{x}^*) + \sum_{i=1}^m \lambda_i^* f_i(\vec{x}^*) + \sum_{i=1}^p \mu_i^* h_i(\vec{x}^*) \stackrel{(2)}{\leq} f_0(\vec{x}^*).$$

Assume strong duality holds:  $p^* = d^*$      $p^* = f_0(\vec{x}^*)$      $d^* = g(\vec{\lambda}^*, \vec{\mu}^*)$

$$\stackrel{(1)}{\Rightarrow} \min_{\vec{x}} \left( f_0(\vec{x}) + \sum_{i=1}^m \lambda_i^* f_i(\vec{x}) + \sum_{i=1}^p \mu_i^* h_i(\vec{x}) \right) = f_0(\vec{x}^*) + \sum_{i=1}^m \lambda_i^* f_i(\vec{x}^*) + \sum_{i=1}^p \mu_i^* h_i(\vec{x}^*) \Rightarrow \vec{x}^* \text{ also minimizes } L.$$

$$\stackrel{(2)}{\Rightarrow} \sum_{i=1}^m \lambda_i^* f_i(\vec{x}^*) = 0 \Rightarrow \begin{cases} \lambda_i^* > 0, f_i(\vec{x}^*) = 0 \\ \lambda_i^* = 0, f_i(\vec{x}^*) > 0 \end{cases} \begin{matrix} \text{not necessarily the only minimizer.} \\ \text{Complementary slackness.} \end{matrix}$$

## KKT Condition.

Strong duality holds; regardless of convexity; differentiable  $f_i$ 's,  $h_i$ 's.

Necessary conditions of  $\vec{x}^*, \vec{\lambda}^*, \vec{\mu}^*$  optimality: Contrapositive: if Not B, Not A.  
then:

$$\textcircled{1} \quad \forall i \in \{1, m\}, f_i(\vec{x}^*) \leq 0.$$

$$\textcircled{2} \quad \forall i \in \{1, p\}, h_i(\vec{x}^*) = 0.$$

$$\textcircled{3} \quad \forall i \in \{1, m\}, \lambda_i^* \geq 0.$$

$$\textcircled{4} \quad \forall i \in \{1, m\}, \lambda_i^* f_i(\vec{x}^*) = 0.$$

$$\textcircled{5} \quad \nabla f_0(\vec{x}^*) + \sum \lambda_i^* \nabla f_i(\vec{x}^*) + \sum \mu_i^* \nabla h_i(\vec{x}^*) = 0 \quad \text{by Main Thm. * respect to } x_i \text{ for } i \in \{1, n\}.$$

fo convex, domain convex.

Convex problems ( $\forall i \in \{1, m\}$ ,  $f_i$  convex;  $h_i$  affine); differentiable  $f_i$ 's,  $h_i$ 's.

## Sufficient conditions

$\vec{x}, \vec{\lambda}, \vec{\mu}$  satisfying KKT; convex;  $h_i$  affine; Slater's condition  $\Leftrightarrow \vec{x}, \vec{\lambda}, \vec{\mu}$  are optimal.

10/27/2022 Thursday.

Proof sufficient condition:

Consider  $L(\vec{x}, \vec{\lambda}, \vec{\mu}) = f_0(\vec{x}) + \sum_{i=1}^m \lambda_i^* f_i(\vec{x}) + \sum_{i=1}^p \mu_i^* h_i(\vec{x}) \rightarrow$  convex fn in  $\vec{x}$ . by  $\textcircled{5}$ ,  $\vec{x}$  is a minimizer

$$g(\vec{x}, \vec{\mu}) = \min_{\vec{\lambda}} L(\vec{x}, \vec{\lambda}, \vec{\mu}) = f_0(\vec{x}) + \sum \lambda_i^* f_i(\vec{x}) + \sum \mu_i^* h_i(\vec{x}) \stackrel{\textcircled{4}}{\leq} \stackrel{\textcircled{5}}{=} f_0(\vec{x})$$

Since  $p^* \geq g(\vec{x}, \vec{\mu})$ ,  $p^* = f_0(\vec{x})$   $\stackrel{p^* \text{ cannot } > f_0(\vec{x})}{\Rightarrow} d^* = g(\vec{x}, \vec{\mu}) = f_0(\vec{x})$  certificate property.

Convex.

Strong duality holds as long as feasible.

Linear Program:

$$\min_{\vec{x}} \vec{c}^T \vec{x} \text{ s.t. } A\vec{x} \leq \vec{b}, \quad \vec{a}_i^T \vec{x} = b_i \Rightarrow \vec{a}_i^T \vec{x} \geq b_i \text{ AND } \vec{a}_i^T \vec{x} \leq b_i$$

$$\cdot \vec{a}_i^T \vec{x} \geq b_i \Rightarrow -\vec{a}_i^T \vec{x} \leq -b_i.$$

$$\text{Standard format: } \min_{\vec{x}} \vec{c}^T \vec{x} \text{ s.t. } \begin{array}{l} A\vec{x} \leq \vec{b} \\ \vec{x} \geq 0 \end{array} \quad \sum_{j=1}^n a_{ij} x_j \leq b_i \Rightarrow \sum_{j=1}^n a_{ij} x_j + s_i = b_i \quad s_i \geq 0.$$

$$\cdot x_i = x_i^+ - x_i^- \text{ where } x_i^+ \geq 0, x_i^- \geq 0.$$

$$\text{e.g. } \min_{x_1, x_2} 2x_1 + 4x_2 \text{ s.t. } \begin{array}{l} 3x_1 + 2x_2 = 14 \\ x_1 \geq 0 \end{array} \Rightarrow \min_{x_1, x_1^+, x_1^-, x_2, x_2^+} 2x_1 + 4x_2^+ - 4x_2^- \text{ s.t. } \begin{array}{l} 3x_1 + 2x_2^+ - 2x_2^- = 14 \\ x_1, x_1^+, x_1^-, x_2, x_2^+ \geq 0 \end{array}$$

$$x_1 + x_1^+ - x_1^- - x_2 = 3$$

slack Variable.

Def: Polyhedron  $\{\vec{x} \in \mathbb{R}^n \mid A\vec{x} \leq \vec{b}\}$   $A \in \mathbb{R}^{m \times n}$   $\vec{b} \in \mathbb{R}^m$ .  
in "standard form"  $\{\vec{x} \in \mathbb{R}^n \mid C\vec{x} = \vec{d}; \vec{x} \geq 0\}$ .

Def: Extreme Points polyhedron = P  $\vec{x} \in P$  is an extreme point (vertex of P) if we cannot find two vectors  $\vec{y}, \vec{z} \neq \vec{x}, \vec{y}, \vec{z} \in P$  and  $\lambda \in (0,1)$  s.t.  $\vec{x} = \lambda \vec{y} + (1-\lambda) \vec{z}$ .

Thm: P has an extreme point iff P does not contain a line.

11/01/2022 Tuesday.

Thm: Consider a LP  $\min_{\vec{x}} \vec{c}^T \vec{x}$  s.t.  $A\vec{x} \leq \vec{b}$ , if P has an extreme point and opt. solution exists and is finite, there exists an optimal solution that is an extreme point of P.

Proof:  $P = \{\vec{x} \mid A\vec{x} \leq \vec{b}\}$ ,  $P^* = \min_{\vec{x} \in P} \vec{c}^T \vec{x}$ .  $Q = \{\vec{x} \mid A\vec{x} \leq \vec{b}, \vec{c}^T \vec{x} = P^*\}$ . Set of all opt. solutions.

Q is also a polyhedron.

P doesn't contain a line,  $Q \subset P \Rightarrow Q$  doesn't contain a line  $\Rightarrow Q$  has an extreme point.

Suppose  $\vec{v}$  is a vertex of Q, while  $\vec{v}$  is a vertex of P.

Prove by contradiction.  $\exists \vec{y}, \vec{z} \in P$ ,  $\vec{y}, \vec{z} \neq \vec{v}$  s.t.  $\lambda \vec{y} + (1-\lambda) \vec{z} = \vec{v}$ .  $\lambda \in (0,1)$ .

$\vec{v} \in Q \Rightarrow \vec{c}^T \vec{v} = P^*$   $\lambda \vec{c}^T \vec{y} + (1-\lambda) \vec{c}^T \vec{z} = \vec{c}^T \vec{v} = P^*$  since  $\vec{c}^T \vec{y} \geq P^*$ ,  $\vec{c}^T \vec{z} \geq P^* \Rightarrow \vec{c}^T \vec{y} = \vec{c}^T \vec{z} = P^*$ .

$\exists \vec{y}, \vec{z} \in Q, \vec{y}, \vec{z} \neq \vec{v}$  s.t.  $\lambda \vec{y} + (1-\lambda) \vec{z} = \vec{v}$  contradicts.

Quadratic programs

convex iff H is PSD.

$$P^* = \min_{\vec{x}} \frac{1}{2} \vec{x}^T H \vec{x} + \vec{c}^T \vec{x} \text{ s.t. } A\vec{x} \leq \vec{b}, \vec{c}^T \vec{x} = d.$$

$$P^* = \min_{\vec{x}} \frac{1}{2} \vec{x}^T H \vec{x} + \vec{c}^T \vec{x} + d. \text{ s.t. } A\vec{x} \leq \vec{b} \text{ s.t. } \vec{c}^T \vec{x} = d.$$

If  $H = H^T$  symmetric, H is PSD  $\Rightarrow QP$  is convex.

If  $H$  is symmetric, at least one negative eigenvalue.

$p^* = -\infty$ , choose  $\tilde{x}$  be the eigenvector corr. to negative eigenvalue.

If  $H$  is symmetric,  $H$  is PSD,  $\tilde{c} \in R(H)$

$$f(\tilde{x}) = \frac{1}{2} \tilde{x}^T H \tilde{x} + \tilde{c}^T \tilde{x} = \frac{1}{2} (\tilde{x} - \tilde{x}_0)^T H (\tilde{x} - \tilde{x}_0) + d = \frac{1}{2} \tilde{x}^T H \tilde{x} + \frac{1}{2} \tilde{x}_0^T H \tilde{x}_0 - \tilde{x}_0^T H \tilde{x} + d.$$

$$\tilde{c} = -H^T \tilde{x}_0 = -H \tilde{x}_0 \quad d = -\frac{1}{2} \tilde{x}_0^T H \tilde{x}_0 \quad \text{Want to choose } \tilde{x} = \tilde{x}_0$$

If  $H$  is invertible,  $\tilde{x}^* = -H^{-1} \tilde{c}$ .

If  $H$  has null space,  $\tilde{x}^* = \text{any } \tilde{x}_0 \text{ s.t. } -H \tilde{x}_0 = \tilde{c}$ .

$$H = U \Sigma V^T \quad H \text{ has rank } r. \quad H = U \begin{bmatrix} \Sigma_{r \times r} & 0 \\ 0 & 0 \end{bmatrix} V^T.$$

$$\text{Moore-Penrose Pseudoinverse} \quad H^+ = V \begin{bmatrix} \Sigma_{r \times r}^{-1} & 0 \\ 0 & 0 \end{bmatrix} U^T \quad HH^+ = U_r U_r^T \quad H^+ H = V_r V_r^T \quad H H^+ H = H.$$

$$\tilde{x}^* = -H^+ \tilde{c} + \tilde{r} \quad \tilde{r} \in N(H)$$

If  $H$  is symmetric,  $H$  is PSD,  $\tilde{c} \notin R(H)$

$$\tilde{c} = -H \tilde{x}_0 + \tilde{r} \quad \tilde{r} \in N(H^T) = N(H)$$

$$f(\tilde{r}) = \frac{1}{2} \tilde{r}^T H \tilde{r} + \tilde{c}^T \tilde{r} = 0 + (-H \tilde{x}_0 + \tilde{r})^T \tilde{r} = -\tilde{x}_0^T H^T \tilde{r} + \tilde{r}^T \tilde{r} = \|\tilde{r}\|_2^2$$

$$f(k\tilde{r}) = k\|\tilde{r}\|_2^2 \quad k \rightarrow -\infty, f(k\tilde{r}) \rightarrow -\infty \Rightarrow p^* = -\infty.$$

Equality constrained OP  $\Rightarrow$  Uncstrained OP.

Application 1: Linear Control LQR.

$$\tilde{x}(t+1) = A \tilde{x}(t) + B u(t). \quad \tilde{x}(t) = A^t \tilde{x}(0) + \sum_{i=0}^{t-1} A^{t-i-1} B \cdot u(i)$$

$$\text{goal: reach } \bar{g} \text{ by time } T. \Rightarrow \min_{\substack{\tilde{x}(t), t=0, \dots, T \\ u(t), t=0, \dots, T-1}} \|\tilde{x}(T) - \bar{g}\|_2^2 + \sum_{t=0}^{T-1} \|u(t)\|_2^2 \quad \text{s.t. } \tilde{x}(t) = A^t \tilde{x}(0) + \sum_{i=0}^{t-1} A^{t-i-1} B \cdot u(i) \quad \text{LP.}$$

Application 2:  $\tilde{y} = \tilde{x} + \text{noise}$ . find  $\tilde{x}$  s.t. it does not change consecutively.

$$D = \begin{bmatrix} 1 & 0 & \dots & 0 \\ 0 & 1 & \dots & 0 \\ \vdots & \vdots & \ddots & \vdots \\ 0 & \dots & 0 & 1 \end{bmatrix} \quad D \tilde{x} = \begin{bmatrix} x_1 - x_0 \\ x_2 - x_1 \\ \vdots \\ x_n - x_{n-1} \end{bmatrix} \quad p^* = \min_{\tilde{x}} \|\tilde{y} - \tilde{x}\|_2^2 \quad \text{s.t. } \text{cond}(D \tilde{x}) \leq k. \quad \text{cond} \rightarrow \# \text{ nonzero entries.}$$

$$\text{relax to } p^* = \min_{\tilde{x}} \|\tilde{y} - \tilde{x}\|_2^2 \quad \text{s.t. } \|D \tilde{x}\|_1 \leq 2. \quad \text{relax to LP. encourage sparsity.}$$

LPCOP convex  $\subset$  QCQP convex  $\subset$  SOCP  $\subset$  {all convex programs}.

Quadratically-Constrained Quadratic Programming. QCQP.

11/03/2022 Thursday.

$$p^* = \min_{\tilde{x}} \tilde{x}^T H_0 \tilde{x} + 2 \tilde{c}^T \tilde{x} + d \quad \text{s.t. } \tilde{x}^T H_i \tilde{x} + 2 \tilde{c}_i^T \tilde{x} + d_i \leq 0 \quad i=1, \dots, m$$

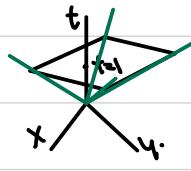
SOCP

$$\tilde{x}^T H_j \tilde{x} + 2 \tilde{c}_j^T \tilde{x} + d_j \leq 0 \quad j=1, \dots, q.$$

Def: A set of points  $C \subseteq \mathbb{R}^n$  is a cone iff  $\lambda \tilde{x} \in C$  if  $\tilde{x} \in C$  AND  $\lambda \geq 0$ . Convex iff  $H_0, H_i \text{ PSD}, H_j = 0$ .

Def: convex cone iff  $\lambda \tilde{x} \in C$  if  $\tilde{x} \in C$  AND  $\theta_1, \theta_2 \geq 0$  AND  $\theta_1 \tilde{x}_1 + \theta_2 \tilde{x}_2 \in C$ .

e.g.  $C = \{(x, y) \mid |x| \leq y\}$ .  $C = \{(x, y) \mid y \geq 0\}$ .



Def: Polyhedron Cone  $\{(\vec{x}, t) \mid A\vec{x} \leq \vec{b}t, t \in \mathbb{R}, t \geq 0\}$ .

Def: Ellipsoidal Cone  $\{(\vec{x}, t) \mid \|A\vec{x} + \vec{b}t\|_2 \leq ct\}$

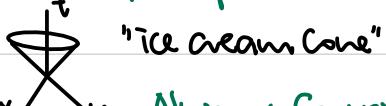
ellipse  $\vec{x}^T P \vec{x} + \vec{q}^T \vec{x} + r \leq 0$   $P$  is PSD.

$$\|A\vec{x} + \vec{b}\|_2^2 \leq c^2 \Rightarrow \vec{x}^T A^T A \vec{x} + 2\vec{b}^T A \vec{x} + \vec{b}^T \vec{b} - c^2 \leq 0$$

PSD

Def: Second-order Cone Special case of Ellipsoidal Cone.

in  $\mathbb{R}^3 \{(\vec{x}_1, \vec{x}_2, t) \mid \sqrt{\vec{x}_1^2 + \vec{x}_2^2} \leq t\}$ .



Define: Second Order Cone Program  $\min_{\vec{x}} \vec{c}^T \vec{x}$  s.t.  $\|A_i \vec{x} + \vec{b}_i\|_2 \leq c_i^T \vec{x} + d_i \quad i=1, 2, \dots, m$  Always Convex.

$\min_{\vec{x}} \vec{c}^T \vec{x}$  s.t.  $\|A_i \vec{x} + \vec{b}_i\|_2 \leq c_i^T \vec{x} + d_i \quad i=1, 2, \dots, m \Rightarrow (A_i \vec{x} + \vec{b}_i, c_i^T \vec{x} + d_i)$  must belong to SOC.

Example:  $\min_{\vec{x}} \sum_{i=1}^m \|A_i \vec{x} - \vec{b}_i\|_2 \iff \min_{\vec{x}, y} \sum_{i=1}^m y_i \text{ s.t. } \|A_i \vec{x} - \vec{b}_i\|_2 = y_i$  relax  $\min_{\vec{x}, y} \sum_{i=1}^m y_i \text{ s.t. } \|A_i \vec{x} - \vec{b}_i\|_2 \leq y_i$  same p.

if  $A_k \vec{x} - \vec{b}_k = y_k - S$   $S \geq 0$ , decrease  $y_k$ .

Example:  $\min_{\vec{x}} \max_{i=1, \dots, m} \|A_i \vec{x} - \vec{b}_i\|_2 \iff \min_{\vec{x}, y} y \text{ s.t. } \|A_i \vec{x} - \vec{b}_i\|_2 \leq y \quad i=1, \dots, m$ .

Facility location problem.

Example: Trilateration GPS.

$$\Delta_i = t_i^R - t_i^T$$

Packet transmission  $t_i^T$  Packet receive  $t_i^R$  offset of time =  $t_i^R - t_i^T$  unknown.

time of flight  $f_i = t_i^{\text{time}} - t_i^T = t_i^R + \delta - t_i^T = \Delta_i + \delta$ .

Distance  $C \cdot f_i = C\Delta_i + C\delta$ .  $C$  = speed of light.

$\|\vec{x} - \vec{a}_i\|_2 = C\Delta_i + C\delta$   $\vec{x}$ : my position  $\vec{a}_i$ : position of satellite  $i$ .

wl 4 satellites  $\|\vec{x} - \vec{a}_1\|_2^2 - \|\vec{x} - \vec{a}_2\|_2^2, \|\vec{x} - \vec{a}_2\|_2^2 - \|\vec{x} - \vec{a}_3\|_2^2, \|\vec{x} - \vec{a}_3\|_2^2 - \|\vec{x} - \vec{a}_4\|_2^2 \Rightarrow 3$  linear equations.

in 2D, 3 unknowns  $\Rightarrow$  Solvable.

wl 3 satellites  $\min_{\vec{x}} \delta$  s.t. ①  $2(\vec{a}_3 - \vec{a}_1)^T \vec{x} + 2C^2(\Delta_3 - \Delta_1)\delta = C^2(\Delta_1^2 - \Delta_3^2) + \|\vec{a}_3\|_2^2 - \|\vec{a}_1\|_2^2$

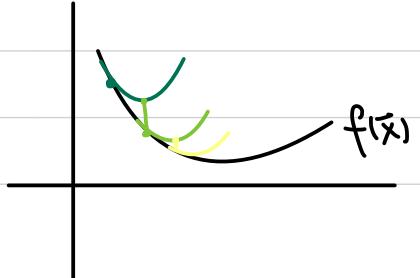
SOCP ②  $2(\vec{a}_3 - \vec{a}_2)^T \vec{x} + 2C^2(\Delta_3 - \Delta_2)\delta = C^2(\Delta_2^2 - \Delta_3^2) + \|\vec{a}_3\|_2^2 - \|\vec{a}_2\|_2^2$

③  $\|\vec{x} - \vec{a}_3\|_2 \leq C\Delta_3 + C\delta$  relaxed. linear obj.  $\Rightarrow \vec{x}^*$  at boundary.

Newtons Method.

$\min_{\vec{x}} f(\vec{x})$  w.r.t.  $\vec{x}_0, \vec{x}_1, \dots$  converging to  $\vec{x}^*$ .

$f(\vec{x} + \vec{v}) = f(\vec{x}) + \nabla f(\vec{x})^T \vec{v} + \frac{1}{2} \vec{v}^T \nabla^2 f(\vec{x}) \vec{v} + \dots$  quadratic approximation.



$\vec{x}_0$  initial point. Hessian is PD (invertible)  $\Rightarrow \vec{v} = -(\nabla^2 f(\vec{x}))^{-1} \nabla f(\vec{x})$ .

Newton Step:  $\vec{x}_{k+1} = \vec{x}_k - (\nabla^2 f(\vec{x}_k))^{-1} \nabla f(\vec{x}_k)$ . No  $\eta$ . go to the actual minimum point instead of that direction if Hessian is PSD  $\Rightarrow$  Quasi-Newton Method.

Newton's Method VS Gradient Descent.

Pros: converge faster. if  $f(\vec{x})$  is quadratic, one step vs. forever.

Cons: need to compute  $\nabla^2$  vs  $\nabla$ .

11/08/2022 Tuesday.

The  $\ell_1$  Norm.  $\ell_1$  Norm (# nonzero elements)  $\Rightarrow$  Not convex  $\Rightarrow$  Relax to  $\ell_1$  Norm.

① min  $\ell_1$  norm:

$\min_{\vec{x}} \|\vec{x}\|_1$  s.t.  $A\vec{x} = \vec{b}$ . Convex objective, linear constraint. Not differentiable everywhere.  $\rightarrow$  full row rank.

$$\|\vec{x}\|_1 = \sum_{i=1}^n |x_i| \quad x_i = x_i^+ - x_i^- \quad |x_i| = x_i^+ + x_i^- \quad x_i^+ = \begin{cases} x_i & \text{if } x_i > 0 \\ 0 & \text{otherwise} \end{cases} \quad x_i^- = \begin{cases} -x_i & \text{if } x_i < 0 \\ 0 & \text{otherwise} \end{cases} \quad x_i^+, x_i^- \geq 0.$$

$$\min_{x^+, x^-} \sum_{i=1}^n x_i^+ + \sum_{i=1}^n x_i^- \quad \text{s.t. } A(\vec{x}^+ - \vec{x}^-) = \vec{b} \quad x_i^+ \geq 0, x_i^- \geq 0.$$

Claim: This new program will always choose only one of  $x_i^+$ ,  $x_i^-$  to be nonzero.

Proof: Suppose  $x_i^+ > 0$  AND  $x_i^- > 0$ . WLOG,  $x_i^+ > x_i^-$

Consider  $x_i^{(\text{new})} = x_i^+ - x_i^-$  and  $x_i^{(\text{old})} = 0$ .

$x_i^+, x_i^-$  couldn't be the optimal.

Still under the constraints But strictly decrease obj.  $\Rightarrow$  contradicts,

② least squares.

$$\min_{\vec{x}} \|\vec{A}\vec{x} - \vec{b}\|_1 \iff \min_{\vec{e}} \|\vec{e}\|_1 \text{ s.t. } \vec{A}\vec{x} - \vec{e} = \vec{b} \quad \text{①} \rightarrow \text{LP.}$$

$$\min_{\vec{x}} \sum_{i=1}^m (\vec{x} - \vec{b}_i)^2 \quad \vec{x}^* = \frac{1}{m} \sum_{i=1}^m \vec{b}_i \quad \text{undefined if } m \text{ is even.}$$

$\ell_1$  Norm Parallel.  $\min_{\vec{x}} \sum_{i=1}^m |\vec{x} - \vec{b}_i| = \text{Median}$  (robust than mean) Scalar only.

Scalar case:  $x \in \mathbb{R}$ ,  $b_i \in \mathbb{R}$   $1 \leq i \leq m$ .  $\min_{x} \sum_{i=1}^m |x - b_i|$

$$|x - b_i| = \begin{cases} x - b_i & \text{if } x > b_i \\ b_i - x & \text{if } x \leq b_i \end{cases} \quad \frac{d}{dx} |x - b_i| = \begin{cases} 1 & \text{if } x > b_i \\ -1 & \text{if } x < b_i \\ \text{DNE if } x = b_i \text{ doesn't exist.} \end{cases}$$

If Not differentiable everywhere, think of critical points ( $\frac{d}{dx} = 0$ ,  $\frac{d}{dx}$  DNE, on boundary).

$$\text{Consider } \frac{d}{dx} \left( \sum_{i=1}^m |x - b_i| \right) = \frac{-7, -5, -3, -1, 1, 3, 5, 7}{b_1, b_2, b_3, b_4, b_5, b_6, b_7} \quad \text{DNE otherwise.}$$

$$\sum_{i=1}^m |x - b_i| = |x - b_1| + \sum_{i=2}^m |x - b_i|. \quad \forall x \in (b_3, b_6], |x - b_3| + |x - b_6| = \text{constant.}$$

### ③ LASSO

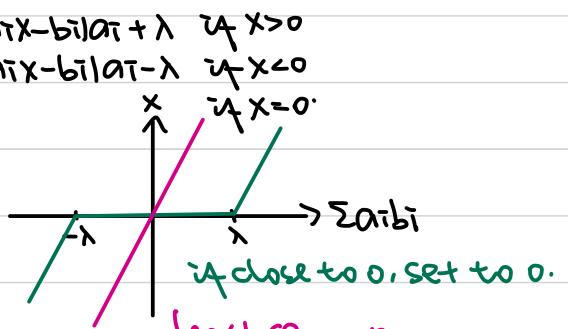
$$\min_{\vec{x}} \frac{1}{2} \|A\vec{x} - b\|_2^2 + \lambda \|\vec{x}\|_1$$

Simple scalar case:  $\min_{x} \frac{1}{2} \sum_{i=1}^n (a_i x - b_i)^2 + \lambda |x|$ .  $\frac{d}{dx} f(x) = \begin{cases} \sum a_i x - b_i + \lambda & \text{if } x > 0 \\ \sum a_i x - b_i - \lambda & \text{if } x < 0 \\ 0 & \text{if } x = 0 \end{cases}$

$$\text{if } x > 0, x_i = \frac{\sum a_i b_i - \lambda}{\sum a_i^2} \text{ if } \sum a_i b_i > \lambda.$$

$$\text{if } x < 0, x_i = \frac{\sum a_i b_i + \lambda}{\sum a_i^2} \text{ if } \sum a_i b_i < -\lambda.$$

Soft thresholding. 1 sparsity.



least square

$$\min_{\vec{x}} \frac{1}{2} \sum_{i=1}^n (a_i x - b_i)^2 \quad \vec{x}^* = \frac{\sum a_i b_i}{\sum a_i^2}$$

11/10/2022 Thursday.

Coordinate Descent. Work well for LASSO.

(sequential) Version.  $f(\vec{x}) = g(\vec{x}) + \sum_{i=1}^n h_i(|x_i|)$   $g$ : convex, differentiable.  $h_i$ : convex.

initial guess:  $\vec{x}^{(0)} = [x_1^{(0)}, x_2^{(0)}, \dots, x_n^{(0)}]$

$x_1^{(1)} = \underset{x_1}{\operatorname{argmin}} f(x_1, x_2^{(0)}, x_3^{(0)}, \dots, x_n^{(0)})$ .

$x_1^{(k)} = \underset{y}{\operatorname{argmin}} f(y, x_2^{(k-1)}, \dots, x_n^{(k-1)})$  fixed.

$x_2^{(k)} = \underset{y}{\operatorname{argmin}} f(x_1^{(k)}, y, x_3^{(k-1)}, \dots, x_n^{(k-1)})$ .

$\vdots$   
 $x_n^{(k)} = \underset{y}{\operatorname{argmin}} f(x_1^{(k)}, x_2^{(k)}, \dots, y)$

Equality Constrained Newton. equality constrained QP  $\rightarrow$  unconstrained QP.

$\min_{\vec{x}} \frac{1}{2} \vec{x}^T H \vec{x} + \vec{c}^T \vec{x} + d$  s.t.  $A\vec{x} = \vec{b}$ .  $H$  is PD, strong duality.

$L(\vec{x}, \vec{v}) = \frac{1}{2} \vec{x}^T H \vec{x} + \vec{c}^T \vec{x} + d + \vec{v}^T (A\vec{x} - \vec{b})$ .

KKT:  $A\vec{x}^* = \vec{b}$ .  $H\vec{x}^* + \vec{c} + A^T \vec{v}^* = 0 \Rightarrow \begin{bmatrix} H & A^T \\ A & 0 \end{bmatrix} \begin{bmatrix} \vec{x}^* \\ \vec{v}^* \end{bmatrix} = \begin{bmatrix} -\vec{c} \\ \vec{b} \end{bmatrix}$  invertible/non-singular  $\Rightarrow$  unique soln.

$\min_{\vec{x}} f(\vec{x}_0 + \vec{v}) = f(\vec{x}_0) + \nabla f(\vec{x}_0)^T \vec{v} + \frac{1}{2} \vec{v}^T D^2 f(\vec{x}_0) \vec{v} + \dots$  s.t.  $A(\vec{x}_0 + \vec{v}) = \vec{b}$ .  $A\vec{v} = 0$ .

Soln:  $\begin{bmatrix} D^2 f(\vec{x}_0) & A^T \\ A & 0 \end{bmatrix} \begin{bmatrix} \vec{v} \\ \vec{v} \end{bmatrix} = \begin{bmatrix} -\nabla f(\vec{x}_0)^T \\ 0 \end{bmatrix}$

11/15/2022 Tuesday.

LQG (Linear Quadratic Gaussian).

Application: Control (Linear Quadratic Regulator).

$$x(t+1) = f(x(t), u(t)).$$

Vertical Rocket. Goal: Maximize height by time T.

$x_1(t)$  = height.  $x_2(t)$  = vertical speed.  $x_3(t)$  = weight of the rocket.  $(x_1(0), x_2(0), x_3(0)) = (0, 0, M)$ .

Force:  $\uparrow$  thrust:  $G \cdot x_3$   $\downarrow$  inertia:  $x_3 \cdot \ddot{x}_1 = x_3 \cdot x_2^2$   $\downarrow$  drag:  $C_D \cdot P(x_1) \cdot x_2^2$   $\downarrow$  gravity:  $g x_3$

$$\dot{x}_1(t) = x_2(t) \quad x_3(t) \quad \dot{x}_2(t) = -C_D P(x_1) x_2^2 - g x_3 + G u(t) \quad \dot{x}_3(t) = -\frac{C_D}{x_3(t)} P(x_1) x_2^2 - g + \frac{G u(t)}{x_3(t)}$$

$\max_{\sum_{t=0}^{N-1}} x_1(t)$  s.t.  $\dot{x}(t) = f(\vec{x}(t), \vec{u}(t))$   $\forall t \in \mathbb{N}$ . linearize / approximate to  $\vec{x}(t+1) = A \vec{x}(t) + B \vec{u}(t)$ .

Cost:  $\sum_{t=0}^{N-1} \frac{1}{2} (\vec{x}_t^T Q \vec{x}_t + \vec{u}_t^T R \vec{u}_t) + \frac{1}{2} \vec{x}_N^T Q_f \vec{x}_N$  ( $Q, R$ : PD matrices. Time invariant).

minimize  $\vec{x}_1, \vec{x}_2, \dots, \vec{x}_N, \vec{u}_1, \vec{u}_2, \dots, \vec{u}_N \sum_{t=0}^{N-1} \frac{1}{2} (\vec{x}_t^T Q \vec{x}_t + \vec{u}_t^T R \vec{u}_t) + \frac{1}{2} \vec{x}_N^T Q_f \vec{x}_N$  s.t.  $\vec{x}_{t+1} = A \vec{x}_t + B \vec{u}_t \quad t = 0, \dots, N-1$ .

Solve by Riccati Equation,

$$L(\vec{x}_0, \dots, \vec{x}_N, \vec{u}_0, \dots, \vec{u}_N, \vec{\lambda}_1, \dots, \vec{\lambda}_N) = \sum_{t=0}^{N-1} \frac{1}{2} (\vec{x}_t^T Q \vec{x}_t + \vec{u}_t^T R \vec{u}_t) + \frac{1}{2} \vec{x}_N^T Q_f \vec{x}_N + \sum_{t=0}^{N-1} \vec{\lambda}_{t+1}^T (A \vec{x}_t + B \vec{u}_t - \vec{x}_{t+1}).$$

KKT condition: ①  $\nabla_{\vec{u}_t} L = R \vec{u}_t + B^T \vec{\lambda}_{t+1} = 0 \quad t = 0, 1, \dots, N-1 \Rightarrow$  ⑥  $\vec{u}_t = -R^{-1} B^T \vec{\lambda}_{t+1}$ .

$$\textcircled{2} \quad \nabla_{\vec{x}_t} L = Q \vec{x}_t + A^T \vec{\lambda}_{t+1} - \vec{\lambda}_t = 0 \quad t = 1, \dots, N-1 \Rightarrow \textcircled{4} \quad \vec{\lambda}_t = Q \vec{x}_t + A^T \vec{\lambda}_{t+1}$$

$$\textcircled{3} \quad \nabla_{\vec{x}_N} L = Q_f \vec{x}_N - \vec{\lambda}_N = 0 \Rightarrow \textcircled{5} \quad \vec{\lambda}_N = Q_f \vec{x}_N. \quad \begin{array}{l} \text{Dynamics of the adjoint sys.} \\ \vec{\lambda} = \text{co-state.} \end{array}$$

Approach: Backward Induction. Goal: find optimal  $\vec{u}_t$ .

Induction Hypothesis:  $\vec{\lambda}_{t+1} = P_{t+1} \cdot \vec{x}_{t+1}$ .

Base Case:  $t=N$ .  $\vec{\lambda}_N = Q_f \cdot \vec{x}_N \Rightarrow P_N = Q_f$ .

Inductive Steps:  $\vec{\lambda}_{t+1} = P_{t+1} \cdot \vec{x}_{t+1} = P_{t+1} \cdot (A \vec{x}_t + B(-R^{-1} B^T \vec{\lambda}_{t+1})) = (I + P_{t+1} B R^{-1} B^T)^{-1} \cdot P_{t+1} A \vec{x}_t$ .

$$\vec{\lambda}_t = A^T \vec{\lambda}_{t+1} + Q \vec{x}_t = \frac{A^T (I + P_{t+1} B R^{-1} B^T)^{-1} \cdot P_{t+1} A + Q}{P_t} \vec{x}_t$$

11/17/2022 Thursday.

do not depend on every single datapoint,  
only ones at the boundary margin  $\Rightarrow$  support vectors.

Support Vector Machine. Classification.

$$\max_{\vec{w}, b, m, m} \text{s.t. } y_i(\vec{w}^T \vec{x}_i + b) > 0 \quad \forall i. \quad \frac{|\vec{w}^T \vec{x}_i + b|}{\|\vec{w}\|_2} \geq m \cdot \forall i \iff \frac{|y_i(\vec{w}^T \vec{x}_i + b)|}{\|\vec{w}\|_2} \geq m.$$

the separating hyperplane w/ greatest margin.

distance of the hyperplane to the closest datapoint.

$$m = \frac{1}{\|\vec{w}\|_2}$$

$$\Leftrightarrow \max_{\vec{w}, b} (\|\vec{w}\|_2)^{-1} \text{ s.t. } y_i(\vec{w}^T \vec{x}_i + b) \geq 1. \Leftrightarrow \min_{\vec{w}, b} \frac{1}{2} \|\vec{w}\|_2^2 \text{ s.t. } y_i(\vec{w}^T \vec{x}_i + b) \geq 1.$$

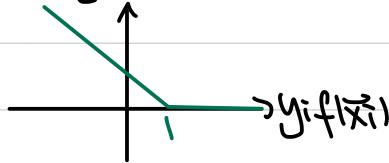
Hard-Margin SVM

Soft-Margin SVM:

$$\min_{\vec{w}, b, \xi_i} \frac{1}{2} \|\vec{w}\|_2^2 + C \sum_{i=1}^n \xi_i \text{ s.t. } y_i(\vec{w}^T \vec{x}_i + b) \geq 1 - \xi_i \quad \xi_i \geq 0.$$

↓  
larger C, harder SVM.

Hinge-loss perspective of Soft-Margin SVM.



$$\text{Hinge}(y_i, \vec{w}^T \vec{x}_i + b) = \max(1 - y_i(\vec{w}^T \vec{x}_i + b), 0).$$

$$\text{Hinge-loss perspective: } \min_{\vec{w}, b} \frac{1}{2} \sum_i (\text{Hinge}(y_i, \vec{w}^T \vec{x}_i + b) + \lambda \|\vec{w}\|_2^2)$$

$$\begin{aligned} \text{Soft-Margin SVM: } & \min_{\vec{w}, b, \xi_i} \frac{1}{2} \|\vec{w}\|_2^2 + C \sum_{i=1}^n \xi_i \text{ s.t. } \xi_i \geq \max(1 - y_i(\vec{w}^T \vec{x}_i + b), 0) \\ & \Leftrightarrow \min_{\vec{w}, b} \frac{1}{2} \|\vec{w}\|_2^2 + C \sum_{i=1}^n \text{Hinge}(y_i, \vec{w}^T \vec{x}_i + b) \quad C = n\lambda. \end{aligned}$$

11/22/2022 Tuesday.

Dual Perspective of Soft-Margin SVM.

$$\begin{aligned} L(\vec{w}, b, \vec{\xi}, \vec{\alpha}, \vec{\beta}) &= \frac{1}{2} \|\vec{w}\|_2^2 + C \sum_{i=1}^n \xi_i + \sum_{i=1}^n \alpha_i ((1 - \xi_i) - y_i(\vec{w}^T \vec{x}_i + b)) + \sum_{i=1}^n \beta_i (-\xi_i) \\ &= \frac{1}{2} \|\vec{w}\|_2^2 - \sum_{i=1}^n \alpha_i y_i (\vec{w}^T \vec{x}_i + b) + \sum_{i=1}^n \alpha_i + \sum_{i=1}^n (C - \alpha_i - \beta_i) \cdot \xi_i \end{aligned}$$

$$\text{First-order KKT: } \nabla_{\vec{w}} L = \vec{w} - \sum_{i=1}^n \alpha_i y_i \vec{x}_i = 0 \Rightarrow \vec{w} = \sum_{i=1}^n \alpha_i y_i \vec{x}_i$$

$$\frac{\partial L}{\partial b} = -\sum_{i=1}^n \alpha_i y_i = 0$$

$$\frac{\partial L}{\partial \xi_i} = (C - \alpha_i - \beta_i) = 0.$$

$$\text{Complementary Slackness: } \alpha_i ((1 - \xi_i) - y_i(\vec{w}^T \vec{x}_i + b)) = 0 \quad \beta_i \xi_i = 0 \quad \forall i.$$

Case ①:  $\alpha_i \neq 0, \alpha_i \neq C, 0 < \alpha_i < C$

$$(1 - \xi_i) - y_i(\vec{w}^T \vec{x}_i + b) = 0 \quad \alpha_i \neq C \Rightarrow \beta_i \neq 0 \Rightarrow \xi_i = 0.$$

$y_i(\vec{w}^T \vec{x}_i + b) = 1 \Rightarrow$  All points are exactly on the margin.

Case ②:  $\alpha_i \neq 0, \alpha_i = C$

$$\beta_i = 0, \xi_i \text{ need not be zero.} \Rightarrow y_i(\vec{w}^T \vec{x}_i + b) = 1 - \xi_i \leq 1$$

Case ③:  $\alpha_i = 0$

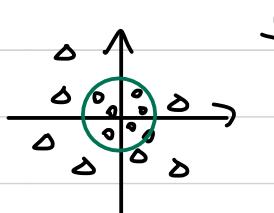
$$\beta_i = C \neq 0 \Rightarrow \xi_i = 0 \Rightarrow \text{No margin violation.}$$

$$\text{Dual: } L = -\frac{1}{2} \|\vec{w}\|_2^2 + \sum_{i=1}^n \alpha_i = -\frac{1}{2} \left( \sum_{i=1}^n \alpha_i y_i \vec{x}_i^T \right) \left( \sum_{i=1}^n \alpha_i y_i \vec{x}_i \right) + \sum_{i=1}^n \alpha_i = -\frac{1}{2} \vec{x}^T \underbrace{\text{diag}(\vec{y})}_{Q} \vec{x} + \sum_{i=1}^n \alpha_i.$$

↑  
by first-order KKT

$$d^* = \max_{\vec{z}, \vec{\beta}} -\frac{1}{2} \vec{z}^T Q \vec{z} + \sum_i \vec{\beta}_i \quad \text{s.t. } \sum_i \alpha_i y_i = 0, C - \alpha_i - \beta_i = 0, \alpha_i \geq 0, \beta_i \geq 0 \quad \forall i \\ 0 \leq \alpha_i \leq C.$$

kernel SVM:



$$f(\vec{x}) = \vec{w} \phi(\vec{x}) + b. \quad \phi: (\vec{x}_1, \vec{x}_2) \mapsto (\vec{x}_1, \vec{x}_2, \vec{x}_1^2 + \vec{x}_2^2).$$

e.g. poly, RBF (radial basis function).

10/11/2022 Tuesday.

Stochastic Gradient Descent.

$$\vec{x}_{k+1} = \vec{x}_k - \eta \nabla f_i(\vec{x}_k). \quad \nabla f(\vec{x}) = \frac{1}{m} \sum_{i=1}^m \nabla f_i(\vec{x})$$

$$\text{Example. } f(\vec{x}) = \frac{1}{m} \sum_{i=1}^m \frac{1}{2} \|\vec{x} - \vec{p}_i\|^2$$

$$\text{SGD. stepsize: } \eta_t = \frac{1}{t}, \vec{x}_0 = 0 \quad \nabla f_i(\vec{x}) = \frac{1}{2} \cdot 2(\vec{x} - \vec{p}_i) = \vec{x} - \vec{p}_i.$$

$$\vec{x}_1 = \vec{x}_0 - (\vec{x}_0 - \vec{p}_1) = \vec{p}_1 \quad \vec{x}_2 = \frac{\vec{p}_1 + \vec{p}_2}{2}$$

Projected Gradient Descent.

$$\vec{x}_{k+1} = \Pi_X(\vec{x}_k - \eta \cdot \nabla f(\vec{x}_k)) \quad \Pi_X(\vec{y}) = \underset{\vec{z} \in X}{\operatorname{argmin}} \|\vec{y} - \vec{z}\|_2^2 \quad \text{Project back into the set } X.$$

Conditional Gradient Descent.

$$\vec{y}_k = \underset{\vec{y} \in X}{\operatorname{argmin}} \nabla f(\vec{x}_k)^T \vec{y} \quad \vec{x}_{k+1} = (1 - \gamma_k) \vec{x}_k + \gamma_k \vec{y}_k \quad \xrightarrow{\text{predetermined sequence.}}$$