

Outline of this course:

1) Linear Algebra

operations, determinants, solution of linear system, block metrics, eigen value, eigen vector, SVD, other decomposition. Vector space, orthogonality principle.

2) Probability theory

definition of probability space, random variables, expectations, conditional prob., conditional expectations, Decision Theory, estimation.

3) Optimization

Calculus of variation

4) Machine learning algorithm

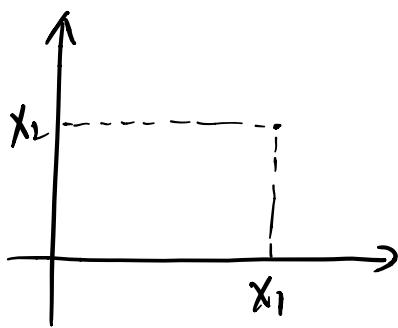
Schedule. Mon. 12:00 - 13:20 13:40 - 15:00

3 5 7 9 11      HW — 1 week

HW HW MT HW HW      MT — 2 days  
↑  
no course

[www.cs.rutgers.edu/~gm463/courses/CS596](http://www.cs.rutgers.edu/~gm463/courses/CS596)

# LINEAR ALGEBRA.



$$\begin{pmatrix} x_1 \\ x_2 \end{pmatrix} \quad \begin{pmatrix} x_2 \\ x_1 \end{pmatrix}$$

order is important.

$$\begin{bmatrix} x_1 \\ x_2 \\ \vdots \\ x_n \end{bmatrix} \text{ Vector } \mathbb{R}^k$$

sequence of real numbers — vector space of infinite dimensions

$\neq X=Y$

$$X + Y = \begin{bmatrix} x_1 + y_1 \\ x_2 + y_2 \\ \vdots \\ x_n + y_n \end{bmatrix}$$

$$\begin{bmatrix} x_1 \\ x_2 \\ \vdots \\ x_n \end{bmatrix} + \begin{bmatrix} y_1 \\ y_2 \\ \vdots \\ y_n \end{bmatrix}$$

$$a * X = \begin{bmatrix} ax_1 \\ ax_2 \\ \vdots \\ ax_n \end{bmatrix}$$

$$ax = b$$

- 1)  $a \neq 0$      $x = \frac{b}{a} = a^{-1}b$
- 2)  $a=0, b=0$     any real #
- 3)  $a=0, b \neq 0$     no solutions

$$a_{11}x_1 + a_{12}x_2 + \dots + a_{1k}x_k = b_1$$

$$a_{21}x_1 + a_{22}x_2 + \dots + a_{2k}x_k = b_2$$

⋮

$$a_{k1}x_1 + a_{k2}x_2 + \dots + a_{kk}x_k = b_k$$

$$A = \begin{bmatrix} a_{11} & a_{12} & \dots & a_{1k} \\ a_{21} & a_{22} & \dots & a_{2k} \\ \vdots & & & \\ a_{k1} & a_{k2} & \dots & a_{kk} \end{bmatrix} \quad X = \begin{bmatrix} x_1 \\ x_2 \\ \vdots \\ x_k \end{bmatrix}, \quad b = \begin{bmatrix} b_1 \\ b_2 \\ \vdots \\ b_k \end{bmatrix}$$

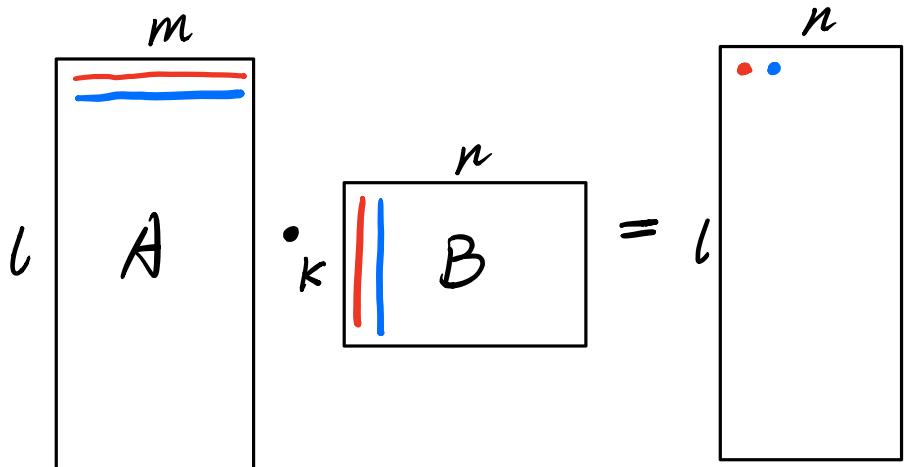
← matics

$$A \cdot X = b \quad \text{Inner product.}$$

$$A_{l \times m} \quad B_{k \times n}$$

$$\underset{\substack{k=1 \\ m=n}}{*} A + B = C \quad \underset{m=k}{*} A \cdot B$$

$$A_{ij} + B_{ij} = C_{ij}$$



$$A = \begin{bmatrix} a_{11} & a_{12} & \dots & a_{1k} \\ a_{21} & a_{22} & \dots & a_{2k} \\ \vdots & \vdots & & \\ a_{k1} & a_{k2} & \dots & a_{kk} \end{bmatrix}$$

$$A^T = \begin{bmatrix} a_{11} & a_{21} & \dots & a_{k1} \\ a_{12} & a_{22} & \dots & a_{k2} \\ \vdots & \vdots & & \vdots \\ a_{k1} & a_{k2} & \dots & a_{kk} \end{bmatrix}$$

$$(A_1 \cdot A_2)^T = A_2^T \cdot A_1^T$$

$$(A_1 \cdot A_2 \cdot \dots \cdot A_k)^T = A_k^T \cdot \dots \cdot A_2^T \cdot A_1^T$$

$$A = \begin{array}{|c|c|c|} \hline & A_{11} & A_{12} & A_{13} \\ \hline & A_{21} & A_{22} & A_{23} \\ \hline & A_{31} & A_{32} & A_{33} \\ \hline & A_{41} & A_{42} & A_{43} \\ \hline \end{array}$$

$$B = \begin{array}{|c|c|} \hline & B_{11} & B_{12} \\ \hline & B_{21} & B_{22} \\ \hline & B_{31} & B_{32} \\ \hline \end{array}$$

A and B have the same dimensions.

BUT. sub-matrices in A and B are different.

To make any operations between A and B,  
both A and B need to be sub-divided into same  
sub-matrices.

$$\begin{array}{|c|c|} \hline & A_{11} & A_{12} \\ \hline & A_{21} & A_{22} \\ \hline \end{array} + \begin{array}{|c|c|} \hline & B_{11} & B_{12} \\ \hline & B_{21} & B_{22} \\ \hline \end{array} = \begin{array}{|c|c|} \hline & A_{11}+B_{11} & A_{12}+B_{12} \\ \hline & A_{21}+B_{21} & A_{22}+B_{22} \\ \hline \end{array}$$

$$A \cdot B = C$$

*Outer product*

$$\begin{pmatrix} \mathbf{a}_1^T \\ \mathbf{a}_2^T \\ \vdots \\ \mathbf{a}_k^T \end{pmatrix} \cdot (\mathbf{b}_1, \mathbf{b}_2, \dots, \mathbf{b}_n) = \begin{bmatrix} \mathbf{a}_1^T \mathbf{b}_1 & \mathbf{a}_1^T \mathbf{b}_2 & \cdots & \mathbf{a}_1^T \mathbf{b}_n \\ \mathbf{a}_2^T \mathbf{b}_1 & \mathbf{a}_2^T \mathbf{b}_2 & \cdots & \mathbf{a}_2^T \mathbf{b}_n \\ \vdots & \vdots & & \vdots \\ \mathbf{a}_k^T \mathbf{b}_1 & \mathbf{a}_k^T \mathbf{b}_2 & \cdots & \mathbf{a}_k^T \mathbf{b}_n \end{bmatrix}$$

$$A_{m \times n} \cdot B_{k \times l} = C_{m \times l}$$

$n=k$

$$[\tilde{\mathbf{a}}_1, \tilde{\mathbf{a}}_2, \dots, \tilde{\mathbf{a}}_n] \cdot \begin{bmatrix} \tilde{\mathbf{b}}_1 \\ \tilde{\mathbf{b}}_2 \\ \vdots \\ \tilde{\mathbf{b}}_k \end{bmatrix} = \tilde{\mathbf{a}}_1 \tilde{\mathbf{b}}_1^T + \dots + \tilde{\mathbf{a}}_n \tilde{\mathbf{b}}_k^T$$

$$A \cdot X = b$$

$$A = \begin{bmatrix} \mathbf{a}_1 & \cdots & \mathbf{a}_k \end{bmatrix} \quad X = \begin{bmatrix} x_1 \\ x_2 \\ \vdots \\ x_k \end{bmatrix} \quad b = \begin{bmatrix} b_1 \\ b_2 \\ \vdots \\ b_k \end{bmatrix}$$

$$(A^T A)X = A^{-1}b$$

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

$$A = \begin{bmatrix} a_{11} & a_{12} & \cdots & a_{1k} \\ \vdots & & & \\ a_{11} & a_{12} & \cancel{a_{1j}} & a_{1k} \\ \vdots & & & \\ a_{kk} & a_{kk} & \cdots & a_{kk} \end{bmatrix}$$

$$A - \text{col}(a_{ij}) - \text{row}(a_{ij}) = M_{ij}$$

By this way, we can reduce the metric from  $x$  to  $x-1$ .  
determine  $|A|$

$$|A| = a_{11}a_{12} + \cdots + a_{1k}a_{1k}$$

$$a_{ij} = (-1)^{i+j} \det(M_{ij})$$

co-efficient.

$$A \begin{bmatrix} A_{11} & A_{12} & \cdots & A_{1k} \\ A_{21} & A_{22} & \cdots & A_{2k} \\ \vdots & \vdots & & \vdots \\ A_{k1} & A_{k2} & \cdots & A_{kk} \end{bmatrix} \xrightarrow{A^T} \begin{bmatrix} A_{11} & A_{21} & \cdots & A_{k1} \\ A_{12} & A_{22} & \cdots & A_{k2} \\ \vdots & \vdots & & \vdots \\ A_{1k} & A_{2k} & \cdots & A_{kk} \end{bmatrix} b = \begin{bmatrix} b_1 \\ b_2 \\ \vdots \\ b_k \end{bmatrix}$$

$$A^{-1} = \frac{A^T}{\det(A)} \quad \text{for } x = A^{-1}b$$

$$x = \frac{A^T}{\det(A)} \cdot b$$

photo 3 cipher

Gauss elimination.

photo 4

$$Ax = b \quad \hat{x}$$

$$\hat{A}\hat{x} = b$$

having a solution does not mean use it.

$$\begin{array}{c}
 \beta_j \alpha_{j1} \cdots \beta_j \alpha_{jk} \\
 \hline
 \beta_i \alpha_{i1} \cdots \beta_i \alpha_{ik}
 \end{array}$$

$\beta_j \alpha_{ji}$  add to  $\beta_i \alpha_{ij}$   
 $= \beta_i \det(A)$

$$\det(A \cdot B) = \det(A) \cdot \det(B)$$

$$\det \begin{vmatrix} A_{11} & 0 & 0 & 0 \\ - & A_{22} & 0 & 0 \\ - & - & A_{33} & 0 \\ - & - & - & A_{44} \end{vmatrix} = \det(A_{11}) \cdots \det(A_{44})$$

Shur's Formula

$$A = \begin{vmatrix} A_{11} & A_{12} \\ A_{21} & A_{22} \end{vmatrix}$$

$$\begin{aligned} \det(A) &= \det(A_{11}) \cdot \det(A_{22} - A_{21} A_{11}^{-1} A_{12}) \\ &= \det(A_{22}) \cdot \det(A_{11} - A_{12} A_{22}^{-1} A_{21}) \end{aligned}$$

$$A = \boxed{\phantom{000}}$$

$$A^T = \boxed{\phantom{000}}$$

$$A \cdot A^T = I$$



rotation matrices

orthonormal matrices  
?

$$\det(A^T) = \det(A)$$

$$\det(A^T \cdot A) = \det(I) = 1$$

$$\det(A^T) = \det(A) = 1$$

*symmetric*

$$\begin{matrix} a_1^T a_1 & a_1^T a_2 & \dots & a_1^T a_k \\ a_2^T a_1 & a_2^T a_2 & \dots & a_2^T a_k \\ \vdots & \vdots & \ddots & \vdots \\ a_k^T a_1 & a_k^T a_2 & \dots & a_k^T a_k \end{matrix}$$

$$b^T a = a^T b$$

Degree of freedom.  $K^2 - \frac{k(k+1)}{2}$  (only for rotation matrices).

$$A x = b$$

$$x = A^{-1} b$$

Matrix Inversion Lemma

$$A^{-1} = A^{-1} B (D^{-1} + C^T A^{-1} B)^{-1} C^T A^{-1}$$

$$2k^2m + 2km^2 + 2m^3 \quad \begin{matrix} K \\ A \end{matrix} + \left( \begin{matrix} m \\ K \\ B \end{matrix} \right) \left( \begin{matrix} m \\ m \\ D \end{matrix} \right) \left( \begin{matrix} m \\ m \\ C^T \end{matrix} \right)$$

$$A = \begin{array}{|c|c|c|} \hline A_{11} & 0 & 0 \\ \hline 0 & A_{22} & 0 \\ \hline 0 & 0 & A_{33} \\ \hline \end{array} \quad A' = \begin{array}{|c|c|c|} \hline A_{11}' & 0 & 0 \\ \hline 0 & A_{22}' & 0 \\ \hline 0 & 0 & A_{33}' \\ \hline \end{array}$$

Shur's Inversim Formula.

$$A = \begin{array}{|c|c|} \hline A_{11} & A_{12} \\ \hline A_{21} & A_{22} \\ \hline \end{array} = \begin{array}{|c|c|} \hline A_{11}^{-1} & \\ \hline & \\ \hline \end{array} + \begin{bmatrix} -E \\ F \\ I \end{bmatrix} \Delta^{-1} \begin{bmatrix} E & F & I \end{bmatrix}$$

$$\Delta = A_{22} - A_{21} A_{11}^{-1} A_{12} \quad E = A_{11}^{-1} A_{12} \quad F = A_{21} A_{11}^{-1}$$

$K^2m \text{ vs } K^3$

$$Ax = 0$$

$$\begin{bmatrix} a_1 & a_2 & \dots & a_k \end{bmatrix}$$

$$\begin{bmatrix} x_1 \\ x_2 \\ \vdots \\ x_k \end{bmatrix} = 0$$

$$x_1 g_1 + x_2 g_2 + \dots + x_k g_k = 0$$

$r_C$  = longest subset of column  
that are linearly independent.

$r_C$  = Column rank     $r_r$  = row rank

$r_C = r_r = \text{matrix rank}$

$$X \rightarrow Y$$

$$X = \begin{bmatrix} x_1 \\ x_2 \\ \vdots \\ x_k \end{bmatrix} \quad Y = \begin{bmatrix} y_1 \\ y_2 \\ \vdots \\ y_m \end{bmatrix}$$

$X_1 \rightarrow Y_1 \quad X_1 + X_2 \rightarrow Y_1 + Y_2$   
 $X_2 \rightarrow Y_2 \quad \alpha X_1 \rightarrow \alpha Y_1$

$Y = AX$  linear mapping  
 $(m \times k)$   
can represent by matrix

$m = k$  vector define direction in space

$X = \alpha_1 X_1 + \dots + \alpha_k X_k$  decomposing  $X$  into  $k$  basis.

$$Y = \beta_1 X_1 + \dots + \beta_k X_k$$

$$\beta_1 X_1 + \dots + \beta_k X_k = A(\alpha_1 X_1 + \dots + \alpha_k X_k)$$

$$= \alpha_1 \lambda_1 X_1 + \dots + \alpha_k \lambda_k X_k$$

$$\begin{bmatrix} \beta_1 \\ \beta_2 \\ \vdots \\ \beta_k \end{bmatrix} = \begin{bmatrix} \lambda_1 & & & \\ & \lambda_2 & & \\ & & \ddots & \\ 0 & & \ddots & \lambda_k \end{bmatrix} \begin{bmatrix} \alpha_1 \\ \alpha_2 \\ \vdots \\ \alpha_k \end{bmatrix}$$

$$AX = \lambda X$$

$$BX = 0 \quad \det(B) = 0$$

$$(A - \lambda I)X = 0$$

$$\det(\lambda I - A) = 0$$

$$\det \begin{pmatrix} \lambda - a_{11} & -a_{12} & \cdots & -a_{1k} \\ -a_{21} & \lambda - a_{22} & \cdots & -a_{2k} \\ \vdots & \ddots & \ddots & \vdots \\ 0 & 0 & \ddots & 0 \\ & & & \lambda - a_{kk} \end{pmatrix} = C_k \lambda^k + C_{k-1} \lambda^{k-1} + \cdots + C_1 \lambda^1 + C_0$$

CHARACTERISTIC  
POLYNOMIAL

$$\det(\alpha B) = \alpha^k \det(B)$$

$$\frac{1}{\lambda^k} \det(\lambda I - A) = C_k + \frac{C_{k-1}}{\lambda} + \cdots + \frac{C_0}{\lambda^k}$$

$$\det(I - \frac{A}{\lambda}) = C_k + \frac{C_{k-1}}{\lambda} + \cdots + \frac{C_0}{\lambda^k} \quad C_k = \det(I) = 1$$

$$C_0 = \det(-A) = (-1)^k \det(A)$$

$$(\lambda - \lambda_1) \cdots (\lambda - \lambda_k) = \lambda^k - (\lambda_1 + \cdots + \lambda_k) \lambda^{k-1} + \cdots$$

$$C_{k-1} = -(\lambda_1 + \cdots + \lambda_k)$$

$B = TAT^{-1}$  similarity transformation  
 $A$  &  $B$  are similar.

$$\begin{aligned} \det(\lambda I - B) &= \det(\lambda I - TAT^{-1}) \\ &= \det(\lambda TT^{-1} - TAT^{-1}) \\ \det(T) \det(T^{-1}) &= 1 = \det(T(\lambda I - A)T^{-1}) \\ T \cdot T^{-1} &= I = \det(T) \det(\lambda I - A) \det(T^{-1}) \\ &= \det(\lambda I - A) \end{aligned}$$

$$\lambda_i \lambda_i = A x_i = A T^{-1} T x_i$$

$$\frac{TAT^{-1}(Tx_i)}{B} = \lambda_i \frac{Tx_i}{Y_i}$$

$$y^T A \stackrel{?}{=} \lambda y^T$$

$$A^T y = \lambda y$$

$$\det(\lambda I - A^T) = \det((\lambda I - A)^T)$$

$$\det(\lambda I - A) =$$

$$\begin{matrix} \lambda_1, \lambda_2, \dots, \lambda_k \\ \diagdown \quad \diagup \quad \diagdown \quad \diagup \\ Y_1 \quad X_1 \quad Y_2 \quad X_2 \quad \quad Y_k \quad X_k \end{matrix}$$

$A_{K \times K}$        $\lambda_1, \lambda_2, \dots, \lambda_r$  eigenvalues

$$\begin{matrix} \downarrow & \downarrow & \downarrow \\ X_1 & X_2 & X_r \end{matrix}$$

linearly independent

$$X_1, X_2, \dots, X_m$$

$$X_{m+1} = \alpha_1 X_1 + \dots + \alpha_m X_m$$

$$X_{m+1}, X_{m+2}, \dots, X_r$$

$$A X_{m+1} = \alpha_1 A X_1 + \dots + \alpha_m A X_m$$

$$\lambda_{m+1} X_{m+1} = \alpha_1 \lambda_1 X_1 + \dots + \alpha_m \lambda_m X_m$$

$$0 = \alpha_1 (\lambda_{m+1} - \lambda_1) X_1 + \alpha_2 (\lambda_{m+1} - \lambda_2) X_2 + \dots + \alpha_m (\lambda_{m+1} - \lambda_m) X_m$$

$$\begin{pmatrix} 1 & 0 \\ 0 & 1 \end{pmatrix} \begin{pmatrix} k_1 \\ k_2 \end{pmatrix} = \begin{pmatrix} k_1 \\ k_2 \end{pmatrix} \quad \begin{pmatrix} 1 \\ 0 \end{pmatrix} \begin{pmatrix} 0 \\ 1 \end{pmatrix}$$

$$\begin{pmatrix} 1 & 1 \\ 0 & 1+\epsilon \end{pmatrix} \begin{pmatrix} k_1 \\ k_2 \end{pmatrix} = \begin{pmatrix} k_1 \\ k_2 \end{pmatrix} \quad \begin{cases} k_1 + k_2 = k_1 \\ k_2 = k_2 \end{cases} \quad \checkmark \quad \begin{pmatrix} 1 \\ 0 \end{pmatrix} \begin{pmatrix} 1+\epsilon \\ 0 \end{pmatrix}$$

if  $\epsilon = 0$ . only one eigen vector.

$x_1, x_2, \dots, x_k$  linearly independent.

$\lambda_1, \lambda_2, \dots, \lambda_k$

$$A[\bar{x}_1 \ \bar{x}_2 \ \dots \ \bar{x}_k] = [\bar{A}\bar{x}_1 \ \bar{A}\bar{x}_2 \ \dots \ \bar{A}\bar{x}_k]$$

$$[\bar{x}_1 \ \bar{x}_2 \ \dots \ \bar{x}_k]^T \begin{bmatrix} \bar{\lambda}_1 x_1 & \bar{\lambda}_2 x_2 & \dots & \bar{\lambda}_k x_k \end{bmatrix}$$

$$\begin{bmatrix} \lambda_1 & & & \\ & \lambda_2 & & \\ & & \ddots & \\ & & & \lambda_k \end{bmatrix}$$

$$A\bar{T} = \bar{T} \cdot \Lambda \Rightarrow A = \bar{T} \Lambda \bar{T}^{-1} \quad \text{diagonalization of } A$$

$$\bar{T}^{-1} A = \Lambda \bar{T}^{-1}$$

left & right  
eigenvectors  
related

$$\begin{bmatrix} x_1^T \\ x_2^T \\ \vdots \end{bmatrix}$$

$$A = \begin{bmatrix} \lambda_1 & & & \\ & \lambda_2 & & \\ & & \ddots & \\ & & & \lambda_k \end{bmatrix} \begin{bmatrix} y_1^T \\ y_2^T \\ \vdots \\ y_k^T \end{bmatrix}$$

$$Y_1^T A = \lambda_1 Y_1^T$$

$$Y_1^T x_j = 0 \quad l \neq j$$

$$Y_1^T x_i = 1$$

symmetric matrix.

$$A^T = A. \quad Ax_i = \lambda_i x_i$$

$$(x_i^*)^T A x_i = \lambda_i (\underbrace{x_i^*}_\|x_i\|^2) \underbrace{x_i}_\text{real number.}$$

$$(x_i)^T A x_i^* = \lambda_i^* \|x_i\|^2$$

$$(x_i^*)^T A^T x_i = \lambda_i^* \|x_i\|^2$$

$$\lambda_i \|x_i\|^2 = \lambda_i^* \|x_i\|^2 \quad x_i \text{ is eigenvector}$$

$$\lambda_i = \lambda_i^* \quad \lambda_i \text{ is real number.} \quad \downarrow$$

$$Ax_i = \lambda_i x_i$$

$$x_i^T A = ((x_i^T A)^T)^T = (A^T \cdot x_i)^T = \lambda_i x_i^T$$

$$A = P \cdot \Lambda \cdot P^T \quad \begin{bmatrix} x_1 & x_2 & \dots & x_k \end{bmatrix}^T \quad P$$

$A$  is symmetric

$$\min_{x \neq 0} \frac{x^T A x}{x^T x}$$

$$\frac{x^T P \cdot \Lambda \cdot P^T x}{x^T P P^T x} \quad y = P^T \cdot x$$

$$\propto \frac{y^T \cdot \Lambda \cdot y}{y^T y}$$

$$Y = \begin{bmatrix} y_1 \\ y_2 \\ \vdots \\ y_k \end{bmatrix} \quad A = \begin{bmatrix} \lambda_1 & & & \\ & \lambda_2 & & 0 \\ & & \ddots & \\ 0 & & & \lambda_k \end{bmatrix}$$

$$[y_1 \ y_2 \ \dots \ y_k] \begin{bmatrix} \lambda_1 & & & \\ & \lambda_2 & & 0 \\ & & \ddots & \\ 0 & & & \lambda_k \end{bmatrix} \begin{bmatrix} y_1 \\ y_2 \\ \vdots \\ y_k \end{bmatrix}$$

$$(\lambda_1 y_1 \ \lambda_2 y_2 \ \dots \ \lambda_k y_k) \begin{bmatrix} y_1 \\ y_2 \\ \vdots \\ y_k \end{bmatrix} = \lambda_1 y_1^2 + \lambda_2 y_2^2 + \dots + \lambda_k y_k^2$$

$$\lambda_{\min} = \min \{\lambda_1, \lambda_2, \dots, \lambda_k\}$$

$$\frac{\lambda_1 y_1^2 + \dots + \lambda_k y_k^2}{y_1^2 + \dots + y_k^2} \geq \lambda_{\min}$$

$$\lambda_1 y_1^2 \geq \lambda_{\min} y_1^2$$

$$\frac{\lambda_{\min} (y_1^2 + \dots + y_k^2)}{(y_1^2 + \dots + y_k^2)}$$

$$A = A^T$$

Positive definite.

$$x \neq 0 \quad x^T A x > 0$$

$$\lambda_{\min} \text{ of } A$$

nonnegative definite

$$x \neq 0 \quad x^T A x \geq 0$$

$$\frac{x^T A x}{x^T x} \geq \lambda_{\min} \geq 0$$

(or > 0)

$$x^T A x \geq \lambda_{\min} \cdot x^T x$$

$$x^T Q x + b x$$

if  $A \neq A^T$        $A$  is square matrix

$$x^T A x \stackrel{?}{\geq} 0$$

$$(x^T A x)^T = x^T A^T x \quad x^T \left(\frac{A+A^T}{2}\right) x$$

$$A \stackrel{?}{\geq} B \quad \text{when } A - B \geq 0$$

$$Ax_i = \lambda_i x_i$$

$$A^2 x_i = A(Ax_i) = \lambda_i Ax_i = \lambda_i^2 x_i$$

i

$$A^n x_i = \lambda_i^n x_i$$

$$A^{-n} x_i = \lambda_i^{-n} x_i.$$

$$A = [y_1^T \dots y_k^T] \begin{bmatrix} \lambda_1 & & \\ & \ddots & \\ & & \lambda_k \end{bmatrix} [x_1 \dots x_k]$$

$$A = \lambda_1 x_1 y_1^T + \lambda_2 x_2 y_2^T + \dots + \lambda_k x_k y_k^T$$

$$A^n = \lambda_1^n x_1 y_1^T + \dots + \lambda_k^n x_k y_k^T$$

$$A^n = \lambda_1^n (x_1 y_1^T + (\frac{\lambda_2}{\lambda_1})^n x_2 y_2^T + \dots)$$

less than 1  
while  $n \uparrow$ .  $(\frac{\lambda_2}{\lambda_1})^n \downarrow$   
all parts disappear. same as following part.

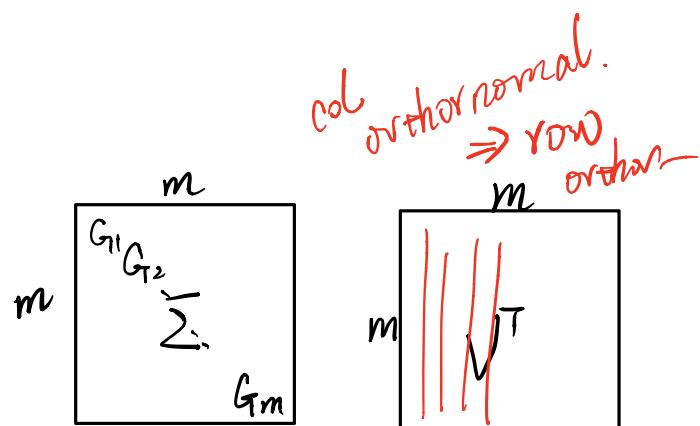
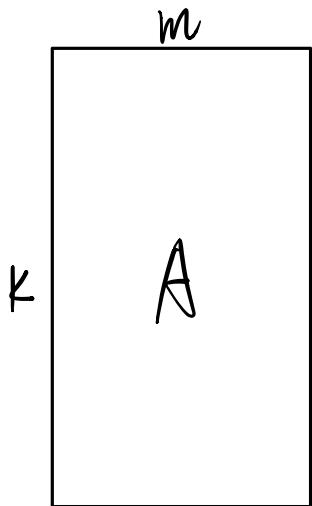
$$y_1 = Ax_0$$

$$y_2 = Ax_1$$

$$y_3 = Ax_2$$

:

$$y_n = \underset{\downarrow}{\text{some scalar.}} \quad Ax_1$$



$$G_1 \geq G_2 \geq G_3 \geq \dots \geq G_m \geq 0$$

$$U^T U = I$$

$$V^T V = I$$

$$A = U \Delta U^{-1}$$

eigenvalue of  $A$

$A$  is symmetric . nonnegative definite.

$$V = U$$

$$SVD(A) = U \Sigma U^T$$

$$\Sigma = \begin{pmatrix} G_1 & & \\ G_2 & \ddots & 0 \\ 0 & \ddots & G_n \end{pmatrix} \quad G_1 > G_2 > \dots > G_n$$

$(u_1, u_2, \dots, u_n)$   
↑  
orthonormal

$$(U_1 \ U_2 \ \dots \ U_n) \cdot \begin{pmatrix} G_1 & & \\ & G_{T_2} & \\ & & \ddots \\ & & & G_{Tn} \end{pmatrix} \cdot \begin{pmatrix} U_1^T \\ U_2^T \\ \vdots \\ U_n^T \end{pmatrix}$$

↓

$$G_1 U_1 U_1^T + \dots + G_n U_n U_n^T = A$$

$$A^n = G_1^n U_1 U_1^T + \dots + G_n^n U_n U_n^T = G_1^n (U_1 U_1^T + (\frac{G_2}{G_1})^n U_2 U_2^T + \dots + (\frac{G_n}{G_1})^n U_n U_n^T)$$

$$\lim_{n \rightarrow \infty} A^n = G_1^n U_1 U_1^T$$

Assume  $x_0^T U_1 \neq 0$

$$z_1 = Ax_0$$

$$x_1 = \frac{z_1}{\|z_1\|}$$

$$z_2 = Ax_1$$

$$x_2 = \frac{z_2}{\|z_2\|}$$

⋮

$$x_n \rightarrow ?$$

$$x_n = \frac{A^n x_0}{\|A^n x_0\|} = \frac{G_1^n (U_1 (U_1^T x_0) + (\frac{G_2}{G_1})^n O(1))}{G_1^n (|U_1^T x_0| + (\frac{G_2}{G_1})^n O(0))}$$

$$x_n \rightarrow \frac{U_1 \cdot (U_1^T x_0)}{|U_1^T x_0|} = U_1 \cdot \text{sgn}(U_1^T x_0)$$

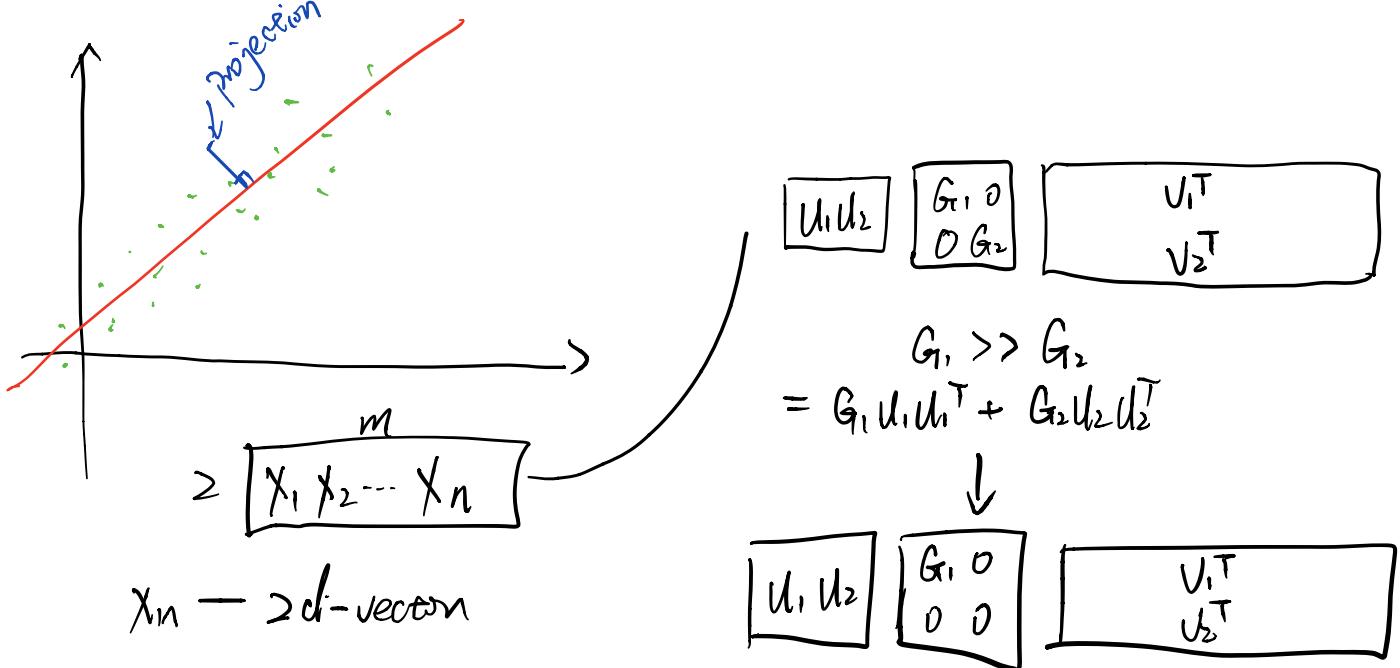
↓

Power method.

$$x_n \sim \pm U_1$$

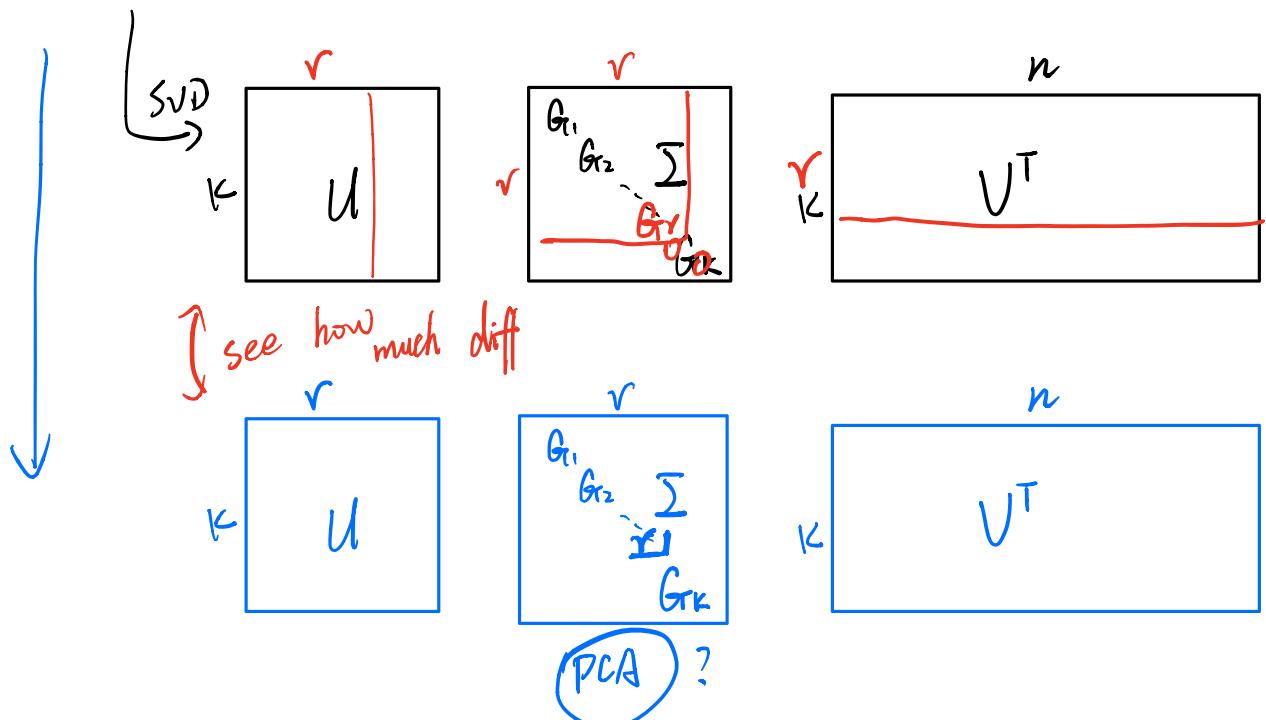
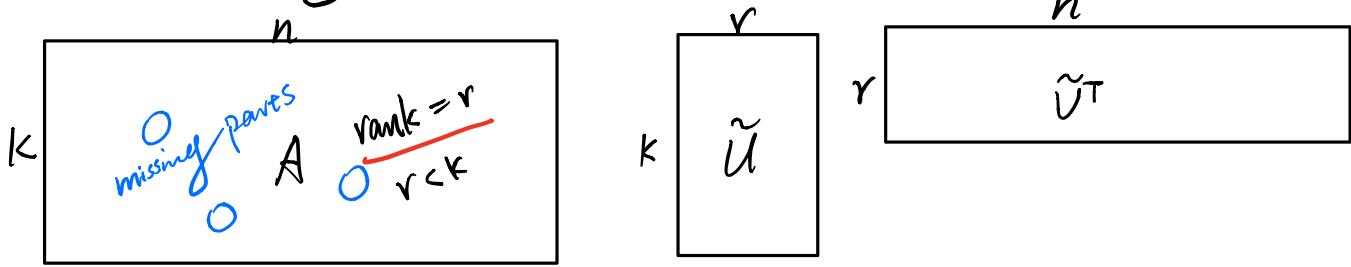
$$z_{n+1} = Ax_n = G_1(\pm U_1) \quad \|z_{n+1}\| = \underbrace{G_1 \|\pm U_1\|}_{\rightarrow 1}$$

Subspace tracking



PCA  
lose some information and keep most data.  
(features)

Matrix filling



# Vector space.

$\oplus \odot$

$V$  a set of vectors.

$x, y$  — vector

$x \oplus y \in V$

$\alpha, \beta$  — real num.

1)  $x \oplus (y \oplus z) = (x \oplus y) \oplus z$  associativity

2)  $x \oplus y = y \oplus x$  commutativity

3)  $x \oplus 0 = x$  identity

if  $V$  satisfy all 9 properties  
then  $V$  is a vector space.

4)  $(-x) \oplus x = 0$

$\downarrow$   
not negative  $x$   
it's a notation.

$x \in V, \alpha \in \mathbb{R}$

1)  $\alpha \odot x \in V$

2)  $\alpha \odot (\beta \odot x) = (\alpha \cdot \beta) \odot x$

3)  $1 \odot x = x$

4)  $\alpha \odot (x \oplus y) = (\alpha \odot x) \oplus (\alpha \odot y)$

5)  $(\alpha + \beta) \odot x = (\alpha \odot x) \oplus (\beta \odot x)$

## HILBERT SPACE

$\rightarrow$

$x, y \in V$

inner product is mapping from vector to  $\mathbb{R}$ .

$\langle x, y \rangle \in \mathbb{R}$

$\langle x, y \rangle = \langle y, x \rangle$  symmetric

$\langle \alpha \odot x, y \rangle = \alpha \cdot \langle x, y \rangle$

$\langle x + z, y \rangle = \langle x, y \rangle + \langle z, y \rangle$

$$\langle x, x \rangle \geq 0 \quad \text{if } \langle x, x \rangle = 0 \Leftrightarrow x = 0$$

$$\mathbb{R}^k \quad x = \begin{bmatrix} x_1 \\ \vdots \\ x_k \end{bmatrix}$$

$$\langle x, y \rangle = x^T y = x_1 y_1 + \dots + x_k y_k$$

$$\langle x, y \rangle = w_1 x_1 y_1 + \dots + w_k x_k y_k$$

$$w_1, w_2, \dots, w_k > 0$$

$$x^T \begin{pmatrix} w_1 & w_2 & \dots & w_k \end{pmatrix} y$$

$x^T \Omega y$  -  $\Omega$ : symmetric & positive definite.

$$x^T \Omega y \geq 0$$

$$\phi(x) = \int \phi(x) \psi(x) dx$$

$$\psi(x) = \int w(x) \phi(x) \phi(x) dx \quad (w(x) > 0)$$

$$\iint \underbrace{\Omega(x,y)}_{\hookrightarrow \text{kernels}} \phi(x) \psi(y) dx dy$$

Expectation.

$$x, y \in \mathcal{V} \rightarrow_{\text{random variables.}} \langle x, y \rangle = \bar{E}[x, y]$$

$$\mathcal{V} . \langle x, y \rangle , \|x\| = \sqrt{\langle x, x \rangle}$$

$$\text{def: } \cos \phi = \frac{\langle x, y \rangle}{\|x\| \cdot \|y\|}$$

Cauchy - Schwarz inequality

$$(\langle x, y \rangle)^2 \leq \|x\|^2 \cdot \|y\|^2 \leftarrow \text{only and only if } y = r \cdot x^{\text{number.}}$$

$$\langle \alpha x + y, \alpha x + y \rangle \geq 0$$

$$= \alpha^2 \langle x, x \rangle + 2\alpha \langle x, y \rangle + \langle y, y \rangle \geq 0 \quad \alpha = -\frac{\langle x, y \rangle}{\langle x, x \rangle}$$

$$\langle y, y \rangle - \frac{(\langle x, y \rangle)^2}{\langle x, x \rangle} \geq 0 \rightarrow \langle y, y \rangle \cdot \langle x, x \rangle \geq (\langle x, y \rangle)^2$$

$$\|y\|^2 \cdot \|x\|^2 \geq (\langle x, y \rangle)^2$$

$$\mathbb{R}^k \quad x^T y$$

$$(x_1 y_1 + \dots + x_k y_k)^2 \leq (x_1^2 + \dots + x_k^2)(y_1^2 + \dots + y_k^2)$$

$$\left( \int \phi(t) \psi(t) dt \right)^2 \leq \left( \int \phi^2(t) dt \right) \left( \int \psi^2(t) dt \right)$$

$$(\bar{E}[x, y])^2 \leq \bar{E}[x^2] \cdot \bar{E}[y^2]$$

$$y = r x.$$

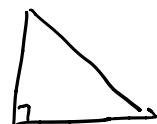
$$x, y \quad \langle x, y \rangle = 0$$

$$\cos \phi = \frac{\langle x, y \rangle}{\|x\| \cdot \|y\|}$$

$$\|x - y\|^2 = \langle x - y, x - y \rangle$$

$$= \langle x, x \rangle + \langle y, y \rangle - 2 \langle x, y \rangle$$

$$\|x - y\|^2 = \|x\|^2 + \|y\|^2$$

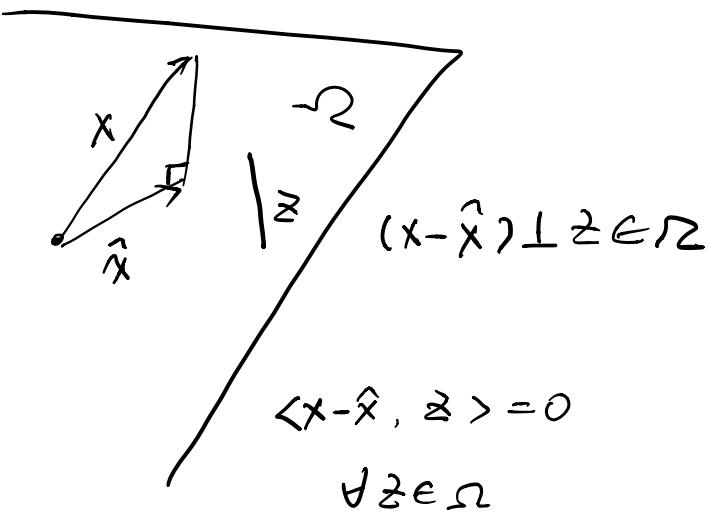


Pythagorean Theorem.

## Orthogonality principle

$$\therefore \langle x, y \rangle$$

$$\Omega \quad x, y \in \Omega \quad \alpha x + \beta y \in \Omega$$



$$\|x - y\|^2 = \langle x - y, x - y \rangle = \langle x - \hat{x} + \hat{x} - y, x - \hat{x} + \hat{x} - y \rangle$$

$$\langle x - \hat{x}, x - \hat{x} \rangle + \langle \hat{x} - y, \hat{x} - y \rangle + 2\langle x - \hat{x}, \hat{x} - y \rangle$$

$$\|x - \hat{x}\|^2 + \|\hat{x} - y\|^2 \geq \|x - \hat{x}\|^2$$

$$\Omega = \{\alpha_1 z_1 + \dots + \alpha_k z_k \mid z_1, z_2, \dots, z_k$$

$$\hat{x} = \alpha_1 z_1 + \dots + \alpha_k z_k$$

$$\langle x - \alpha_1 z_1 - \dots - \alpha_k z_k, z_1 \rangle = 0$$

$\vdots$

$\vdots$

$\vdots$

$\vdots$

$$\begin{bmatrix} \langle x, z_1 \rangle \\ \langle x, z_2 \rangle \\ \vdots \\ \langle x, z_k \rangle \end{bmatrix} = \begin{bmatrix} \langle z_1, z_1 \rangle & \langle z_1, z_2 \rangle & \dots \\ \langle z_2, z_1 \rangle & \langle z_2, z_2 \rangle & \dots \\ \vdots & \vdots & \ddots \\ \langle z_k, z_1 \rangle & \langle z_k, z_2 \rangle & \dots & \langle z_k, z_k \rangle \end{bmatrix} \begin{bmatrix} \alpha_1 \\ \alpha_2 \\ \vdots \\ \alpha_k \end{bmatrix}$$

solve this equation ↑

Application: Kalman filter

Total probability

$$A_1, A_2, \dots, A_k \quad A_i \cap A_j = \emptyset \quad (i \neq j)$$

$$A_1 \cup A_2 \cup \dots \cup A_k = \textcircled{H}$$

event B.  $P(B) = \sum_{i=1}^k \underbrace{P(B|A_i)}_{\text{density}} \underbrace{P(A_i)}_{\text{probability}}$   $B = \{x < m < x + dx\}$

random events  $H_1, H_2 \dots H_k$

$$\begin{array}{lll} H_1, f_1(x), \boxed{P(H_1)} \\ H_2, f_2(x), \boxed{P(H_2)} \\ \vdots \\ H_k, f_k(x), \boxed{P(H_k)} \end{array}$$

$\Sigma = 1$

$$\begin{aligned} f(x)dx &= \sum_{i=1}^k P(x < m \leq x + dx | H_i) P(H_i) \\ &= \sum_{i=1}^k f_i(x) P(H_i) \end{aligned}$$

$$f(x) = f_1(x|A_1) P(H_1) + f_2(x|A_2) P(H_2) + \dots + f_k(x|A_k) P(H_k)$$

Posterior Probability

$$P(A_j|B) = \frac{P(B|A_j) P(A_j)}{\sum_{i=1}^k P(B|A_i) P(A_i)}$$

$$P(H_j | m=x) = \frac{f_j(x) P(H_j)}{\sum_{i=1}^k f_i(x) P(H_i)}$$

(0) none of these guys 10%

$$H_0: f_0(x), P(H_0)$$

(1) George 50%

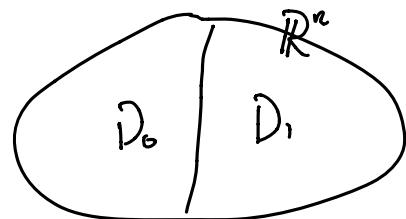
$$H_1: f_1(x), P(H_1)$$

(2) Jim 20%

find a decision strategy.

(3) John 20%

$$x = [0.1, 0.33, \dots] \quad D_0. \quad x \rightarrow D_0 \text{ or } D_1$$



Randomized

$$\delta_0(x), \delta_1(x)$$



$$\delta_0(x) = \begin{cases} 1 & . x \in H_0 \\ 0 & . x \in H_1 \end{cases} \quad \delta_1(x) = \begin{cases} 1 & . x \in H_1 \\ 0 & . x \in H_0 \end{cases}$$

Bayesian Approach

$$H_0 = f_0(x), \bar{P}(H_0) \quad \{D_0, H_0\} \quad C_{00}$$

$$H_1 = f_1(x), \bar{P}(H_1) \quad \{D_0, H_1\} \quad C_{01}$$

$$\{D_1, H_0\} \quad C_{10}$$

$$\{D_1, H_1\} \quad C_{11}$$

$$C(\delta_0, \delta_1) = C_{00} \bar{P}(D_0, H_0) + C_{01} \bar{P}(D_0, H_1) + C_{10} \bar{P}(D_1, H_0) + C_{11} \bar{P}(D_1, H_1)$$

for a specific data

optimal  $\neq$  best      optimal is average minimal dist.

$$\bar{P}(D_i \& H_j) = \bar{P}(D_i | H_j) \cdot \bar{P}(H_j) \quad \bar{P}(D_i | H_j) = \int \delta_i(x) f_j(x) dx$$

$$\begin{aligned}
 e(\delta_0, \delta_1) &= \sum_{i=0}^1 \sum_{j=0}^1 C_{ij} \bar{P}(H_j) \int \delta_i(x) f_j(x) dx \\
 &= \int (\delta_0(x) \left\{ \frac{C_{00} \bar{P}(H_0) f_0(x) + C_{01} \bar{P}(H_1) f_1(x)}{C_0(x)} \right\} + \delta_1(x) \left\{ \frac{C_{10} \bar{P}(H_0) f_0(x)}{C_1(x)} \right. \\
 &\quad \left. + C_{11} \bar{P}(H_1) f_1(x) \right\}) dx \\
 &= \int (\delta_0(x) C_0(x) + \delta_1(x) C_1(x)) dx
 \end{aligned}$$

$$\begin{aligned}
 C_1 > C_0, \quad \delta_1 = 0, \quad \delta_0 = 1 & \quad \delta_0 C_0 \geq \delta_0 \min(C_0, C_1) \\
 C_1 < C_0, \quad \delta_1 = 1, \quad \delta_0 = 0 & \quad \delta_1 C_1 \geq \delta_1 \min(C_0, C_1) \\
 C_1 = C_0, \quad \text{throw a coin} & \quad \delta_0 C_0 + \delta_1 C_1 \geq \min(C_0, C_1)
 \end{aligned}$$

$$\frac{f_1(x)}{f_0(x)} \stackrel{H_1}{>} \frac{(C_{10} - C_{00}) \bar{P}(H_0)}{(C_{01} - C_{11}) \bar{P}(H_1)} \quad \Leftarrow \text{Likelihood Ratio Test}$$

minimal average cost.

$$H_0 : f_0(x), \bar{P}(H_0)$$

$$H_1 : f_1(x), \bar{P}(H_1)$$

$$e(\delta_0, \delta_1) = \sum_{i=0}^1 \sum_{j=0}^1 C_{ij} \bar{P}(D_i \& H_j)$$

$$\frac{f_1(x)}{f_0(x)} \stackrel{H_1}{>} \frac{(C_0 - C_{00}) \bar{P}(H_0)}{(C_{01} - C_{11}) \bar{P}(H_1)}$$

$$C_{00} = C_{11} = 0$$

$$\bar{P}(D_0 \& H_1) + \bar{P}(D_1 \& H_0)$$

$$C_{01} = C_{10} = 1$$

$$\frac{f_1(x)}{f_0(x)} \stackrel{H_1}{\geq} \frac{\bar{P}(H_0)}{\bar{P}(H_1)}$$

$$\bar{P}(H_1) f_1(x) \stackrel{H_1}{\geq} \bar{P}(H_0) f_0(x)$$

$$\frac{\bar{P}(H_1) f_1(x)}{C(x)} \stackrel{H_1}{\leq} \frac{\bar{P}(H_0) f_0(x)}{C(x)}$$

$$\frac{\bar{P}(H_1) f_1(x)}{\bar{P}(H_1) f_1(x) + \bar{P}(H_0) f_0(x)} \stackrel{H_1}{\leq} \frac{\bar{P}(H_0) f_0(x)}{\bar{P}(H_1) f_1(x) + \bar{P}(H_0) f_0(x)}$$

$$\bar{P}(H_1|x) \stackrel{H_1}{\geq} \bar{P}(H_0|x)$$

$$\text{Sign } (\bar{P}(H_1|x) - \bar{P}(H_0|x)) \stackrel{H_1}{>} \stackrel{H_0}{<} 0$$

\* present & future,

$$H_0 = f_0(x), \bar{P}(H_0) \quad x_1^0, x_2^0 \dots x_{n_0}^0 \quad \text{classification.}$$

$$H_1 = f_1(x), \bar{P}(H_1) \quad x_1^1, x_2^1 \dots x_{n_1}^1$$

$$U(x, \theta)$$

$$\frac{1}{n_0 + n_1} \left( \sum_{i=1}^{n_0} (1 + U(x_i^0, \theta))^2 + \sum_{i=1}^{n_1} (1 - U(x_i^1, \theta))^2 \right)$$

$$\frac{n_0}{n_0 + n_1} \frac{1}{n_0} \sum_{i=1}^{n_0} (1 + U(x_i^0, \theta))^2$$

$$n_0, n_1 \rightarrow \infty$$

$$\bar{P}(H_0) + \bar{E}_0 [(1 + U(x, \theta))^2] \xrightarrow[\text{neural network}]{\text{large}} \bar{P}(H_0) + \bar{E}_0 [(1 + U(x))^2]$$

$$\min_u \overline{P}(H_0) E_0 [(1 + u(x))^2] + \overline{P}(H_1) E [(1 - u(x))^2]$$

$$u(x) = \frac{\overline{P}(H_1|x) - \overline{P}(H_0|x)}{\overline{P}(H_1|x) + \overline{P}(H_0|x)} \rightarrow \text{equal to Bayesian.}$$

Hinge - cost

\* only two classification.

Cross - Entropy

multiple hypothesis

$$H_0 = f_0(x), \overline{P}(H_0)$$

$$S_0(x), S_1(x), \dots, S_{k-1}(x)$$

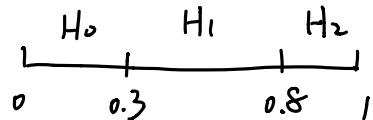
$$H_1 = f_1(x), \overline{P}(H_1)$$

:

$$S_0(x), S_1(x), S_2(x)$$

$$0.3 \quad 0.5 \quad 0.2$$

$$H_{k-1} = f_{k-1}(x), \overline{P}(H_{k-1})$$



$$\{D_i \& H_j\} C_{ij}$$

$$E(S_0, S_1, \dots, S_{k-1}) = \sum_{i=0}^{k-1} \sum_{j=0}^{k-1} C_{ij} \overline{P}(D_i | H_j) \overline{P}(H_j)$$

$$\int S_i(x) f_j(x) dx$$

$$\frac{\int \sum_{i=0}^{k-1} S_i(x) \left( \sum_{j=0}^{k-1} C_{ij} \overline{P}(H_j) f_j(x) \right) dx}{C_i(x)}$$

$$= \int \left( \sum_{i=0}^{k-1} S_i(x) C_i(x) \right) dx$$

$$\delta_0 \cdot c_0 + \delta_1 \cdot c_1 + \dots + \delta_{k-1} \cdot c_{k-1} \geq \min \{c_0, c_1, \dots, c_{k-1}\}$$

$$\delta_0 c_0 \geq \min \{c_0, c_1, \dots, c_{k-1}\} \quad \delta_0$$

$$\delta_1 c_1 \geq \delta_1 \quad \dots$$

⋮

$$\delta_{k-1} c_{k-1} \geq \delta_{k-1} \quad \dots$$

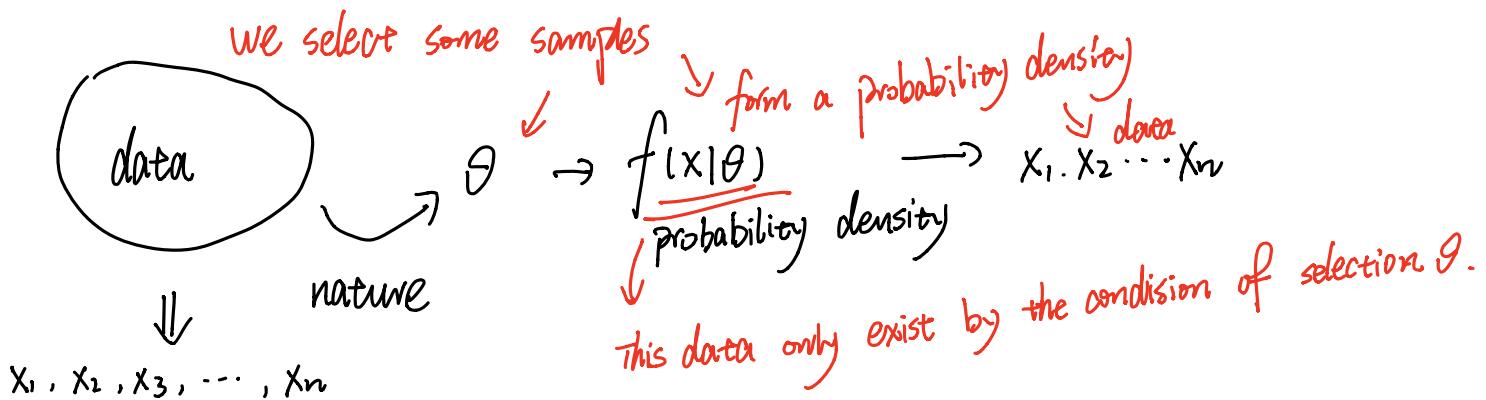
$$c_0 = \min \rightarrow \delta_0 = 1, \delta_1 = \delta_2 = \dots = \delta_{k-1} = 0$$

$$c_1 = \min \rightarrow \delta_1 = 1, \delta_0 = \delta_2 = \dots = \delta_{k-1} = 0$$

## PARAMETER ESTIMATION

$$x_1, x_2, \dots, x_m \quad \theta$$

  
Corr(x, y)



we don't know the data.

Estimate is a function of data.

It comes from nature.

We want to estimate it.

$X \rightarrow \theta(X)$

data      estimation

$x_1, x_2, \dots, x_n$

$$\hat{\mu} = \frac{x_1 + x_2 + \dots + x_n}{n}$$

good one

} both estimation.

$$\hat{\mu}' = \frac{x_1^2 + x_2^2 + x_3^2 + \dots + x_n^2}{n}$$

bad one

1.  $f(x|\theta)$  we can describe the behavior of  $\theta$ .

$\pi(\theta)$  probability density. We treat  $\theta$  as random.  
prior density  
hard to know

$$f(x|\theta) \cdot \pi(\theta) = f(x, \theta)$$

$$f(\theta|x) = \frac{f(x, \theta)}{f(x)} = \frac{f(x, \theta)}{\int f(x, \theta) d\theta}$$

Bayesian Estimation Theory

from  $\theta$  data  $x, \theta$  true  $\hat{\theta}$  estimate  
 $C(\theta, \hat{\theta})$  cost function

function of a function of  $x$ .

joint pdf

$$\text{Average cost} = \iint C(\theta, \hat{\theta}(x)) f(x, \theta) dx d\theta = \underline{e(\hat{\theta})}$$

$$\iint C(\theta, \hat{\theta}(x)) f(\theta|x) \cdot f(x) dx d\theta$$

$$\int \left( \int C(\theta, \hat{\theta}(x)) f(\theta|x) d\theta \right) f(x) dx$$

$$\int \left( \int C(\theta, u) f(\theta|x) d\theta \right) f(x) dx = \int G(x, u) f(x) dx$$

$$\text{number} \rightarrow \min_{x, u} G(x, u) \quad \frac{dx}{\partial x} \quad \frac{du}{\partial u}$$

$$\text{function} \rightarrow \min_u G(x, u) = \phi(x) \quad G(x, u) \geq \phi(x) \quad G(x, u_0(x)) = \phi(x)$$

$$G(x, u) = (u - \log(1+x^2))^2 + x^4$$

$$\min_{u,x} G(x, u) = 0 \quad u = x = 0$$

$$\min_u G(x, u) = x^4 \quad u(x) = \log(1+x^2)$$

$$E(\hat{\theta}) = \int \left( \underbrace{\int C(\theta, \hat{\theta}(x)) f(\theta|x) d\theta}_{G(x, \hat{\theta}(x))} \right) f(x) dx \\ \geq \phi(x)$$

$$E(\hat{\theta}) = \int G(x, \hat{\theta}(x)) f(x) dx \geq \underbrace{\int \phi(x) f(x) dx}_{\text{a known number}}$$

find  $\hat{\theta}$  satisfy above.

$$\hat{\theta}_o(x) = \arg \min_u G(x, u)$$

$$G(x, u) = \int C(\theta, u) f(\theta|x) d\theta.$$

$$\nabla_u G(x, u) = 0$$

$$\int \nabla_u C(\theta, u) f(\theta|x) d\theta = 0 \quad C(\theta, \hat{\theta}) = \|\theta - \hat{\theta}\|^2 \quad \text{MMSE}$$

minimal min square estimator

$$(\theta - u)^T (\theta - u) = \underline{\theta^T \theta - 2u^T \theta + u^T u}$$

$$\nabla_u \phi(u_1, u_2, \dots, u_k) = \begin{pmatrix} \frac{\partial \phi}{\partial u_1} \\ \frac{\partial \phi}{\partial u_2} \\ \vdots \\ \frac{\partial \phi}{\partial u_k} \end{pmatrix} \quad \begin{array}{l} \Downarrow \text{vector} \\ -2\theta + 2u \\ = \underline{-2(u - \theta)} \\ \text{vector} \end{array}$$

$$(1-2) \int \underline{(u - \theta)} f(\theta|x) d\theta = \underline{0} \quad \text{vector}$$

$$\cancel{\int_0^u f(\theta|x) d\theta} = \int \theta f(\theta|x) d\theta \quad \hat{\theta}(x) = \underline{\mathbb{E}[\theta|x]}$$

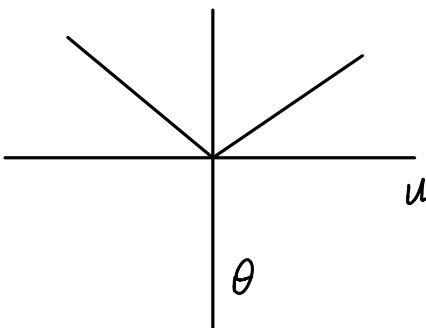
$$u = \underline{\mathbb{E}[\theta|x]}$$

$\theta$  scalar

$$C(\theta, \hat{\theta}) = |\theta - \hat{\theta}|$$

$$C(\theta, u) = |u - \theta|$$

$$\frac{d|u-\theta|}{du} = \text{sgn}(u-\theta)$$



$$\int \text{sgn}(u-\theta) f(\theta|x) d\theta = \int_{-\infty}^u 1 \cdot f(\theta|x) d\theta + \int_u^\infty f(\theta|x) d\theta$$

cdf

$$\bar{F}(u|x) - (1 - \bar{F}(u|x)) = 0$$

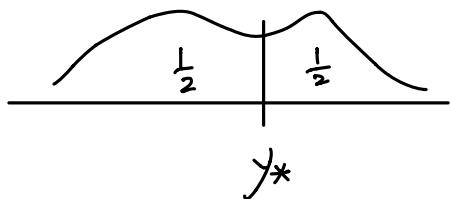
$$\bar{F}(u|x) = \frac{1}{2}$$

$y, f(y), F(y)$

$$\underset{\text{median}}{\leftarrow} F(y_*) = \frac{1}{2}$$

$$C(\theta, \hat{\theta}) = |\theta_1 - \hat{\theta}_1| + |\theta_2 - \hat{\theta}_2| + \dots$$

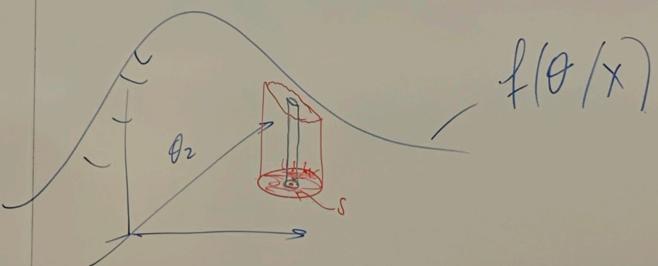
↑ median of ↑ median of



MAP Estimator

$$C(\theta, \hat{\theta}) = \begin{cases} 0 & \|\theta - \hat{\theta}\| = \delta \\ 1 & \text{otherwise} \end{cases}$$

$$\begin{aligned} G(x, u) &= \int C(\theta, u) f(\theta|x) d\theta \\ &= \int_0^\infty f_\theta(u|x) d\theta + \int_1^\infty f_\theta(\theta|x) d\theta \\ &\quad \|\theta-u\| \leq \delta \quad \|\theta-u\| > \delta \\ &\quad + \int_{\|\theta-u\| \leq \delta} f_\theta(\theta|x) d\theta - \int_{\|\theta-u\| > \delta} f_\theta(\theta|x) d\theta \end{aligned}$$



$$= 1 - \int f(\theta|x) d\theta$$

$$\|\theta - u\| \leq \delta$$

$\delta$  is very small.

$$\approx B_\delta \cdot f(u|x)$$

$$\hat{\theta}_{MAP}(x) = \arg \max_{\theta} f(\theta|x)$$

$$\arg \max_{\theta} \frac{f(x, \theta)}{f(x)} = \arg \max_{\theta} f(x, \theta)$$

$$\arg \max_{\theta} f(x|\theta) \pi(\theta)$$

$$\pi(\theta) = \begin{cases} \frac{1}{V(\Omega)} & \theta \in \Omega \\ 0 & \text{otherwise.} \end{cases}$$

$$\arg \max_{\theta} f(x|\theta) \frac{1}{V(\Omega)}$$

number

$$\hat{\theta}_{MLE}(x)$$

$$\hat{\theta}(x) \quad \text{unbiased estimator}$$

$$\text{average} : \int \hat{\theta}(x) f(x|\theta) dx = \theta$$

$$\int (\hat{\theta}(x) - \theta)(\hat{\theta}(x) - \theta)^T f(x|\theta) dx = \sum \hat{\theta} \geq (I_{\hat{\theta}}(\theta))^{-1}$$

Fisher information matrix

Grammer-Rao Lower Bound

CRLB

$$\int \left( \frac{\nabla_{\theta} f(x|\theta)}{f(x|\theta)} \right) \left( \quad \right)^T f(x|\theta) dx = I_{F2}(\theta)$$

$x_1, x_2, \dots, x_n$  independent

$x_i \sim N(\mu, 1)$  Gaussian.

$\bar{x}(M) \sim N(\mu_0, \sigma^2)$

$$f(x_1, x_2, \dots, x_n / M) = \frac{1}{\sqrt{2\pi n}} e^{-\frac{1}{2} \sum_{i=1}^n (x_i - M)^2} = \left( \frac{1}{\sqrt{2\pi n}} \right)^n \cdot e^{-\frac{1}{2} \sum_{i=1}^n (x_i - M)^2}$$

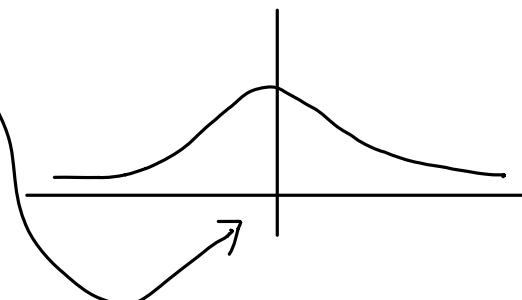
$$\min_M \sum_{i=1}^n (x_i - M)^2 \rightarrow \sum_{i=1}^n (-2)(x_i - M) = 0$$

$$\frac{1}{n} \sum_{i=1}^n M = M$$

$$f(M/x_1, x_2, \dots, x_n) = N\left(\frac{\sum_{i=1}^n x_i + \mu_0}{n + \frac{1}{\sigma^2}}, \frac{1}{n + \frac{1}{\sigma^2}}\right)$$

$$M_{MMSE} = \bar{E}[M/x] = \frac{\sum_{i=1}^n x_i + \mu_0}{n + \frac{1}{\sigma^2}}$$

$$M_{\text{median}} =$$



## Optimal problem

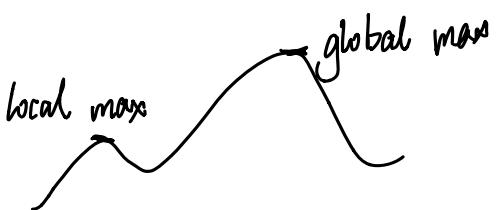
- 1) function optimize with constraint
- 2) linear programming
- 3) optimize functional (function of function)

Calculus of variations

$$g(x_1, x_2, \dots, x_n) = g(x) \quad x = \begin{pmatrix} x_1 \\ x_2 \\ \vdots \\ x_n \end{pmatrix}$$

$\max_x g(x)$  exist a  $x_*$  that  $g(x_*) \geq g(x) \rightarrow$  global maximum.

**supreme**  $\sup_x g(x)$   $x_*$  may not exist but there's a max limit.



there's no alg. to find global max.  
It's always a local max.

Hessian matrix =  $H_{ij} = \frac{\partial^2 g}{\partial x_i \partial x_j}$  matrix function of  $X$ .

$H(x_*) \rightarrow$  positive definite (min)  
negative definite (max).  
neither (subtle points)

pos def matrix

positive def:  $x^T Q x > 0$

$x \neq 0$

$$x = x_* + \delta$$

$$= 0$$

$$g(x) = g(x_* + \delta) = g(x_*) + (\nabla_x g(x_*))^T \delta + \frac{1}{2} \delta^T H(x_*) \delta$$

$$g(x_* + \delta) - g(x_*) = \frac{1}{2} \delta^T H(x_*) \delta \text{ compare } 0$$

$$\max_x g(x_1, x_2, \dots, x_n) = g(x)$$

$\downarrow x$  has to satisfy

$$w(x) = 0$$

$$\max_{x,y} g(x, y) \quad x=5$$

$$x+y=1 \quad y=1-x \quad \max_x g(x, 1-x)$$

$$\max_{x_1, x_2} g(x_1, x_2)$$

$$\max_{x_1} \left\{ \max_{x_2} g(x_1, x_2) \right\}$$

$$\max_{x_2} \left\{ \max_{x_1} g(x_1, x_2) \right\}$$

It's like going on the ridge while climbing a mountain.

These three are equal.

Lagrange multiplier method

$$\max_x g(x) \quad \begin{cases} w_1(x) = 0 \\ w_2(x) = 0 \end{cases}$$

$$\tilde{g}(x) = g(x) + \lambda_1 w_1(x) + \lambda_2 w_2(x)$$

Lagrange factors.

$$\underset{x}{\text{maximize}} \quad \tilde{g}(x) \rightarrow x^*(\lambda_1, \lambda_2)$$

$$w_1(x^*(\lambda_1, \lambda_2)) = 0 \Rightarrow \lambda_1^*, \lambda_2^* \rightarrow x^*(\lambda_1^*, \lambda_2^*)$$

$$w_2(x^*(\lambda_1, \lambda_2)) = 0 \quad \text{possible solution.}$$

$$\mathcal{X} = \{x : \begin{cases} w_1(x) = 0 \\ w_2(x) = 0 \end{cases}\}$$

$$\max_{x \in \mathcal{X}} g(x)$$

$$g(x) = g(x) + \lambda_1^* w_1(x) + \lambda_2^* w_2(x) = \tilde{g}(x)$$

$$\tilde{g}(x) \leq \tilde{g}(x^*(\lambda_1^*, \lambda_2^*))$$

$$= g(x^*(\lambda_1^*, \lambda_2^*)) + \lambda_1^* w_1(x^*(\lambda_1^*, \lambda_2^*)) + \lambda_2^* w_2(x^*(\lambda_1^*, \lambda_2^*))$$

Example.

$Q$  positive definite

$$g(x) = x^T Q x > 0 \quad x \neq 0$$

$$\min_x x^T Q x = 0$$

$$x = 0$$

$$A^T x = b$$

$$\begin{cases} a_1^T x = b_1 \\ a_2^T x = b_2 \\ \vdots \\ a_k^T x = b_k \end{cases} \quad k = \text{dimension of } x.$$

$$\text{find } \min_x x^T Q x.$$

$$x^T Q x + \lambda_1(a_1^T x - b_1) + \lambda_2(a_2^T x - b_2) + \dots + \lambda_k(a_k^T x - b_k)$$

$$\text{take gradient} = 2Qx + \lambda_1 a_1 + \lambda_2 a_2 + \dots + \lambda_k a_k = 0$$

$$[a_1 \ a_2 \ \dots \ a_k] \begin{bmatrix} \lambda_1 \\ \lambda_2 \\ \vdots \\ \lambda_k \end{bmatrix} \lambda$$

$$2Qx + A\lambda = 0 \rightarrow x^* = -\frac{1}{2}Q^{-1}A\lambda$$

$$A^T x = -\frac{1}{2}A^T Q^{-1} A \lambda = b \Rightarrow \lambda^* = -2(A^T Q^{-1} A)^{-1} b$$

$$x^* = -\frac{1}{2}Q^{-1}A \cdot -2(A^T Q^{-1} A)^{-1} b$$

$$= Q^{-1}A(A^T Q^{-1} A)^{-1} b$$

## LINEAR PROGRAMMING.

$$x_1, x_2, \dots, x_n$$

$$\min_{x_1, \dots, x_m} \{g_1 x_1 + g_2 x_2 + \dots + g_m x_m\}$$

$$G = \begin{bmatrix} g_1 \\ g_2 \\ \vdots \\ g_m \end{bmatrix} \quad X = \begin{bmatrix} x_1 \\ x_2 \\ \vdots \\ x_m \end{bmatrix}$$

Constraints:

$$\alpha_{11} x_1 + \dots + \alpha_{1m} x_m = b_1$$

linear      :

$$A_e x = b_e$$

equality

$$\alpha_{m1} x_1 + \dots + \alpha_{mm} x_m = b_m$$

$$\alpha_{11} x_1 + \dots + \alpha_{1n} x_n \leq \beta_1$$

$$\min_x G^T x$$

non linear    :

$$A x \leq b$$

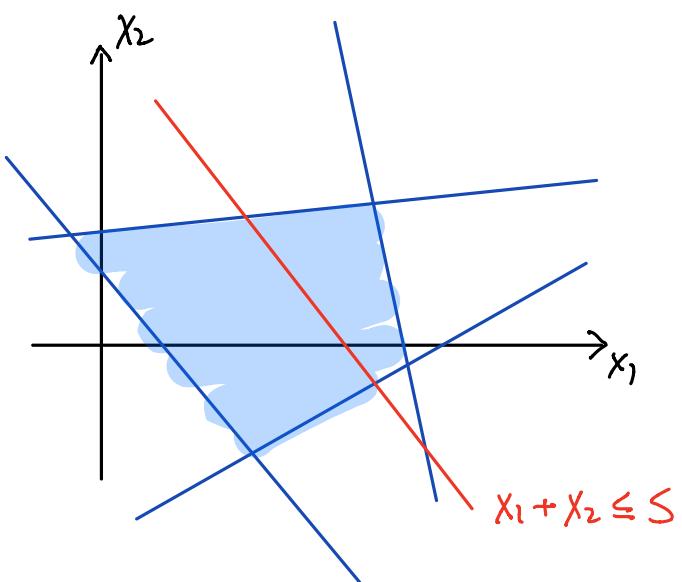
$$\alpha_{k1} x_1 + \dots + \alpha_{kn} x_n \leq \beta_k$$

$$\min_{x_1, x_2} \{x_1 + x_2\}$$

$$x_1 + 2x_2 \leq 2$$

$$3x_1 + x_2 \leq 0$$

⋮



$$k_1 = |C_{11}x_1 + \dots + C_{1n}x_n|$$

$$k_2 = |C_{21}x_1 + \dots + C_{2n}x_n| \quad Ax = b$$

$$k_s = |C_{s1}x_1 + \dots + C_{sn}x_n| \quad Ax \leq b$$

$$\min_{x_1, \dots, x_n} \left\{ \max \{k_1, k_2, \dots, k_s\} \right\}$$

$$-\delta \leq C_{11}x_1 + \dots + C_{1n}x_n \leq \delta$$

⋮

$$-\delta \leq C_{s1}x_1 + \dots + C_{sn}x_n \leq \delta$$

$$g_1x_1 + g_2x_2 + \dots + g_nx_n + g_s \cdot \delta$$

## CALCULUS OF VARIATION

$$F(x, y)$$

$$g(x) \quad \min_{g(x)} \int_{x_1}^{x_2} F(x, g(x)) dx \quad g_o(x) + \epsilon \delta(x)$$

$$\Omega(\epsilon) = \int_{x_1}^{x_2} F(x, g_o(x) + \epsilon \delta(x)) dx \geq \Omega(0)$$

$$\Omega'(0) = 0$$

$$\Omega'(\epsilon) = \int_{x_1}^{x_2} \frac{d}{d\epsilon} F(x, g_o(x) + \epsilon \delta(x)) dx$$

$$= \int_{x_1}^{x_2} \frac{\partial}{\partial y} F(x, g_o(x)) \underline{\delta(x)} dx = 0$$

any function.

$$\int_{x_1}^{x_2} \phi(x) \cdot \delta(x) dx = 0 \Rightarrow \phi(x) = 0$$

$$\frac{\partial}{\partial y} \bar{F}(x, g_0(x)) = 0$$

$$\frac{\partial}{\partial y} \bar{F}(x, y) = \bar{F}_y(x, y) = 0$$

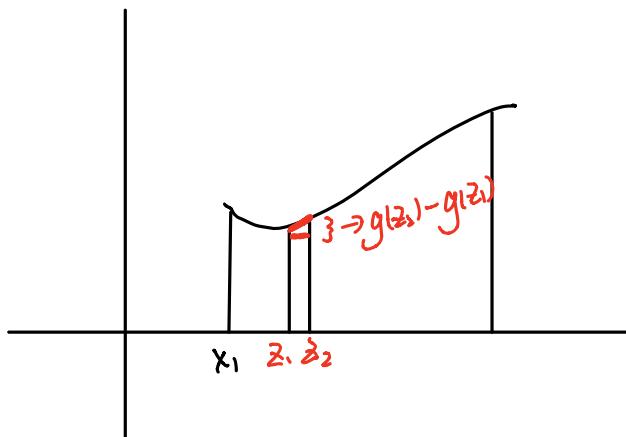
$$F(x, y, z) \quad g(x) = g_0(x) + \epsilon \delta(x)$$

$$\min_{g(x)} \int_{x_1}^{x_2} \bar{F}(x, g(x), g'(x)) dx.$$

$$g(x_1) = \alpha_1$$

$$g(x_2) = \alpha_2$$

$$\frac{d}{dy} \bar{F}(x, g_0(x), g'_0(x)) - \frac{d}{dy} (\frac{\partial}{\partial \epsilon} \bar{F}(x, g_0(x), g'_0(x))) = 0.$$



$$\sqrt{(g(z_2) - g(z_1))^2 + (z_2 - z_1)^2}$$

$$\sum_n \sqrt{1 + \left(\frac{g(z_2) - g(z_1)}{z_2 - z_1}\right)^2} (z_2 - z_1)$$

$$l(g) = \int_{x_1}^{x_2} \sqrt{[g'(x)]^2 + 1} dx$$

$$F(x, y, z) = \sqrt{1+z^2}$$

$$\frac{d}{dx} \left( \frac{g'(x)}{\sqrt{1+[g'(x)]^2}} \right) = 0$$

$$\frac{g'(x)}{\sqrt{1+[g'(x)]^2}} = C \quad g(x) = \sqrt{\frac{C^2}{1-C^2}} x + b.$$

SVM : approximate non-linear function

kernels

1) Approx function

2) classification

image compare

Kernels approximate dirac delta function.

Dirac delta function

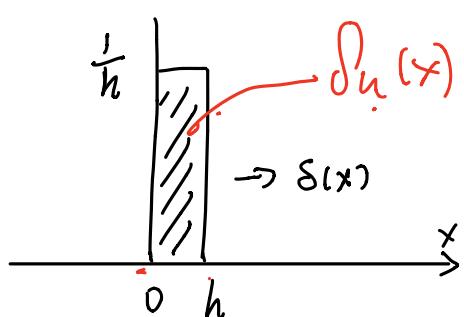
$\delta(x)$



$$\int_{-\varepsilon}^{\varepsilon} \delta(x) dx = 1$$

$$\int \phi(x) \delta(x) dx = \phi(0) \rightarrow \text{a way of definition}$$

$$\int \phi(x) \delta'(x) dx = -\phi'(0)$$



$$\int_{-\infty}^{\infty} g(x) dx = 1. \quad g(x, h) = g(\frac{x}{h}) \frac{1}{h}$$

$$\int_{-\infty}^{\infty} \phi(x) \delta_h(x) dx = \cancel{\int_0^h \phi(x) dx} \xrightarrow{h \rightarrow 0} \frac{1}{h} \int_0^h \phi(0) dx = \phi(0) \frac{1}{h} \int_0^h dx$$

$$g(x) = \frac{e^{-\frac{x^2}{2}}}{\sqrt{2\pi}} \quad (\text{gaussian approximation})$$

$$\int \phi(x) \delta(x) dx = \phi(0)$$

$$\delta(x) = \delta(x_1) \cdot \delta(x_2) \cdots \delta(x_k) \quad \text{multi-dimension of } \delta(x)$$

$$\delta(x_1) = \frac{e^{-\frac{1}{2} \frac{x_1^2}{h^2}}}{\sqrt{2\pi h^2}}$$

$$\delta(x, h) = \frac{e^{-\frac{1}{2} \frac{x_1^2}{h^2}}}{\sqrt{2\pi h^2}} \cdot \frac{e^{-\frac{1}{2} \frac{x_2^2}{h^2}}}{\sqrt{2\pi h^2}} \cdots = \frac{e^{-\frac{1}{2h^2} \|x\|^2}}{\sqrt{(2\pi)^k (h^2)^k}}$$

$$\int \phi(x) \delta(\underbrace{x - x_0}_z) dx = \phi(x_0)$$

$$\int \phi(z + x_0) \delta(z) dz = \phi(x_0) \rightarrow x = z + x_0$$

$$\phi(x_0) = \int \phi(x) \delta(x - x_0) dx \quad K(x, h) \text{ approximation of } \delta(x)$$

$$\approx \int \phi(x) K(x - x_0, h) dx.$$

$$\int_{-\infty}^{+\infty} f(x) \cdot \delta(x_0 - x) dx = f(x_0)$$

↓

$$f(x_0) = \bar{E}_x [\delta(y - x)] \approx \bar{E}_{x_0} [K(x_0 - x, h)]$$

$$\mathcal{X} : x_1, x_2, x_3 \cdots x_k$$

$$\bar{E}[x] \approx \frac{x_1 + x_2 + \cdots + x_k}{k}$$

$$\bar{E}[x^2] \approx \frac{x_1^2 + x_2^2 + \cdots + x_k^2}{k}$$

$$\bar{E}[g(x)] \approx \frac{g(x_1) + g(x_2) + \dots + g(x_k)}{k}$$

$$\bar{E}_x[k(x, h)] \approx \frac{k(x_1, h) + k(x_2, h) + \dots + k(x_k, h)}{k}$$

$$\bar{E}_x[k(y-x, h)] \approx \frac{k(y-x_1, h) + k(y-x_2, h) + \dots + k(y-x_k, h)}{k}$$

$$f(x) \approx \frac{k(x-x_1, h) + k(x-x_2, h) + \dots + k(x-x_k, h)}{k}$$

kernel function  
approximation of  $\delta(x)$

classification

$$\begin{matrix} \text{class} \\ 1 & \{x_1^1, x_2^1, \dots, x_k^1\} \\ 2 & \{x_1^2, x_2^2, \dots, x_k^2\} \end{matrix}$$

$$\frac{f_1(x)}{f_2(x)} \geq 1. \quad (\text{likelihood ratio})$$

greater than 1  $\rightarrow X^1$   
smaller than 1  $\rightarrow X^2$

OR. maybe make estimation  
of the entire ratio,

## MERCER KERNELS

$k(x, y)$  — scalar function

1)  $k(x, y) = k(y, x)$

$\phi(x)$  scalar function

2) positive definite :  $\iint \phi(x) \phi(y) k(x, y) dx dy > 0$ .  $\phi(x) \neq 0$

$$\iint \phi(x) \phi(y) k(x, y) dx dy = 0 \text{ if and only if } \phi(x) = 0$$

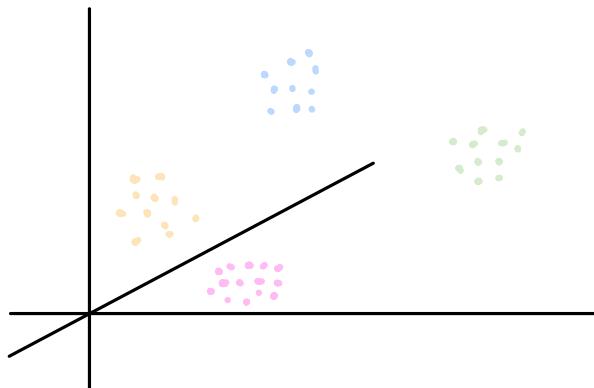
different kernels.

$$k(x, y) = (x^T y + C)^2$$

$$k(x, y) = e^{-\frac{1}{2h} \|x-y\|^2}$$

$$k(x, y) = e^{-\frac{1}{2h} \|x-y\|_L^2}$$

$$k(x, y) = e^{-\frac{1}{h} \|x-y\|_L} \quad \|z\|_L = |z_1| + |z_2| + \dots + |z_k|$$



$$x_0 \in \mathbb{R}^n \rightarrow k(x, x_0)$$

$$\phi(x) \in \mathcal{V}$$

$$\phi(x) = \alpha_1 k(x, x_1) + \dots + \alpha_m k(x, x_m)$$

↗  
vector space.

$$\phi(x) = \alpha_1 k(x, x_1) + \dots + \alpha_m k(x, x_m)$$

$$\psi(x) = \beta_1 k(x, x'_1) + \dots + \beta_{m'} k(x, x'_{m'})$$

$$\langle \phi, \psi \rangle \triangleq \sum_{i=1}^m \sum_{j=1}^{m'} \alpha_i \cdot \beta_j k(x_i, x'_j)$$

$$\langle \phi, \phi \rangle = \sum_{i=1}^m \sum_{j=1}^m \alpha_i \cdot \alpha_j \underbrace{k(x_i, x_j)}_{?} \geq 0$$

$$\phi(x) \in \mathcal{V}$$

$$\phi(x) = \alpha_1 k(x, x_1) + \dots + \alpha_m k(x, x_m)$$

$$k(x, x_0)$$

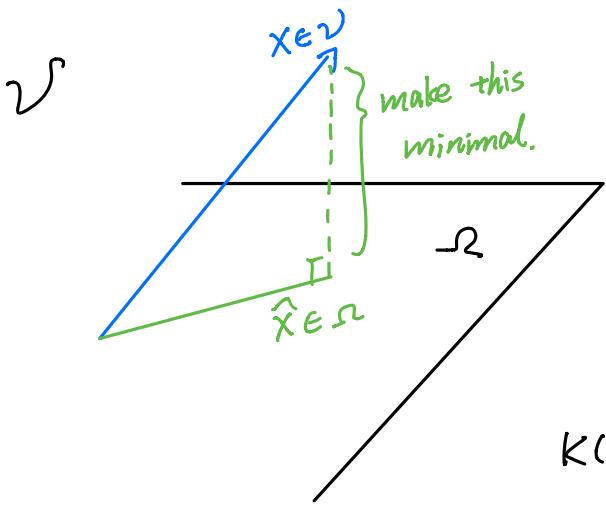
$$\langle k(x, x_0), \phi(x) \rangle = \underbrace{\alpha_1 \langle k(x, x_1), k(x, x_0) \rangle +}_{\dots} +$$

$$\dots + \alpha_m \langle k(x, x_m), k(x, x_0) \rangle /$$

$$\alpha_1 k(x_0, x_1) + \dots + \alpha_m k(x_0, x_m) = \phi(x_0)$$

$$\begin{bmatrix} k(x_1, y_1) & k(x_1, y_2) & \dots \\ \vdots & \ddots & \dots & \dots \\ k(x_m, y_1) & k(x_m, y_2) & \dots \end{bmatrix}$$

$k(x, y)$  is positive definite.



$$\mathcal{L} = \{\alpha_1 z_1 + \dots + \alpha_n z_n\}$$

$$z_1, z_2, \dots, z_n \in \mathcal{Z}$$

$$x - \hat{x} \perp z_1, z_2, \dots, z_n$$

$$K(x, x_1), K(x, x_2), \dots, K(x, x_n)$$

$$z_1, z_2, \dots, z_n$$

$$\mathcal{L} = \{\alpha_1 K(x, x_1) + \dots + \alpha_n K(x, x_n)\}$$

$$\phi(x) \in \mathcal{V} \quad \phi(x) = \beta_1 K(x, x'_1) + \dots + \beta_m K(x, x'_m)$$

$$\phi(x) - \hat{\phi}(x) \perp K(x, x_1), K(x, x_2), \dots, K(x, x_n)$$

$$\langle \phi(x) - \hat{\phi}(x), K(x, x_i) \rangle = 0 \quad \text{Reresenter Theorem.}$$

$$\underline{\phi(x_i)} = \langle \phi(x), K(x, x_i) \rangle = \langle \hat{\phi}(x), K(x, x_i) \rangle = \underline{\hat{\phi}(x_i)}$$

$$\phi(x_i) = \hat{\phi}(x_i)$$

$(y_1, x_1), (y_2, x_2), \dots, (y_k, x_k)$  — data

$$\min_{\phi(x) \in \mathcal{V}} \sum_{i=1}^k l(y_i, \hat{\phi}(x_i)) \quad \phi(x) \in \mathcal{V}$$

$$l(y, \phi(x)) = (y - \phi(x))^2$$

$\phi(x) = \alpha_1 K(x, x'_1) + \alpha_2 K(x, x'_2) + \dots + \alpha_m K(x, x'_m)$  every kernel function define a  $\mathcal{V}$  / vector space

$$\phi(x) = \alpha_1 K(x, x_1) + \alpha_2 K(x, x_2) + \dots + \alpha_n K(x, x_n)$$

$$\min_{\alpha_1, \dots, \alpha_n} \sum_{i=1}^k (y_i - \underline{\alpha_1 K(x_i, x_1)} - \underline{\alpha_2 K(x_i, x_2)} - \dots - \underline{\alpha_n K(x_i, x_n)})^2$$

Known numbers.

EASY TO OPTIMISE.

$$\phi(x) = \alpha_1^* K(x, x_1) + \dots + \alpha_n^* K(x, x_n) \quad \text{problem comes to linear.}$$

## Learning Algorithm.

$$G(\theta) = \begin{bmatrix} G_1(\theta_1, \dots, \theta_k) \\ G_2(\theta_1, \dots, \theta_k) \\ \vdots \\ G_K(\theta_1, \dots, \theta_k) \end{bmatrix}$$

$$G(\theta) = 0$$

Assume  $G(\theta)$  is continuous in  $\theta$ .

$$V_t = V_{t-1} - \mu \overset{\text{scalar}}{G}(V_{t-1})$$

$$V_t \rightarrow V_\infty$$

$$\lim_{t \rightarrow \infty} V_t = \lim_{t \rightarrow \infty} V_{t-1} - \mu \lim_{t \rightarrow \infty} G(V_{t-1})$$

$$V_\infty \quad V_\infty \quad G(V_\infty)$$

$$V_\infty = V_\infty - \mu G(V_\infty) \rightarrow G(V_\infty) = 0$$

$$g(\theta) = g(\theta_1, \theta_2, \dots, \theta_k)$$

$$\min_{\theta} g(\theta)$$

$$\nabla_{\theta} g(\theta) = 0$$

$$\nabla_{\theta} g(\theta) = \begin{bmatrix} \frac{\partial g}{\partial \theta_1}(\theta_1, \theta_2, \dots, \theta_k) \\ \vdots \\ \frac{\partial g}{\partial \theta_k}(\theta_1, \theta_2, \dots, \theta_k) \end{bmatrix}$$

## Gradient Descent

$$V_t = V_{t-1} - \underbrace{\mu}_{\text{step size}} \nabla_{\theta} g(V_{t-1})$$

$$H(z, \theta) = \begin{bmatrix} H_1(z_1, z_2, \dots, z_n, \theta_1, \theta_2, \dots, \theta_k) \\ \vdots \\ H_n(z_1, z_2, \dots, z_n, \theta_1, \theta_2, \dots, \theta_k) \end{bmatrix}$$

Assume.

$z$  is random

pdf of  $z = f(z)$

$$G(\theta) = \mathbb{E}_z [H(z, \theta)] = 0 \quad G(\theta) = \int H(z, \theta) \underbrace{f(z)}_{\text{we don't know}} dz$$

We have to know data.  $z_1, z_2, \dots$

$$V_t = V_{t-1} - \mu \underbrace{\mathbb{E}_z [H(z, V_{t-1})]}_{G(V_{t-1})} \rightarrow \theta_t = \theta_{t-1} - \mu H(\underbrace{z_t}_{\text{data point}}, \theta_{t-1})$$

$$G(\theta) = 0 \quad \text{deterministic}$$

random ( $z_t$  is random)

$$\theta_t = \bar{\theta}_t + v_t$$

$$\bar{\theta}_t = \bar{\theta}_{t-1} - \mu \overline{H(z_t, \theta_{t-1})} \rightarrow \bar{\theta}_t = \bar{\theta}_{t-1} - \overline{\mu H(z_t, \underline{\bar{\theta}_{t-1}})}$$

not the same, but only  
cause an error of  $O(\mu^2)$

$$\theta_t \rightarrow \theta_\infty = \underline{\theta_\infty} + \underline{v_\infty}$$

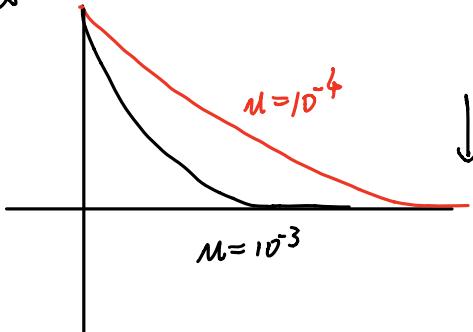
solution to  
 $G(\theta) = 0$

small bias  
we want to make  
this small.

$$\theta_t = \theta_{t-1} - \mu H(z_t, \theta_{t-1})$$

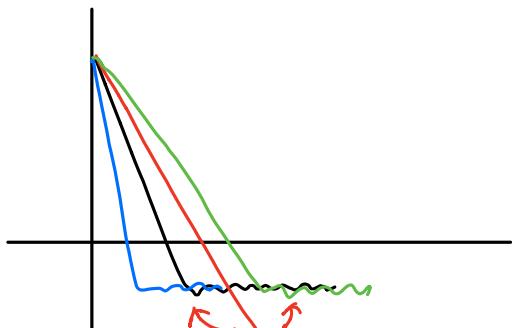
$$v_\infty$$

$$\|\theta_t - \theta_\infty\|^2$$



↓ converge

← we can not evaluate Alg.  
by this.



$$\Delta_1 = \mu_1 C_1 + O(\mu^2)$$

$$\Delta_2 \quad \mu'_1 \cdot C_1 = \Delta_2 \quad \mu'_1 = \frac{\Delta_2}{C_1} = \frac{\Delta_2}{\Delta_1} \mu_1$$

blue better than BLACK.

green worse than BLACK.

$$\theta_t = \theta_{t-1} - \mu H(z_t, \theta_{t-1})$$

$$\underbrace{\mathbb{E}_z [H(z, \theta)]}_G(\theta) = 0$$

→

we are trying to solve this.

$$h(z, \theta) \quad g(\theta) = \mathbb{E}_z [h(z, \theta)] \quad \text{xpdf} \quad \checkmark \text{data: } z_1, z_2, \dots$$

$$\theta_t = \theta_{t-1} - \underline{\mu} \nabla_{\theta} h(z_t, \theta_{t-1})$$

**SG**      **learning rate**  
 (stochastic gradient)      (step size)

$$\nabla_{\theta} g(\theta) = 0$$

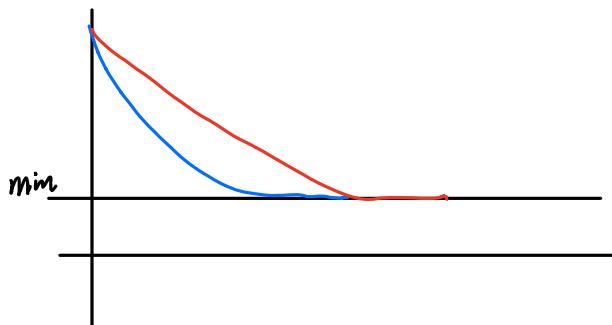
$$G(\theta) = \mathbb{E}_z [\underbrace{\nabla_{\theta} h(z, \theta)}_{H(z, \theta)}]$$

$$G(\theta) = \int H(z, \theta) f(z) dz$$

$$\theta_t^i = \bar{\theta}_t^i + v_t^i \quad g(\theta_t^i)$$

$\hookrightarrow v_\infty$

$$g(\theta_t^i) = g(v_\infty) + (\nabla_{\theta} g(v_\infty))^T v_t + \frac{1}{2} v_t^T (\nabla_{\theta}^2 g(v_\infty)) v_t$$



Avg(blue-min) vs Avg(red-min)  
 we need make sure  $m$  is the same.  
 time step).

## LINEAR REGRESSION.

$$y_t = \theta_*^T x_t + w_t$$

$w_t$

$x_t \xrightarrow[\theta_*]{} \theta^T x_t \oplus y_t$

$\{w_t\}$  is independent from  $\{x_t\}$   
 need to estimate  $\theta_*$ .

$$g(\theta) = \mathbb{E}_z [(y - \theta^T x)^2]$$

$$= \mathbb{E} [( \theta_* - \theta )^T X + w ]^2$$

$$z = \begin{bmatrix} y \\ x \end{bmatrix}$$

$$(\theta_* - \theta)^T \overline{X X^T} (\theta_* - \theta) + 2(\theta_* - \theta)^T \overline{X w} + \overline{w^2}$$

$\sigma_w^2$

$$g(\theta) = (\theta_* - \theta)^T R_x (\theta_* - \theta) + \underline{\sigma_w^2} \text{ don't know}$$

data:  $(y_1, x_1), (y_2, x_2) \dots$

$$\theta_t = \theta_{t-1} - \mu (y_t - \theta_{t-1}^T x_t) x_t$$

LMS

$$g(\theta) = \mathbb{E}_z [ \underbrace{(y - \theta^T x)^2}_{\downarrow} ]$$

$$h(z, \theta) = z = \begin{bmatrix} y \\ x \end{bmatrix}$$

$$\mathbb{E} [(y_t - \theta^T x_t) x_t] = 0$$

$$\theta_t^2 = \theta_{t-1}^2 - \mu \text{sign}(y_t - (\theta_{t-1}^2)^T x_t) x_t$$

$$\mathbb{E} [\text{sign}(y_t - \theta_*^T x_t) \cdot x_t] = 0$$

$$\mathbb{E} [\text{sign}(w_t) x_t] = \mathbb{E} [\text{sign}(w_*)] \cdot \mathbb{E} [x_t]$$

$$\theta_t^3 = \theta_{t-1}^3 - \mu (y_t - (\theta_{t-1}^3)^T x_t) \text{sign}(x_t)$$

$$g(\theta) = \mathbb{E}_z [h(z, \theta)] \quad h(z, \theta), p(z, \theta)$$

$$g(\theta) = (\underbrace{\mathbb{E}_z [h(z, \theta)]}_{\alpha})^2 + (\underbrace{\mathbb{E}_z [p(z, \theta)]}_{\beta})^2 \quad \text{Tensorflow can not solve}$$

$$\nabla_\theta g(\theta) = 2\alpha \mathbb{E}_z [\nabla_\theta h(z, \theta)] + 2\beta \mathbb{E}_z [\nabla_\theta p(z, \theta)] = 0$$

$$\mathbb{E}_z [\alpha \nabla_\theta h(z, \theta) + \beta \nabla_\theta p(z, \theta)] = 0$$

$$\alpha - \mathbb{E}_z [h(z, \theta)] = 0$$

$$\mathbb{E}_z [\alpha - h(z, \theta)] = 0$$

$$\mathbb{E}_z [\beta - p(z, \theta)] = 0$$

$$\mathbb{E}_z [H(z, \theta)] = 0$$

self-defined Alg.  
to solve this.

$$\theta_t = \theta_{t-1} - \mu H(z_t, \theta_{t-1})$$

$$\begin{bmatrix} \theta \\ \alpha \\ \beta \end{bmatrix}_t = \begin{bmatrix} \theta \\ \alpha \\ \beta \end{bmatrix}_{t-1} - \mu \begin{bmatrix} \alpha_{t-1} \nabla_\theta h(z_t, \theta_{t-1}) + \beta_{t-1} p(z_t, \theta_{t-1}) \\ \alpha_{t-1} - h(z_t, \theta_{t-1}) \\ \beta_{t-1} - p(z_t, \theta_{t-1}) \end{bmatrix}$$

$$v_t = v_{t-1} - \mu \bar{E}_z [H(z, v_{t-1})]$$

A:  $\theta_t = \theta_{t-1} - \mu \frac{[H(z_t, \theta_{t-1}) + H(z_{t-1}, \theta_{t-1}) + \dots + H(z_{t-m}, \theta_{t-1})]}{m+1}$

B:  $\theta_t = \theta_{t-1} - \mu H(z_t, \theta_{t-1})$

Converge almost the same.

A will update every  $m+1$  steps.

B will update every step.

Even if A converge faster, it will update slower.

By using GPU, A is better.

$$G(\theta) = 0$$

$$G(\theta) = \begin{bmatrix} G_1(\theta_1, \dots, \theta_K) \\ G_2(\theta_1, \dots, \theta_K) \\ \vdots \\ G_K(\theta_1, \dots, \theta_K) \end{bmatrix} = 0$$

solvble  
algorithm.  
random

$$G(\theta) = \bar{E}_z [H(z, \theta)] \text{ but have data.}$$

$$\theta_t = \theta_{t-1} - \mu H(z_t, \theta_{t-1})$$

$$\theta_t = \bar{\theta}_t + v_t$$

$$G(\bar{\theta}_{t-1}) \rightarrow \theta_\infty$$

$$\bar{\theta}_t = \bar{\theta}_{t-1} - \mu \bar{E}_z [H(z_t, \bar{\theta}_{t-1})] + O(t^2)$$

$$\theta_t = \theta_{t-1} - \mu G(\theta_{t-1})$$

$$\{\theta_t\} \rightarrow \theta_\infty \text{ continue}$$

$G(\theta) = 0$   
solution for this.

small error  
not important

mean + noise

$$\min g(\theta)$$

$$\nabla_{\theta} g(\theta) = 0$$

$$\theta_t = \theta_{t-1} - \mu \nabla_{\theta} g(\theta_{t-1})$$

$$\mathbb{E} [\|\mu_{\theta}\|^2] \sim \mu$$

$\mu$  too small

converge too slow.

$$g(\theta) = \mathbb{E}_z [h(z, \theta)] \text{ don't know distribution}$$

but have data.

solvable algorithm,  $\theta_t = \theta_{t-1} - \mu \nabla_{\theta} h(z_t, \theta_{t-1})$

$$\theta_t = \theta_0 + v_t$$

$$\theta_t \rightarrow \theta_\infty$$

$\mu$  can reduce noise.

but slow down convergence.

$$\Omega(\theta, w), \quad \Omega(\theta, 0) = \theta$$

$$\text{if } \Omega(\theta, w) = \theta + w$$

$$\theta_t = \Omega(\theta_{t-1}, \mu G(\theta_{t-1}))$$

$$\theta_t = \theta_{t-1} - \mu G(\theta_{t-1})$$

$$\approx \Omega(\theta_{t-1}, 0) + (\underbrace{J_w \Omega(\theta_{t-1}, 0)}_{\text{Jacobi matrix}})^T \mu G(\theta_{t-1})$$

$$= \theta_{t-1} + \mu (J_w \Omega(\theta_{t-1}))^T G(\theta_{t-1})$$

$$(J_w \Omega(\theta, 0))^T G(\theta) = 0$$

$$G(\theta) = \mathbb{E}_z [h(z, \theta)]$$

$$\theta_t = \Omega(\theta_{t-1}, \mu H(z_t, \theta_{t-1}))$$

$$\theta_t = \theta_{t-1} + \mu H(z_t, \theta_{t-1})$$

PCA.

$$R = \mathbb{E} [z z^T] \quad z, z_2, \dots$$

$$R \longrightarrow$$

$$R_t = \frac{z_1 z_1^T + \dots + z_t z_t^T}{t}$$

$$u_0 = \frac{R u_0}{\|R u_0\|}$$

$$R_t = R_{t-1} + z_t z_t^T = R_{t-1} + \mu (z_t z_t^T - R_{t-1})$$

$$u_t = \frac{R_t u_{t-1}}{\|R_t u_{t-1}\|} u_{t-1}$$

$$u_t = \frac{R_t \cdot u_{t-1}}{\|R_t u_{t-1}\|} \quad \text{Complexity} = k^2$$

$$\underline{M \cdot z_t z_t^T + (1-\mu) z_{t-1} z_{t-1}^T + \dots + (1-\mu)^k z_0 z_0^T}$$

evolve by time.

$$u_t = \frac{R u_{t-1}}{\|R u_{t-1}\|}$$

$$v_t = \frac{(I + \mu R) v_{t-1}}{\|(I + \mu R) v_{t-1}\|}$$

$u$  finds eigen largest vector  $R$ .

$v$  ...  $\dots$   $\dots$   $\dots$   $I + \mu R$ .

$$\lambda u = Ru$$

$$(I + \mu R)u = (I + \mu \lambda)u.$$

Converge to same thing.

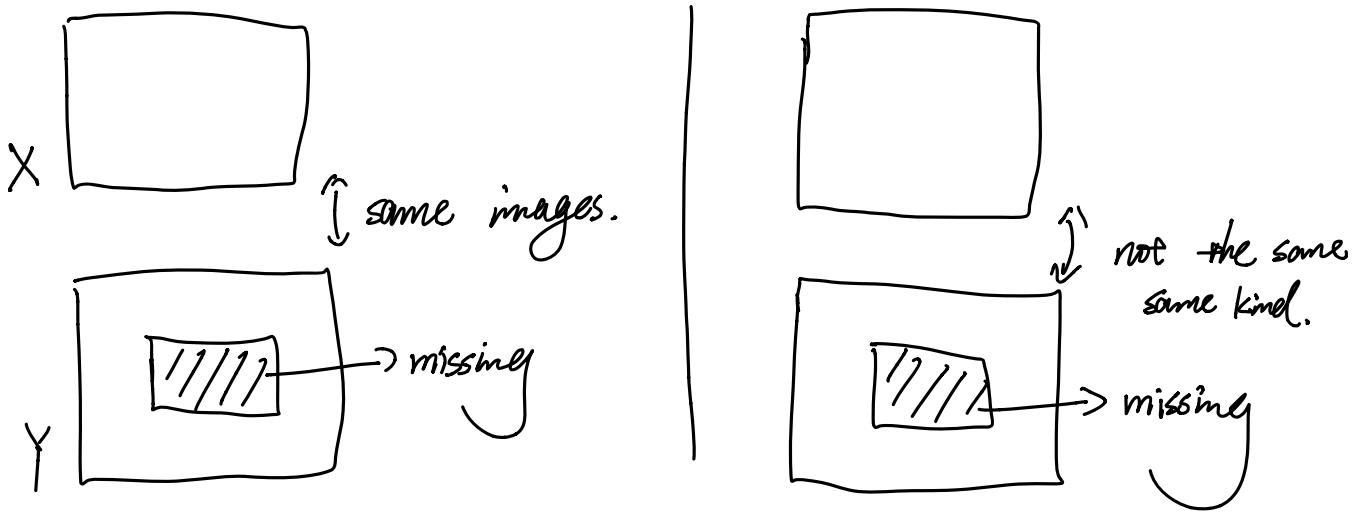
$$\Omega(\theta, w) \rightarrow \Omega(J, w) = \frac{(I + w)J}{\|(I + w)J\|}$$

$$R = \bar{E}[z z^T]$$

$$J_t = \frac{(I + \mu z_t z_t^T) J + 1}{\|(I + \mu z_t z_t^T) J + 1\|} \quad O(h) \quad \text{find largest eigen vector}$$

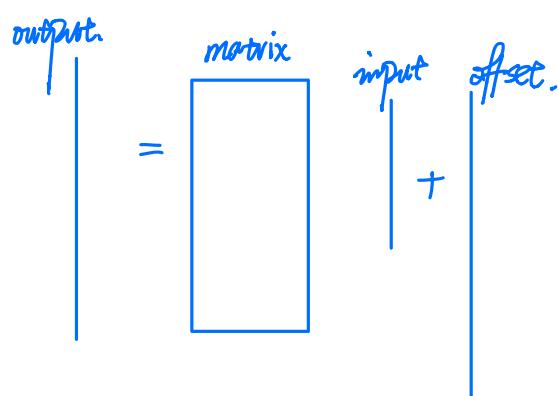
$$v_t = \frac{(I - \mu R) v_{t-1}}{\|(I - \mu R) v_{t-1}\|}$$

$$u_t = \frac{R^{-1} v_{t-1}}{\|R^{-1} v_{t-1}\|} \quad \text{find smallest eigen vector.}$$



$$\begin{aligned}
 & X \quad (Y, x) \quad \hat{Y} = \phi(x) \\
 \downarrow & Y = \phi(x) \quad \mathbb{E}[C(Y, \hat{Y})] \text{ loss case.} \\
 & Y \quad = \mathbb{E}_{\hat{Y}}[C(Y, \phi(x))] \\
 & \min_{\phi(x)} \mathbb{E}_{\hat{Y}}[C(Y, \phi(x))] \quad \text{Don't know distribution we have data.}
 \end{aligned}$$

$\phi(x) = \phi(x, \theta)$  collection of all parameters.  
output of neural network.



$$\min_{\theta} \mathbb{E}_{\hat{Y}}[C(Y, \phi(x, \theta))]$$

$$\theta_t = \theta_{t-1} - \mu \nabla_{\theta} C(Y_t, \phi(X_t, \theta_{t-1}))$$

along layers  
↓  
Back propagations.

Assume it converge to  $\theta^*$ .

last point  $\hat{Y} = \phi(X, \theta^*)$   
lost image

classification.

data.  $X_1, X_2 \dots$   
table  $y_1, y_2 \dots$

$\hat{y} = \phi(x)$   
↓  
(able. new  $x$ )

$$\sum_{i=1}^k (1 - \phi(x_i))^2 + \sum_{j=1}^k (1 + \phi(x_j))^2$$

error for  $x_i$       error for  $x_j$

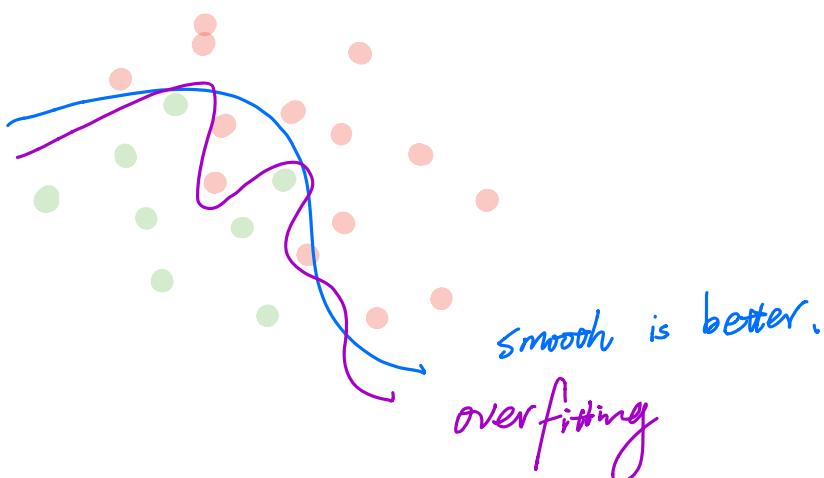
$$\bar{E}_1 [(1 - \phi(x^1))^2] + \bar{E}_2 [(1 + \phi(x^2))^2]$$

$$\bar{E}_1 [(1 - \phi(x^1, \theta))^2] + \bar{E}_2 [(1 + \phi(x^2, \theta))^2] \quad \text{total average cost}$$

$$\theta_t = \theta_{t-1} - \mu \left\{ \nabla_{\theta} (1 - \phi(x_t^1, \theta_{t-1}))^2 + \nabla_{\theta} (1 + \phi(x_t^2, \theta_{t-1}))^2 \right\}.$$

class1  $(\underline{y_1^1}, \underline{x_1^1}), (\underline{y_2^1}, \underline{x_2^1})$   
class2  $(\underline{y_1^2}, \underline{x_1^2}), (\underline{y_2^2}, \underline{x_2^2})$

$\phi(x, \theta_*)$       boundary  $\rightarrow \phi(x, \theta_*) = 0$   
table      input data



GANS.

$$f(x) \\ h(z) \rightarrow z$$

$$\underline{G(z)} = Y$$

*non-linear transformation.*

$$x \quad f(x) \\ z \quad h(z) \\ G(z) = Y \rightarrow g(Y)$$

random  $\downarrow$   
also random.

scalar function  $D(x)$

$$\bar{E}_x[D(x)] - \bar{E}_y[D(Y)] = 0$$

*D is arbitrary.*

$$\max_{D(x)} \{ \bar{E}_x[D(x)] - \bar{E}_y[D(Y)] \}$$

$$= \max_{D(x)} \{ \bar{E}_x[D(x)] - \bar{E}_z[D(G(z))] \}.$$

$$\min_{G(z)} \left\{ \max_{D(x)} \{ \bar{E}_x[D(x)] - \bar{E}_z[D(G(z))] \} \right\}$$

$$\phi(z) \quad \psi(z)$$

$$\min_G \max_D \{ \bar{E}_x[\phi(D(x))] + \bar{E}_z[\psi(D(G(z)))] \}$$

$$\phi(z) = \log(1-z)$$

$$\psi(z) = \log(z)$$

$$G(z) \sim G(z, \theta)$$

$$D(x) \sim D(x, \theta)$$

$$\nu_t = \nu_{t-1} + \mu \nabla_\nu X_t z_t \quad \nu_t \rightarrow \nu_\infty$$

$$\theta_t = \theta_{t-1} - \mu \nabla_\theta X_t z_t \quad \theta_t \rightarrow \theta_\infty$$