

# ADVANCED MACHINE LEARNING

9/11/23

→ well known conferences for ML

- ICML
- Neurips
- ICLR
- AAAI
- CVPR
- ACL

↑  
Reproduce the results in papers.

## ① A Statistical Learning framework.

Domain set -  $x$

Label set -  $y$

Training data:  $S = ((x_1, y_1), (x_2, y_2), \dots, (x_m, y_m))$ .

from  $x \times y$ .

Learner's output:  $h: x \rightarrow y$ .

(Source: 2nd book -

Learning from Data)

## ② Data generation model

prob. dist. over  $x$  is  $D$ .

Labelling func<sup>2</sup>:  $f: x \rightarrow y$ .

## ③ Measures of success.

$$(over \atop over) h_{\text{dif}}(h) = P_{x \sim D} [h(x) \neq f(x)] \rightarrow \text{loss function}$$

↓  
over  
entire  
distribution

learn  $\rightarrow h$   
fixed  $\rightarrow f$

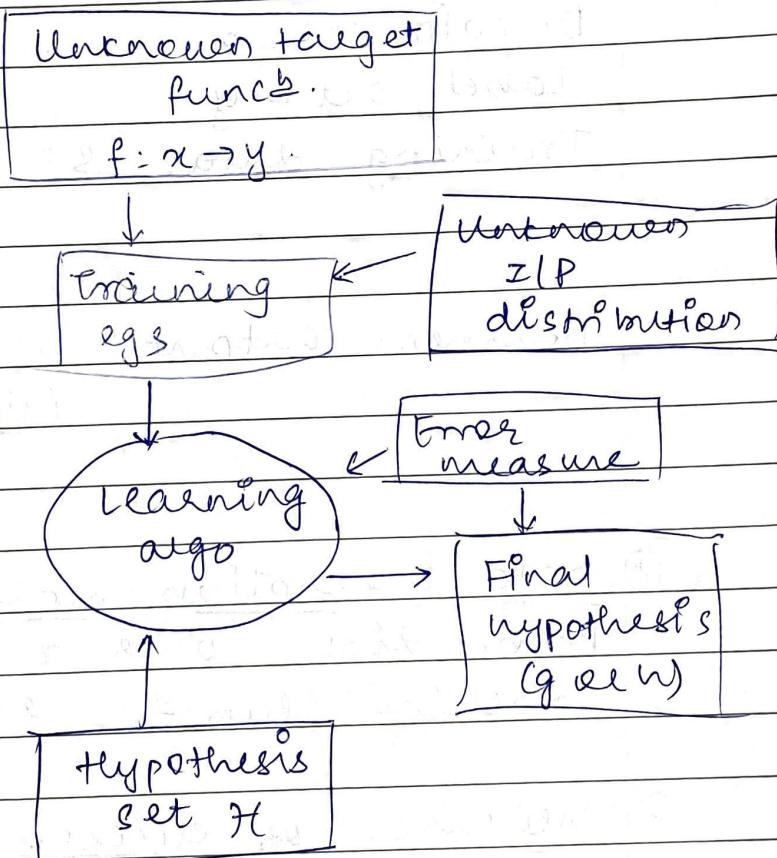
④ ERM framework. (Empirical Risk Minimization)

Training error  $\rightarrow$   
 (empirical error)  $L_s(h) = \frac{|\{i \in [m] : h(x_i) \neq y_i\}|}{m}$

Search for a sol<sup>n</sup> that works well on available data

$\rightarrow$  I.I.D assumption

$\rightarrow$  Goal is that func<sup>t</sup> performs well on test data.



16/11/23

Lec: 2

- / -

$$E_{in}(h) = \frac{1}{N} \sum_{i=1}^N [h(x_i) \neq f(x_i)]$$

$$E_{out}(h) = P[h(x) \neq f(x)]$$

① Märklee Experiment

A bin with red & green marbles.

Prob. of picking red  $\rightarrow \mu$

" green  $\rightarrow 1-\mu$

Problem  $\rightarrow \mu$  is unknown

N - total no. of marbles.

v - ~~prob.~~ proportion of red marbles out of N

$$\mu = \frac{v}{N}$$

If 'N' is large enough,  
then  $\mu$  is near to  
empirical mean

② Hoeffding's Inequality  
type of Chernoff Bound.

$$P[|v - \mu| > \epsilon] \leq 2e^{-2\epsilon^2 N}, \epsilon > 0.$$

$\uparrow \quad \downarrow$   
empirical mean original mean

$\rightarrow$  e.g. If  $\epsilon = 0.1$ , what is no. of samples?

$\epsilon$  - failure probability

$\rightarrow$  Relate this to learning.

For fixed hypothesis,

$$P[|E_{in}(h) - E_{out}(h)| > \epsilon] \leq 2e^{-2\epsilon^2 N}, \epsilon > 0$$

(h=hypothesis)

## ④ Union Bound

Events  $\rightarrow A \& B$

$$P(A \cup B) \leq P(A) + P(B)$$

$\rightarrow$  For a fixed set of hypothesis of size  $M$

$$P[|E_{in}(h) - E_{out}(g)| > \epsilon] \leq 2M e^{-2\epsilon^2 N}, \epsilon > 0.$$

(derive)

with probability atleast  $1 - \delta$

$$E_{out}(g) \leq E_{in}(g) + \sqrt{\frac{1}{2N} \ln \frac{2M}{\delta}}$$

Generalization Bound

$$\text{let } 2M e^{-2\epsilon^2 N} = \delta.$$

$$e^{-2\epsilon^2 N} = \frac{1}{2M} \delta.$$

$$-2\epsilon^2 N = \frac{\ln(\delta)}{2M}$$

$$\therefore \delta = \sqrt{\frac{1}{2N} \ln \left( \frac{2M}{\delta} \right)}$$

$$-2\epsilon^2 N = \ln \frac{\delta}{2M}$$

$$N = \frac{1}{2\epsilon^2} \ln \left( \frac{2M}{\delta} \right)$$

$\mathcal{H}$ - hypothesis space

④ PAC learning  $\rightarrow$  Probably Approximately Correct

It is PAC learnable if

$$- m_{\mathcal{H}}: (0, 1)^2 \rightarrow N$$

$$- \epsilon, \gamma \in (0, 1), D \text{ over } X \times Y$$

$\rightarrow$  when running the learning algo.

④ Define

Dichotomies:

$$\mathcal{H}(x_1, \dots, x_N) = \{ h(x_1), \dots, h(x_N) \mid h \in \mathcal{H} \}$$

10/11/23

with probability atleast  $1 - \delta$

$$E_{out}(g) \leq E_{in}(g) + \sqrt{\frac{1}{2N} \ln \frac{2M}{\delta}} \quad \begin{matrix} \leftarrow \text{generalization} \\ \text{bound} \end{matrix}$$

$\rightarrow$  What if  $M = \infty$ ?

④ Growth function - for a hypothesis set  $\mathcal{H}$  is

$$m_{\mathcal{H}}(N) = \max_{x_1, \dots, x_N} |\{h \in \mathcal{H} \mid h(x_1, \dots, x_N)\}|$$

$$\therefore m_{\mathcal{H}}(N) \leq 2^N$$

④ Shattering: If  $m_{\mathcal{H}}(N) = 2^N$ . (hypothesis set can generate all possible dichotomies)

④ VC-Dimension: of a hypothesis set  $\mathcal{H}$ , denoted by  $d_{VC}(\mathcal{H})$  or  $d_V$  is the largest value of  $N$  for which  $m_{\mathcal{H}}(N) = 2^N$

# Learning Theory

111

→ If  $m_H(N) = 2^N$  for  $\#H$ ,  $dvc(H) = \infty$

$$VC(\text{hyperplanes}) = d + 1$$

⇒ Sauer's Lemma

→ VC-Generalization bound

- more complex hypothesis - doesn't generalize well

Evaluation

① Reproducibility challenge - write reproducibility report

② Project

paperswithcode.com/

rc 2022

Newips 2021/22

ICML 2021/22

AAAI 2021/22

ACM-FACCT 2021/22

reflect of  
8-10 pages

choose any 1 paper

- study algorithmic

- choose papers with more

① refl. paper project empirical implement  
atian

- Replicate results

(on same datasets)

- can you do something more?

◦ Tweak algo

◦ More dataset

◦ Hyperparameter tunings

- Replicate based on library changing

→ git repositoris  
for this  
distill. pub

20/1/23

## Bias Variance tradeoff and PLA

VC Generalization Bound.

$$E_{\text{out}}(g) \leq E_{\text{in}}(g) + \sqrt{\frac{8}{N} \ln \frac{4m_H(2N)}{8}}$$

The VC bound is loose.

- Hoeffding's Inequality Woes:-

$$m_H(N) \leq \sum_{i=1}^{d_{VC}} \binom{N}{i}$$

$$m_H(N) \leq N^{d_{VC}} + 1$$

→ Hoeffding's Inequality has some slack.

$m_H(N)$  gives a worst case estimate

→ Bounding  $m_H(N)$  by polynomial of order

→ (Not easy to get VC dimension of decision tree)

More VC dimension = more complex

→ Rule of thumb → keep training data

10 times the VC dimension

23/1/23

— / —

## ⑩ Predictive & Generative Models

Bias & Variance:-

$$E_{\text{out}}(g^{(D)}) = E_x[(g^{(D)}(x) - f(x))^2]$$

$E_x$  denotes expected value w.r.t.  $x$

Imp.:  $g^{(D)}$  is dependent on  $(D)$  (data)

$$E_D [E_{\text{out}}(g^{(D)})] \quad (\text{Out sample error})$$

$$= E_D [E_x [(g^{(D)}(x) - f(x))^2]] \quad \swarrow (a-b)^2$$

$$= E_x [E_D [(g^{(D)}(x) - f(x))^2]]$$

$$= E_x [E_D [g^{(D)}(x)^2] - 2E_D [g^{(D)}(x)] f(x) + f(x)^2]$$

$$E_D(g^D(x)) = \bar{g}(x)$$

$$E_x [E_D(g^D(x)^2) - 2\bar{g}(x)f(x) + f(x)^2]$$

$$= E_x [E_D(g^D(x)^2) - \bar{g}(x)^2 + \bar{g}(x)^2]$$

$$\quad \quad \quad - 2\bar{g}(x)f(x) + f(x)^2 \\ \quad \quad \quad (\bar{g}(x) - f(x))^2$$

$$\therefore = E_D[(g^{(D)}(x) - \bar{g}(x))^2] + (\bar{g}(x) - f(x))^2$$

$$= \text{Var}(x) + \text{bias}$$

K/

— / —

Idea - we can break down our error into bias & variance.

→ It also has irreducible error, but we assume here that data is noise-free.

$$\therefore E_D [E_{out}(g^{(D)})] = E_X [\text{bias}(x) + \text{var}(x)]$$

Here, along with H algorithm, A also matters.

→ Total least square (orthogonal least square)  
    ↳ tries to reduce the L<sub>2</sub> distance between points & line

→ OLS - minimize the vertical offset.

### ④ PLA algorithm (Perceptron Learning Algorithm)

$$w^{(0)} = (0, 0, 0, 0, 0) \quad (\text{initial})$$

while there is a misclassified point

$$x(t), y(t)$$

$$w(t+1) = w(t) + y(t) \cdot x(t)$$

$$y(n) \\ = \text{sgn}(w^T x)$$

- If a point is classified correctly, then,  $((w^T x)(y)) > 0$ .

→ If data is linearly separable, PLA gives 0 errors.

11

27/1/23.

Convex Optimization in ML }  
Primal Dual Formulations } SVM  
KKT conditions } problem

## Optimization for ML

continuous  
Non convex  
↓  
Convex

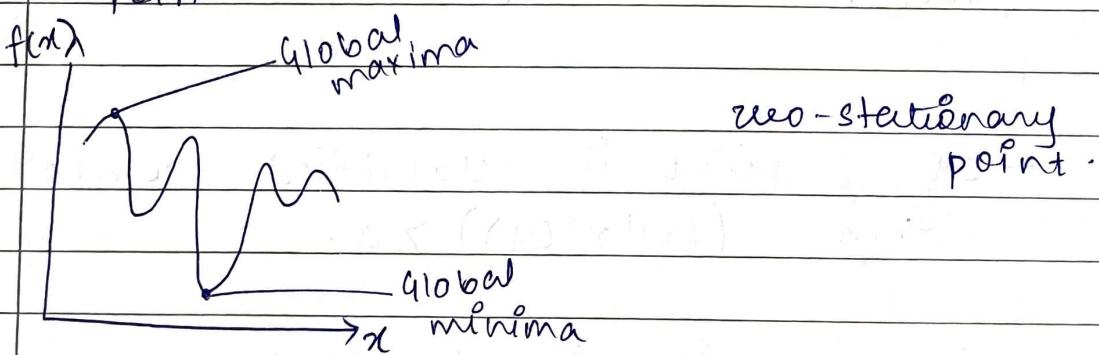
Combinational  
Submodularity  
Subset Selection

Gradient Descent  $\rightarrow$  move in direction of -ve gradient

[Optimizations for ML - Nishanth Veshma]

→ In convex  $f(x)$ , there is high chance to reach global minima with more no. of iterations.

→ In practice, we use SGD in non-convex func<sup>n</sup>, which can also get stuck at local minima or saddle point.



$$f: \mathbb{R}^m \rightarrow \mathbb{R}$$

J.

J

— / —

Unconstrained  
 $\min f(x)$

Constrained  
 $\min f(x)$

such that

$$x \in D = \mathbb{R}^m + \dots$$

s.t.

$$x \in \mathbb{R}^m$$

Linear regression:-

$$\min_{\substack{x \in \mathbb{R}^m \\ x \neq 0}} \|Ax - b\|^2 \quad \leftarrow \text{unconstrained optimization}$$

Ridge regression:-

$$\min_{\substack{x \in \mathbb{R}^m \\ x \neq 0}} \|Ax - b\|^2 + \lambda \|x\|_2^2 \quad \text{is diag}$$

Lasso regression:-

$$\min_{\substack{x \in \mathbb{R}^m \\ x \neq 0}} \|Ax - b\|_2^2 + \lambda \|x\|_1$$

shape?

Least Absolute Shrinkage & Selection Operator

Cond? :-

i) Let  $f: \mathbb{R}^m \rightarrow \mathbb{R}$  continuous & differentiable

ii) For local minima

$$\nabla f(x) = \begin{bmatrix} \frac{\partial f(x)}{\partial x_1} \\ \vdots \\ \frac{\partial f(x)}{\partial x_n} \end{bmatrix} \quad \text{eg. } x_1^2 + 2x_2$$

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

Positive definite  $\rightarrow ?$

II      semi-II     $\rightarrow ?$

Hessian matrix - matrix of partial derivatives

6/2/23

U.G.

1/1

→ Eg. of unconstrained optimization problem → Linear regression

→ Eg. of constrained → Lasso, Ridge

constraints

$$\|x\|_1$$

$$\|x\|_2^2$$

$H(x)$

↑

positive

$$x^T A x$$

→ gives scalar

$$x^T A x > 0, \forall x \neq 0$$

positive semi-definite

→ If all eigen values are true then.

$$Ax = \lambda x$$

→ Convex set -  $X \subseteq \mathbb{R}^n$  is said to be convex

if  $\forall x, y \in X$

$$\lambda x + (1-\lambda) y \in X, \forall 0 \leq \lambda \leq 1$$

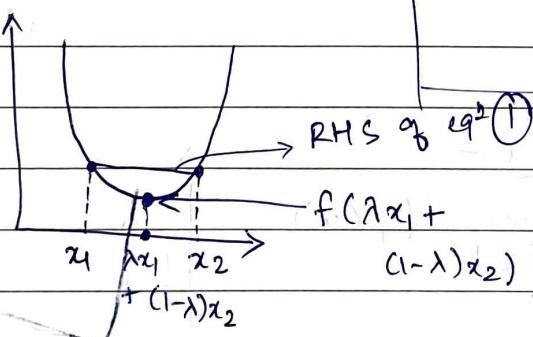
→ Convex function:- If  $f: \mathbb{R}^n \rightarrow \mathbb{R}$

(def<sup>2</sup>)

$$\mathbb{R}^n \subseteq X$$

$$\lambda x_1, x_2$$

$$f(\lambda x_1 + (1-\lambda) x_2) \leq \lambda f(x_1) + (1-\lambda) f(x_2)$$



∴ f is convex func

If in eq<sup>2</sup> ①, ' $<$ '

sign  $\rightarrow$  strictly convex

max (0,  $\infty$ )

$\|x\|_1$

→ If  $f(x)$  is convex  $\Leftrightarrow H(x)$  is P.S.D.

$\forall x \in X$

→  $f(\cdot)$  is strictly convex

$\Leftrightarrow H(x)$  is PD  $\forall x \in X$ .

→ If func<sup>2</sup> is convex, local minima = global minima

Eg:

①  $w^T x + b$  on  $R^n$  - convex

② -ve log likelihood func<sup>2</sup> - "

③ quadratice func<sup>2</sup>  $\Rightarrow f(x) = \frac{1}{2} x^T A x + b^T x + c$  en  $R^n$   
(A is PSD).

→ check if  $\max: f(x) = \max(x_1, \dots, x_n)$   
is convex?

### ④ Operations preserving convexity

①  $f(x)$  is convex

$\Rightarrow \alpha f(x)$  is also convex,  $\alpha > 0$

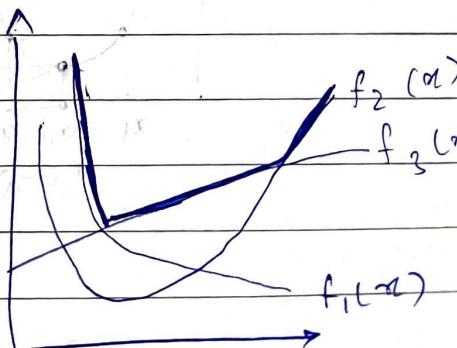
② Sum:-  $f_1(x) + f_2(x) = f(x)$

If  $f_1(x)$  &  $f_2(x)$  are convex,  
 $f(x)$  is also convex.

### ③ Pointwise max:-

$$f(x) = \max_{1 \leq i \leq k} f_i(x)$$

convex



smooth func?

— / / —

④ log sum exponential

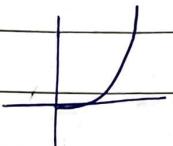
$$f(x) = \log \sum_{i=1}^n \exp(x_i)$$

⑤ convex functions by composition

$h(x) \rightarrow h$  is convex

case i)  $g(h(x))$

convex + non decreasing



case ii)  $g(h(x))$

concave

convex + non increasing



8/2/23

AML

→ Unconstrained optimization

→ Lagrange multiplier

## ④ constrained optimization problem.

$$f_0: \mathbb{R}^n \rightarrow \mathbb{R}$$

$$\min f_0(x)$$

$$\text{such that } f_i(x) \leq 0 \quad i \in [1, m]$$

$$g_j(x) = 0 \quad j \in [1, l]$$

no. of constraints

$$f_0(x) + \sum_{i=1}^m \lambda_i f_i(x) + \sum_{j=1}^l \mu_j g_j(x)$$

Lagrange  
multipliers  
problem

$$\|Ax - b\|_2^2 + \lambda \|x\|_2^2, \quad \lambda > 0$$

$$\|Ax - b\|_2^2 \text{ s.t. } \|x\|_2^2 \leq R$$

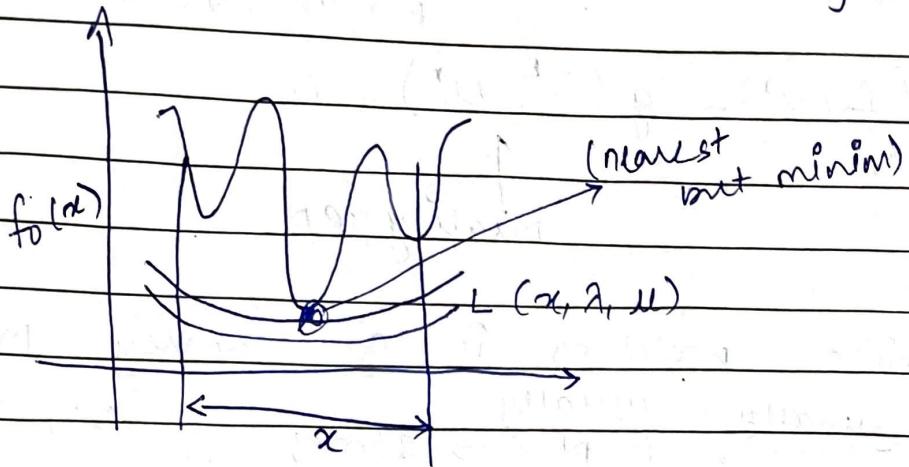
→ where  $\lambda_i, \mu_j$  are called Lagrangian multipliers.

$$L(x, \lambda, \mu) = f_0(x) + \sum_{i=1}^m \lambda_i f_i(x) + \sum_{j=1}^l \mu_j g_j(x)$$

→ what if  $x$  is a feasible point?

↓  
means it  
satisfies  
constraints)

For feasible  $x$ ,  $L(x, \lambda, \mu) \leq f_0(x)$



$$\begin{bmatrix} x_1 \\ x_2 \end{bmatrix} = \begin{bmatrix} 2 \\ 3 \end{bmatrix}$$

$$L(\lambda, \mu) = 2\lambda_1 u_1^2 + 3\lambda_2 u_2^2$$

$$\left. \begin{array}{l} \min L(x, \lambda, \mu) \text{ such that } \lambda > 0 \\ \text{lower bound on optimal solution} \end{array} \right\} \leq f_0(x^*)$$

$$\max_{\lambda, \mu} \left[ \min_x L(x, \lambda, \mu) \right] \leq f_0(x^*)$$

$(\lambda, \mu) \in \text{Domain s.t. } \lambda > 0$

$$\left[ \min_x L(x, \lambda, \mu) \right] = g(\lambda, \mu)$$

dual function

→ Primal problem & Dual problem

$$\max_{(\lambda, \mu)} g(\lambda, \mu)$$

$\lambda > 0$

$$g(\lambda, \mu) = \min_x L(x, \lambda, \mu)$$

weak duality theorem  $\rightarrow g(\lambda, \mu) \leq [g(\lambda^*, \mu^*) \leq f_0(x^*)]$

## Weak duality theorem - statement

Date \_\_\_\_\_  
Page \_\_\_\_\_

$f_0(x^*)$  - optimal value of  
primal function

$$(f_0(x^*) - g(\lambda^*, u^*)) \geq 0$$

↑  
duality gap

→ Primal problem is a convex problem

usually implies → strong duality holds

(not always)

↓  
(Optimum value of dual = value of primal)

$$[f_0(x^*) = 0 ??]$$

constraints  
should  
correctly  
some  
cond

Regularity  
constraint

KKT conditions :-  $x^*$  - point of  
primal optimal

(i) Primal feasibility

$\rightarrow x^*$  is feasible sol<sup>1</sup> for  
primal problem.

$\lambda^*, u^*$  - point of  
dual optimal  
with zero  
duality gap

$$f_i(x^*) \leq 0 \quad i \in [1, m]$$

$$f_j(x^*) \geq 0 \quad j \in [1, l]$$

(ii) Dual feasibility

$$\lambda^* \geq 0$$

(e.g. support vectors)

(iii) Complementary slackness:-

$$\lambda^* f_i(x^*) = 0 \quad \forall i$$

iv) Lagrange optimality(func<sup>2</sup> of x)

$$\nabla_x L(x^*, \lambda^*, \mu^*) = 0$$

→ Any tuple  $(x, \lambda, \mu)$ 

$$\begin{matrix} \uparrow & \uparrow & \uparrow \\ R^n & R^m & R^k \end{matrix}$$

is called a KKT point if it satisfies  
KKT cond.→ When you have zero duality gap,  
 $(x^*, \lambda^*, \mu^*) \Rightarrow$  KKT pointIf something is KKT point, it will  
always be optimal (if func<sup>2</sup> is convex)

## ④ General stat.

If duality gap = 0

 $x^*$  - Primal optimal $\lambda^*$  & dual optimal $\mu^*$

9/3/23

AML

Date  
PageScribe  
notesApplication for SVM:-

Optimizations

book

① Hard Margin SVM

(what is it?)

LDL

→ used in  
NLP for  
topic modelling

what is margin?

why we maximize margin?

Supporting hyperplanes

→ Maximum margin classifier

-1 &amp; +1 labels are used

## → Primal SVM

$$\min_{(w, w_0)} \frac{1}{2} \|w\|_2^2 \text{ such that } y_i(w^T x_i + w_0) \geq 1$$

+ 1, i ∈ [1, n]

sign  $(w^T x_{\text{new}} + w_0)$  → whatever  
the sign,  
assign that label  
(+1 or -1)

② Soft margin SVM → slack

$$(w^T x_i + w_0) \geq 1 - \xi_i \quad (\xi_i)$$

$$1 - \xi_i$$

allows some misclassifications

⇒ Hinge Loss ?

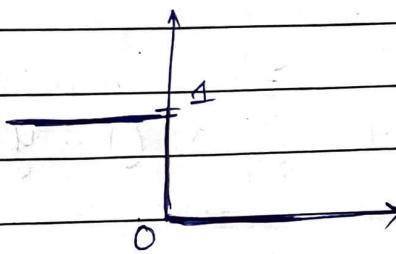
$$y_i^o (w^T x_i^o + w_0) < 1.$$

$$0 < 1 - y_i^o (w^T x_i^o + w_0)$$

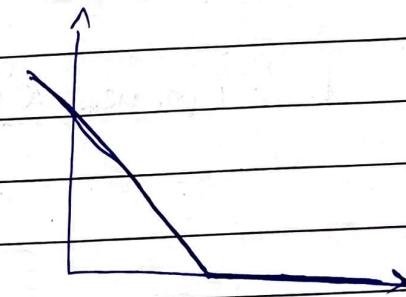
Hinge

$$\text{loss}_i \left[ (x, y), (w, w_0) \right] = \max \left\{ 0, 1 - y_i^o (w^T x_i^o + w_0) \right\}$$

⇒ Zero-one loss



Hinge



⇒ Applying KKT conditions :-

$$\min_{(w, w_0)} \frac{1}{2} \|w\|_2^2 \quad \text{such that } y_i^o (w^T x_i^o + w_0) \geq 1 \quad \forall i, i \in [1, n]$$

(i) Primal Feasibility

$$1 - y_i^o (w^T x_i^o + w_0) \leq 0$$

$w^{**}, w_0^{**}$

primal variables →  $w$  &  $w_0$

$$L(w, w_0, \alpha) = \frac{1}{2} \|w\|_2^2 + \sum_{i=1}^n \alpha_i (1 - y_i^o (w^T x_i^o + w_0))$$

(ii) Dual feasibility

$$\alpha_i^* \geq 0 \quad \forall i$$

↑  
only these are  
support vectors that  
satisfy this condition

(iii) Complementary slackness :-

$$\alpha_i^* (1 - y_i^* (w^T x_i^* + w_0)) = 0$$

(iv) Lagrange optimality :-

$$L(w, w_0, \alpha) = \frac{1}{2} \|w\|_2^2 + \sum_{i=1}^n \alpha_i^* (1 - y_i^* (w^T x_i^* + w_0))$$

$$= \frac{1}{2} w^T w + \sum_{i=1}^n \alpha_i^* (1 - y_i^* (w^T x_i^* + w_0))$$

$$\nabla_w L(w, w_0, \alpha) = 0$$

(differentiating scalar gives vector)

$$\nabla_w L(w, w_0, \alpha) = w$$

$$\sum_{i=1}^n \alpha_i^* - \sum_{i=1}^n \alpha_i^* y_i^* (w^T x_i^* + w_0)$$

$$\nabla_w L(w, w_0, \alpha) = w - \sum_{i=1}^n \alpha_i^* y_i^* x_i^* = 0$$

(partial derivative w.r.t. w)

$$\Rightarrow w = \sum_{i=1}^n \alpha_i^* y_i^* x_i^*$$

cond 1

gives hyperplane in terms of  
coeffs & s.v.s

$$L(w, w_0, \alpha) = \frac{1}{2} w^T w + \sum_{i=1}^n \alpha_i^* (1 - y_i^* (w^T x_i + w_0))$$

$$\nabla_{w_0} L(w, w_0, \alpha) = 0 + 0 - \sum_{i=1}^n \alpha_i^* y_i$$

$$\therefore \sum_{i=1}^n \alpha_i^* y_i = 0 \quad \text{cond } 2$$

B

$\alpha_i^*$  for  $i \in C^+$

$\alpha_j^*$  for  $j \in C^-$

$$\sum_{i \in C^+} \alpha_i^* = \sum_{j \in C^-} \alpha_j^*$$

→ Any e.g. for which dual var ( $\alpha$ ) is non-zero, will lie on supporting hyperplane.

⇒ Defining support vectors:-

$$\{\alpha_i^* \mid \alpha_i^* > 0\}$$

SV are most difficult to classify

→ Getting dual form from primal:-

$$\min_{(w, w_0)} \frac{1}{2} \|w\|_2^2 \quad \text{s.t. } 1 - y_i^* (w^T x_i + w_0) \leq 0 \quad \forall i$$

$$L(w, w_0, \alpha) = \frac{1}{2} \|w\|^2 + \sum_{i=1}^n \alpha_i^* (1 - y_i^* (w^T x_i + w_0))$$

$$g(\lambda, \mu) = g(\alpha) = \min_{w, w_0} L(w, w_0, \alpha)$$

$$w^* = \sum_{i=1}^n \alpha_i^* y_i^* x_i^*, \quad w^* \text{ satisfies}$$

$$\sum \alpha_i^* y_i^* \geq 0$$

$$g(x) = \frac{1}{2} w^{*T} w^* + \sum_{i=1}^n \alpha_i^* - \sum_{i=1}^n \alpha_i^* y_i^* w^{*T} x_i^*$$

$$- \sum_{i=1}^n \alpha_i^* y_i^* w_0^* = 0$$

$$= \frac{1}{2} w^{*T} w^* + \sum_{i=1}^n \alpha_i^* - w^{*T} \sum_{i=1}^n \alpha_i^* y_i^* x_i^* - 0$$

$$= \frac{1}{2} w^{*T} w^* + \cancel{\sum_{i=1}^n \alpha_i^*}$$

$$\sum_{i=1}^n \alpha_i^* - w^{*T} w^*$$

$$= -\frac{1}{2} w^{*T} w^* + \sum_{i=1}^n \alpha_i^*$$

$$= -\frac{1}{2} \left( \sum_{i=1}^n \alpha_i^* y_i^* x_i^* \right)^T \cdot \left( \sum_{j=1}^n \alpha_j^* y_j^* x_j^* \right)$$

$$+ \sum_{i=1}^n \alpha_i^*$$

$$= \max \left[ -\frac{1}{2} \sum_{i=1}^n \sum_{j=1}^n \alpha_i^* \alpha_j^* y_i^* y_j^* x_i^T x_j^* \right]$$

$$+ \sum_{i=1}^n \alpha_i^*$$

Dual form of  
hard margin

s.t.  $\alpha_i^* \geq 0 \quad \forall i$

$$\sum_{i=1}^n \alpha_i^* y_i^* = 0$$

✓ SVM

curve notes end

10/2/23

### ④ Hard Margin SVM

Primal  $\Rightarrow$

$$\min_{(w, w_0)} \frac{1}{2} \|w\|_2^2 \quad \text{s.t. } y_i (w^T x_i + w_0) \geq 1 \quad \forall i$$

Dual  $\Rightarrow$

$$\max \frac{-1}{2} \sum_{i=1}^n \sum_{j=1}^n \alpha_i \alpha_j y_i y_j x_i^T x_j + \sum_{i=1}^n \alpha_i$$

s.t.  $\alpha_i \geq 0 \quad \forall i$   
 $\& \sum_{i=1}^n \alpha_i y_i = 0$

$\rightarrow$  Kernel SVM  $x \in \mathbb{R}^d$   
 $x \xrightarrow{\Phi} z \in \mathbb{R}^n, n > d$

$$\Phi(x) \in \mathbb{R}^n$$

$\uparrow$   
 $\mathbb{R}^d$

$x_i \mapsto \Phi$   
 $z_i = \Phi(x_i)$

$$\text{label}(x_{\text{test}}) = \text{sign} [w^{*T} \Phi(x_{\text{test}}) + w_0^{*}]$$

$$w^* = \sum_{i=1}^m \alpha_i^* y_i \Phi(x_i)$$

$$= \text{sign} \left[ \sum_{i=1}^m \alpha_i^* y_i \Phi(x_i)^T \Phi(x_{\text{test}}) + w_0^* \right]$$

kernel

m data points  
p dimensions

$\Phi$  = feature func<sup>2</sup>

Date \_\_\_\_\_

Page \_\_\_\_\_

Q How to implement SVM in high dimension without actually getting  $\Phi(x_i)$

Mercer's theorem :-

Kernel func<sup>2</sup> :-  $K: \mathcal{X} \times \mathcal{X} \rightarrow \mathbb{R}$

A symmetric function  $K(x_i, x_j) = K(x_j, x_i)$  can be expressed as  $K(x_i, x_j) = \Phi^T(x_i) \Phi(x_j)$  for some  $\Phi$ .

If the matrix,

$$K = \begin{bmatrix} K(x_1, x_1) & K(x_1, x_2) & \dots & K(x_1, x_m) \\ K(x_2, x_1) & \vdots & & \\ \vdots & & & \\ K(x_m, x_1) & \dots & \dots & K(x_m, x_m) \end{bmatrix}_{m \times m}$$

$K$  is Positive definite matrix

Replacing  $\Phi^T(x_i) \Phi(x_j) = K(x_i, x_j)$  as

$$K(x_i, x_{\text{test}})$$

is called Kernel trick.

① Examples of Kernel

① Linear Kernel

$$K(x_i, x_j) = x_i^T x_j$$

$$= x_i^T I x_j$$

② Polynomial kernel

$$K(x_i, x_j) = (1 + x_i^T x_j)^t ; t > 0$$

### ③ Gaussian kernel (RBF kernel)

Radial Basis Function

$$k(x_i^o, x_j^o) = \exp\left(-\frac{\|x_i^o - x_j^o\|_2^2}{2\sigma^2}\right)$$

$$k_1(x_i, x_j) \neq k_2(x_i, x_j)$$

$$K = k_1(x_i, x_j), k_2(x_i, x_j)$$

$$k(x_i, x_j) = c k(x_i, x_j), c > 0$$

$\Rightarrow$  If  $x_1 > x_2$ ,  $\exp(x_1) > \exp(x_2)$   
(monotonically increasing)

### ④ Soft-SVM

#### ① L1-SVM

$$\min_{(w, w_0, \xi)} \frac{1}{2} \|w\|_2^2 + \frac{1}{n} \sum_{i=1}^n \xi_i$$

$$\text{s.t. } y_i^o(w^T x_i^o + w_0) \geq 1 - \xi_i^o$$

$$\xi_i^o \geq 0 \forall i$$

#### ② L2-SVM

$$\min_{(w, w_0, \xi)} \frac{1}{2} \|w\|_2^2 + \frac{1}{n} \sum_{i=1}^n \xi_i^2$$

#### ③ L1C-SVM

$$+ \frac{c}{n} \sum_{i=1}^n \xi_i^o + \frac{1}{2} \|w\|^2$$

#### ④ L2C-SVM

$$\frac{c}{n} \sum_{i=1}^n \xi_i^o + \frac{1}{2} \|w\|^2$$

$$e^x = 1 + x + \frac{x^2}{2!} + \frac{x^3}{3!} + \dots$$

Date \_\_\_\_\_  
Page \_\_\_\_\_

Let  $x = R$

$$k(x_i^0, x_j^0) = \exp\left(-\frac{(x_i^0 - x_j^0)^2}{2\sigma^2}\right)$$

$$= \exp\left(-\frac{x_i^{02}}{2\sigma^2} - \frac{x_j^{02}}{2\sigma^2} + \frac{2x_i^0 x_j^0}{2\sigma^2}\right)$$

$$= \exp\left(-\frac{x_i^{02}}{2\sigma^2}\right) \exp\left(-\frac{x_j^{02}}{2\sigma^2}\right) * \exp\left(\frac{2x_i^0 x_j^0}{2\sigma^2}\right)$$

$$= \exp\left(-\frac{x_i^{02}}{2\sigma^2}\right) \exp\left(-\frac{x_j^{02}}{2\sigma^2}\right) \left[1 + \frac{x_i^0 x_j^0}{\sigma^2} + \frac{(x_i^0 x_j^0)^2}{2!(\sigma^2)^2} + \dots\right]$$

$$\exp\left(-\frac{x_i^2}{2\sigma^2}\right) \exp\left(-\frac{x_j^2}{2\sigma^2}\right) \left[1, \frac{x_i^0}{\sigma}, \frac{x_i^{02}}{\sqrt{2}\sigma^2}, \dots \rightarrow \left\langle 1, \frac{x_i^0}{\sigma}, \dots \right\rangle\right]$$

$$= \exp\left(-\frac{x_i^{02}}{2\sigma^2}\right) \left[1, \frac{x_i^0}{\sigma}, \frac{x_i^{02}}{\sqrt{2}\sigma^2}, \dots\right]$$

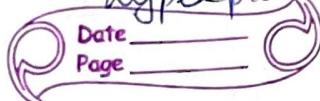
in terms of  $x_i^0$   
in infinite  
dimensions

$$\phi(x_i^0)$$

13/2/23

separating & supporting  
hyperplanes

SVC - default = RBF



SVD → Singular Value Decomposition

### (\*) Eigen Value Decomposition

A (non-zero) vector  $v$  of  $n$  dimension is an eigen vector of the  $n \times n$  matrix  $A$  if it satisfies

$$Av = \lambda v$$

PCA → (variance-covariance matrix)

Uses SVD.

SVD used for:-

- (i) Dimensionality reduction (PCA)
- (ii) Calculating pseudo-inverse  $(ATA)^{-1}AT$

'etc'

(SVD.pdf) - Sir's notes.

$$A = V \begin{pmatrix} \Lambda \\ 0 \end{pmatrix} V^T$$

(diagonal matrix)

(defn of  
eigen vector  
&  
eigen vector)

$$Av = \lambda v$$

$\downarrow$  vectors

$\lambda$  = scalar (eigen value)

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

$$Ax$$

$$\downarrow$$
  
 $x \in \mathbb{R}^d$

$$\rightarrow Ax \in \mathbb{R}^{n \times 1}$$

→ If  $A$  is a square matrix ( $n \times n$ ) with  $n$  linearly independent eigenvectors,  $A$  can be factorised as

$$A = Q \Lambda Q^{-1}$$

$$Q = \begin{bmatrix} | & | & \dots & | \\ v_1 & v_2 & \dots & v_n \\ | & | & \dots & | \end{bmatrix}$$

- Linear comb<sup>o</sup> of vectors
- Basis (size of basis  $\ell$ )
- Null space, column space (or range space)

$Ax$

$$\begin{bmatrix} & 1 \\ A & \begin{bmatrix} a_1 & \dots & a_d \end{bmatrix} \end{bmatrix} \quad x = \begin{bmatrix} x_1 \\ x_2 \\ \vdots \\ x_d \end{bmatrix} \quad \sum_{j=1}^d x_j^0 \underbrace{a_j}_{\substack{\text{vector} \\ \uparrow \\ \text{scalar}}} \quad Ax = \sum_{j=1}^d x_j^0 a_j$$

→ Rank of matrix

$$Ax = 0$$

→ All  $x$  that satisfy  $Ax = 0 \rightarrow$  Null space

dim (range space) + dim (null space)  
= total dimension

→ Every  $n \times n$  real symmetric matrix  
⇒ eigenvalues are real & eigen vectors  
can be chosen orthonormal

$$Av = \lambda v$$

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

$$\therefore A = Q \Lambda Q^{-1}$$

$$A = Q \Lambda Q^T \quad (Q^{-1} = Q^T) \quad \text{where?}$$

(linear algebra book) - ref

Date \_\_\_\_\_  
Page \_\_\_\_\_

$\rightarrow$  (orthonormal?)  $\rightarrow$  k vectors,

$$x_i^{\circ T} x_j^{\circ} = 0 \quad \forall i \neq j$$

→ (difference)

$\psi_1$  &  $\psi_2$  are orthogonal  
 $\&$  orthonormal

$$x_i^T x_j = 1 \quad \forall i = j$$

→ Algebraic multiplicity :- No. of repetitions of a particular eigen value is its algebraic multiplicity.

→ Geometric multiplicity :- No. of linearly independent eigenvectors associated with it. i.e. dimension of null space of  $A - \lambda I$ .

$$AM(e) \geq GM(e)$$

$\Rightarrow$  Issues with EVD - (eigenvalue decomposition)

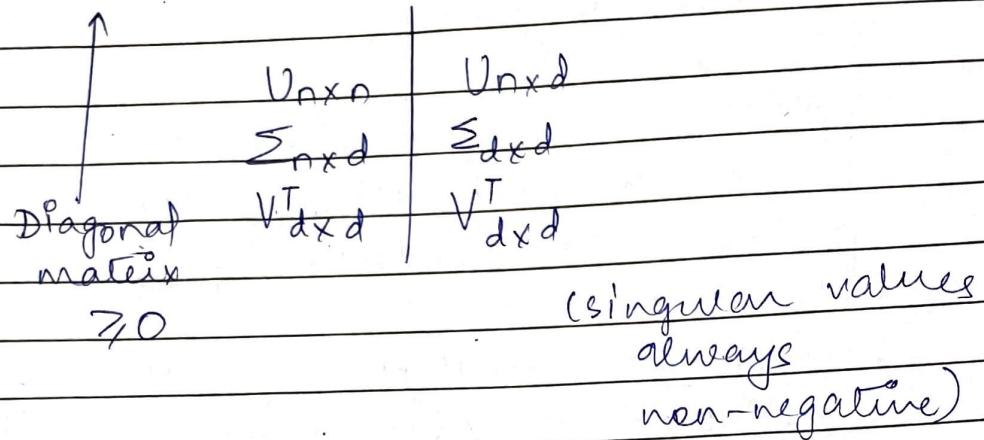
- Only applicable to square matrices
- May be complex eigenvalues

		An $n \times d$	
		$n \times n$	real symm
		Matrices	
$n \times n$	matrices		

SVD

Any matrix  $A$  ( $n \times d$ ) can be decomposed as

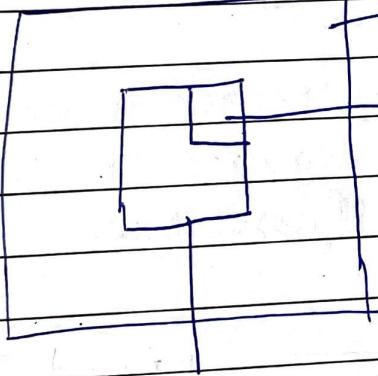
$$A = U \Sigma V^T$$



$$U^T U = I_n$$

$$V^T V = I_d$$

20/2/23



Rectangular  
matrices

Real & symmetric

SVD (contd.)

Spectral theorem

→ Columns of  $U$  and  $V$

are orthonormal

$\Sigma$  is diagonal matrix with  
non-negative real entries.

for SVD.

[Book - Foundations  
for DS by

Ravi Kanau

Hopcroft]

$A^{10 \times 3}$

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

$U^{10 \times 3}$

$\Sigma^{3 \times 3}$

$V^T^{3 \times 3}$

Date \_\_\_\_\_

Page \_\_\_\_\_

$$[U] \begin{bmatrix} \sigma_1 & & \\ & \sigma_2 & \\ & & \sigma_3 \end{bmatrix} [V^T]$$

→ Power iteration methods  
 $(AV = A^2V)$

$A^T A \rightarrow$  square symm.

Suppose we consider the square symm matrix  $\rightarrow A^T A$ .

Let  $x$  be eigenvector of  $A^T A$  and  $\lambda$  be its eigenvalue.

$$A^T A x = \lambda x.$$

Multiplying by  $A$  (on both sides),

$$A A^T (A x) = \lambda (A x) \Rightarrow (A A^T y) = \lambda y$$

$$A A^T = V D V^T$$

$$A A^T = V D U^T$$

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

$$= V \Sigma^T U^T U \Sigma V^T \quad (AB)^T = B^T A^T$$

$$= V \Sigma^T \Sigma V^T$$

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

$$= U \Sigma V^T V \Sigma^T U^T$$

$$= U \Sigma^T \Sigma U^T$$

$$= U \Sigma^2 U^T$$

$U$ -formed by eigenvectors of  $A A^T$

$V$ -formed by  $u$  "  $A^T A$

rank of matrix  
full rank matrix

Singular values  $\rightarrow$  square roots of eigenvalues of  $ATA$ .

Singular vectors:-

Consider rows of  $A$  as  $n$  points in  $d$  dimensions. Consider best fit line through origin. Let  $v$  be unit vector along the line.

The length of projection of  $a_i$  ( $i$ th row of  $A$ ) onto  $v$  is  $|a_i \cdot v|$

Sum of length of squared projections is  $\|Av\|^2$ .

$$\sum \beta_i^2 = \sum_{i=1}^n (a_i \cdot v)^2$$

written in terms of matrix  $A$

Maximizing  $\|Av\|^2$  is best fit line.

First singular vector of  $A$

$$v_1 = \arg \max \|Av\|_2$$

$$\|v\|_2 = 1$$

$$\sigma_1(A) = \|Av_1\|_2$$

$$\text{Why, } v_2 = \arg \max_{V \perp V_1, \|V\|_2=1} \|AV\|_2$$

Similarly find  $v_3, v_4, \dots, v_d$

$$\sigma_3(A) = \|Av_2\|_2$$

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

Imp result - Let  $A$  be  $n \times d$  matrix where

$v_1, \dots, v_k$  are singular vectors.

For  $1 \leq k \leq r$ , let  $V_k$  be subspace

spanned by  $V : v_1, \dots, v_k$ , then

for each  $k$ ,  $V$  is best fit  $k$ -dimensional

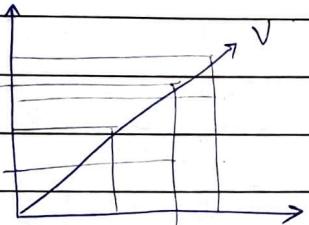
subspace for  $A$ .

23/2/20

Q: How many singular

vectors orthogonal?

rank  
of  
matrix



$$V = \{v_1, v_2\}$$

$$\text{Span}(V) = \{z \mid z = \alpha_1 v_1 + \alpha_2 v_2\}$$

$$V_k = \text{span}(v_1, \dots, v_k)$$

$$\text{Define vectors } u_i^0 = \frac{1}{\sigma_i^0(A)} Av_i$$

$u_1, u_2, \dots, u_d$  are called left singular

vectors of  $A$ .

$$\therefore A = \sum_{i=1}^d \sigma_i^0 u_i^0 v_i^{0 T}$$

scalars       $n \times 1$        $1 \times d$

$$\left\{ \begin{array}{l} A = U \Sigma V^T \\ n \times d \quad d \times d \end{array} \right.$$

$n \times d$   
(full)  
 $d \times d$   
 $n \times d$

matrix formed by 1 outer product - rank 1

$$Av = \sigma u \quad (\text{any } Av \text{ is scaled version of } u)$$

→ Induced Norm

→ General matrix Norm

Date \_\_\_\_\_  
Page \_\_\_\_\_

→ Frobenius Norm:  $\|A\|_F = \sqrt{\sum_{j,k} a_{jk}^2} = \sqrt{\sum_{i=1}^{\text{rank}} \sigma_i^2}$  (singular values)

Spectral Norm (2 norm) - highest singular value of matrix

$$\|A\|_2 = \sigma_1(A)$$

$$\sum \sigma_i^2(A) = \|A\|_F^2$$

→ Frobenius Norm used for low rank approximation.

### ④ Low-Rank Approximation Problem

$$\min_{\tilde{A} \in \mathbb{R}^{n \times d}} \|A - \tilde{A}\|_F^2 \quad \text{s.t. } \text{rank}(\tilde{A}) = k, \quad k < d$$

Best  $\tilde{A}$

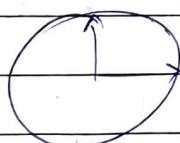
$$\tilde{A} = \sum_{i=1}^k \sigma_i u_i v_i^T$$

### ⑤ Geometric Interpretation

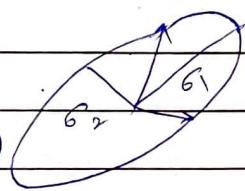
$$A = U \Sigma V^T$$

$$\therefore A\mathbf{x} = U \Sigma V^T \mathbf{x}$$

$A$   
 $(U\Sigma V^T)$



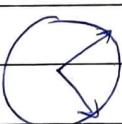
$A$



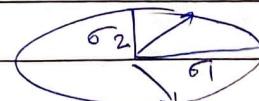
$V^T$

$U$

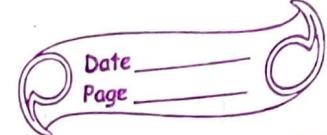
Chang  
ox or  
y's



$\Sigma$



(Scaling)



24/2/2023 ~~Related Work~~

## (1) The Class Imbalance Problem

→ No. of obs<sup>n</sup> in one class label << no. of obs<sup>n</sup> in other class label

Issue →

- Conventional classifiers designed to optimize accuracy
- Biases performance towards majority class.
- More pronounced the imbalance, more pronounced is the issue
- 

(Accuracy not correct measure in this case)

Solutions

### (1) Random undersampling (of majority class).

adv:- simple, efficient

disadv:- - throws away lot of data

- not representative of test data

### (2) Random oversampling (copies of minority class)

adv:- - no data loss

- perform better than undersampling (empirically)

disadv:- prone to overfitting

⇒ Ensemble learning  $\xrightarrow{\text{Bagging}} \text{Bootstrapping}$  (with replacement)  
 $\xrightarrow{\text{Boosting}}$

③ SMOTE: Synthetic Minority oversampling technique

- Avoid overfitting due to exact replicas of minority class samples.
- Subset of minority class is taken
- New synthetic data samples similar to subset are created & added.

→ other techniques/solve

- Penalties based models
- Class weight based models
- Try to solve for different evaluation measure
- Ensemble based models

→ Parameter estimation

$(x, y)$

Discriminative

Generative

Model  $\rightarrow P(Y|X)$

Model  $P(X, Y)$

→ Task of ML → prediction

→ Inference → (find data patterns)

Frequentist approach  $\rightarrow$  run exp. multiple times & come to conclusion  
 Bayesian approach

→ uses prior prob.

$$P(\theta|D) \propto P(D|\theta) \cdot P(\theta)$$

posterior

likelihood

prior

MAP

MLE

→ tries to maximize likelihood

$\theta = \arg \max$

⑦

## Bayesian

Date \_\_\_\_\_  
Page \_\_\_\_\_

27/2/2023

$$P(\theta | D, \alpha) \propto P(D | \theta, \alpha) \cdot P(\theta | \alpha)$$

Posterior      Likelihood      Priors

↓

MLE tries to maximize this.

↳ Point estimates → MAP, MLE(?)

$$P(\theta | \{(x_1, y_1), (x_2, y_2), \dots, (x_n, y_n)\}, \alpha)$$

$$\propto P\{y_1, y_2, \dots, y_n | (x_1, x_2, \dots, x_n), \theta, \alpha\}$$

$$\cdot P(\theta | \alpha)$$

$$P(y_1, y_2, \dots, y_n | x_1, x_2, \dots, x_n, \theta)$$

$$= L(\theta)$$

$$= \prod_{i=1}^n P(y_i | x_i, \theta)$$

[ $\theta$  = parameter]

→ ~~exp~~ conjugate prior (?)

$$\hat{\theta}_{MLE} = \operatorname{argmax}_{\theta} \prod_{i=1}^n P(y_i | x_i, \theta)$$

log likelihood → why? → bcoz log is  
monotonically increasing function  
& applying log converts it to summation,

$$= \operatorname{argmax}_{\theta} \sum_{i=1}^n \log(P(y_i | x_i, \theta))$$

$$= \operatorname{argmin}_{\theta} \sum_{i=1}^n -\log(P(y_i | x_i, \theta))$$

-ve  
log  
likelihood

$$l(\theta; (x, y)) = -\log P(y|x, \theta) \quad (\text{loss on single point})$$

Represents loss

$$\text{Loss}(\theta, (x, y)) = (\theta^T x - y)^2 \quad (\text{MSE})$$

$$P(y_i^o | x_i^o, (w, \beta)) = N(y_i^o | w^T x_i^o, \frac{1}{\beta})$$

$\downarrow$

$\theta$  (Normal distribution)

$$\mathbb{E} \epsilon \sim N(0, \frac{1}{\beta}) \quad (\text{error})$$

$$\begin{aligned} \text{Likelihood} &= \prod_{i=1}^n P(y_i^o | x_i^o, w, \beta) \\ &= \prod_{i=1}^n N(y_i^o | w^T x_i^o, \frac{1}{\beta}) \\ &= \prod_{i=1}^n \frac{\sqrt{\beta}}{\sqrt{2\pi}} \exp \left\{ -\left( \frac{1}{2} \beta (y_i^o - w^T x_i^o) \right)^2 \right\} \end{aligned}$$

$$X = \begin{bmatrix} x_1 \\ \vdots \\ x_n \end{bmatrix} \sim N(\mu, \sigma^2)$$

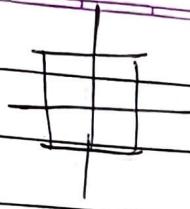
$$\text{Mean vector, } \mu = \begin{bmatrix} \mu_1 \\ \vdots \\ \mu_n \end{bmatrix}$$

Variance-  $= \sum$   
 covariance  
 matrix

multivariate Gaussian

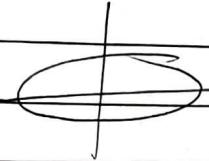
$$f_x(x) = \frac{1}{(2\pi)^{n/2} |\Sigma|^{1/2}} \exp \left( -\frac{1}{2} (x - \mu)^T \Sigma^{-1} (x - \mu) \right)$$

max norm  $\rightarrow$



Mahalanobis distance

- geometrical concept



$$\max_{(w, \beta)} \prod_{i=1}^n \left( \frac{1}{\sqrt{2\pi}} \sqrt{\beta} \exp \left( -\frac{1}{2} \beta (y_i - w^T x_i)^2 \right) \right)$$

$$-\text{ve neg likelihood} = \sum_{i=1}^n \left( -\log \sqrt{\beta} + \log \sqrt{2\pi} \right. \\ \left. + \frac{1}{2} \beta (y_i - w^T x_i)^2 \right)$$

min -ve neg likelihood

$$\frac{\partial L(w, \beta)}{\partial w} = 0 \Rightarrow \sum_{i=1}^n -\frac{1}{2} \beta \cdot 2(y_i - w^T x_i) \cdot x_i = 0$$

$$\frac{\partial L(w, \beta)}{\partial \beta} = 0 \Rightarrow$$

$$\sum_{i=1}^n \left[ -\frac{1}{2\beta} + \frac{1}{2} (y_i - w^T x_i)^2 \right] = 0$$

$$w = (x^T x)^{-1} x^T y$$

$$\frac{1}{\beta} = \frac{1}{n} \sum (y_i - w^T x_i)^2$$

$\nwarrow$  unexplained variance

$$P(y_i | x_i^o, w, \beta)$$

$P(y_i | x_i^o)$  sampled from normal dist.



$$P(y_i | x_i^o) \sim N(w^T x_i^o, \frac{1}{\beta})$$

$$\begin{bmatrix} y_i = w^T x_i^o \\ + \varepsilon \end{bmatrix}$$

→ -ve log likelihood

$$\frac{\partial L(w, \beta)}{\partial w} = 0 \Rightarrow w_{MLE} = (X^T X)^{-1} X^T Y$$

Gauss-  
ian  
noise  
with  
 $\downarrow$   
 $N(0, \beta)$

$$\frac{\partial L(w, \beta)}{\beta} = 0 \Rightarrow \beta = \frac{1}{n} \sum (y_i - w^T x_i^o)^2$$

(unexplained  
variance)

(Frequentist approach)

→ For RN, we find expectation & variance.

⇒ Gauss-Markov theorem :-



(Best in terms of variance)

$$E_D [w_{MLE}] = w^*$$

$$E_D [(X^T X)^{-1} X^T Y] = (X^T X)^{-1} X^T E_D [Y] \quad \text{--- ①}$$

(locally  
randomness  
only)

$$Y = Xw + \varepsilon$$

$$\varepsilon = \begin{bmatrix} \varepsilon_1 \\ \vdots \\ \varepsilon_n \end{bmatrix} \sim N(0, \frac{1}{\beta})$$

$$E[Y] = E[Xw] + E[\varepsilon]$$

(due to linearity  
of expectation)

$$= E[Xw] + 0 \quad (\varepsilon \sim N(0, \frac{1}{\beta}))$$

Replacing in ①

$$\begin{aligned} E_D[(X^T X)^{-1} X^T Y] &= (X^T X)^{-1} X^T E_D[Y] \\ &= \underline{(X^T X)^{-1} X^T X w} \\ &= Iw \\ &= w \end{aligned}$$

$$\therefore E_D [w_{MLE}] = w^* \quad (\because \text{It's unbiased dist})$$

$$\begin{aligned}
 E_D \left[ \frac{1}{\beta_{MLE}} \right] &= E_D \left[ \frac{1}{n} (x w_{MLE} - y)^T (x w_{MLE} - y) \right] \\
 &= \frac{1}{n} E_D \left[ \underbrace{\frac{w_{MLE}^T}{\downarrow} x^T x}_{\text{I}} \underbrace{w_{MLE}}_{\text{II}} - 2 w_{MLE}^T x^T y + y^T y \right] \\
 &= \frac{1}{n} E_D \left[ Y^T x (x^T x)^{-1} \cdot (x^T x) (x^T x)^{-1} x^T y \right. \\
 &\quad \left. - 2 \cancel{Y^T x (x^T x)^{-1} x^T y} + Y^T y \right] \\
 &= \frac{1}{n} E_D \left[ -Y^T x (x^T x)^{-1} x^T y + Y^T y \right] \\
 &= \frac{1}{n} E_D \left[ Y^T \underbrace{\left( I - x (x^T x)^{-1} x^T \right)}_Z Y \right] \\
 &\quad (x^T A x + x^T B x) \\
 &\quad = x^T (A + B) x \\
 &= \frac{1}{n} E_D [Y^T Z Y] \\
 &= \frac{1}{n} E_D [(x w + \varepsilon)^T Z (x w + \varepsilon)] \\
 &= \frac{1}{n} E_D [w^T x^T Z x w + \underset{0}{\underset{\parallel}{\varepsilon^T Z x w}} + w^T x^T Z \varepsilon + \underset{0}{\underset{\parallel}{\varepsilon^T Z \varepsilon}}] \\
 &= \frac{1}{n} [w^T x^T Z x w + 0 + 0 + \sum_{i,j} z_{i,j}^0 + E[\varepsilon_i^0, \varepsilon_j^0]]
 \end{aligned}$$

( $Z$  is  $n \times n$  matrix)

$$\mathbf{\Sigma}^T \mathbf{z} \mathbf{\Sigma} = [\varepsilon_1 \ \varepsilon_2 \ \dots \ \varepsilon_n] \begin{bmatrix} z_{11} & z_{12} & \dots \\ & \ddots & \\ & & z_{nn} \end{bmatrix} \begin{bmatrix} \varepsilon_1 \\ \vdots \\ \varepsilon_n \end{bmatrix}$$

$$= \sum_{i,j} z_{ij} \cdot \varepsilon_i \varepsilon_j$$

5

$$E[\varepsilon_i \varepsilon_j] = 0 \text{ when } i \neq j$$

$$\sum_i \sum_j P_{\mathbf{z}}(\varepsilon_i \varepsilon_j) \cdot \varepsilon_i \varepsilon_j$$

↓

$$P_{\mathbf{z}}(\varepsilon_i) P_{\mathbf{z}}(\varepsilon_j)$$

$$\sum_i P_{\mathbf{z}}(\varepsilon_i) \cdot \sum_j P_{\mathbf{z}}(\varepsilon_j)$$

From ②, (contd)

$$= \frac{1}{n} [w^T X^T Z X w + 0 + 0 + \frac{1}{B} \sum_{i=1}^B \text{diag}(z)]$$

$$= \frac{1}{n} [w^T X^T (I - X(X^T X)^{-1} X^T) X w + \frac{1}{B} \sum_{i=1}^B \text{diag}(z)]$$

$$= \frac{1}{n} [w^T X^T X w - w^T X^T X (X^T X)^{-1} X w + \frac{1}{B} \sum_{i=1}^B \text{diag}(z)]$$

$$= \frac{1}{n} \left( \frac{1}{B} \sum_{i=1}^B \text{diag}(z) \right)$$

$$= \frac{1}{n} \left[ \frac{1}{B} \text{tr}(z) \right]$$

(sum of diagonals  
or  
trace)

$$= \frac{1}{n} \left[ \frac{1}{\beta} (n - \text{tr}(x^T x (x^T x)^{-1} x^T)) \right]$$

$$= \frac{1}{n} \left[ \frac{1}{\beta} \cdot (n - \text{tr}(x^T x (x^T x)^{-1})) \right]$$

→ d dimension

$$= \frac{1}{n} \left[ \frac{1}{\beta} (n-d) \right]$$

(trace of identity matrix  $d \times d$ )

$$\mathbb{E}_{\beta_{MLP}} \left[ \frac{1}{\beta} \right] = \left(1 - \frac{d}{n}\right) \frac{1}{\beta}$$

$$= d$$

$$\begin{bmatrix} 1 & 0 & \dots & 0 \\ 0 & 1 & \dots & 0 \\ \vdots & \vdots & \ddots & \vdots \\ 0 & 0 & \dots & 1 \end{bmatrix}$$

# ADVANCED ML (CLASS NOTES)

classmate

Date

Page

13/3/28

- ⑩ MLE for linear regression

$$E[w_{MLE}] = w.$$

$$E\left[\frac{1}{\beta_{MLE}}\right] = \left(1 - \frac{\alpha}{n}\right) \frac{1}{\beta}$$

Gram-Schmidt theorem

$w_{MLE} \rightarrow \text{vector}$

$w_{MLE}$  is  
BLUE

$$\text{var}[w_{MLE}] < \text{var}[\tilde{w}]$$

where

$\tilde{w}$  is any  
other linear  
estimate

Best Linear  
Unbiased  
Estimator

$$A - B \geq 0$$

means

$A - B$  is the semi-definite

$$\Rightarrow \text{var}[w_{MLE}] - \text{var}[\tilde{w}] \geq 0$$

then,  $\text{var}[w_{MLE}]$  is best variance.

$$\text{var}(w_{MLE})$$

$$= E[(w_{MLE} - w)(w_{MLE} - w)^T]$$

$$\text{var}(x)$$

$$= E[(x - u)(x - u)^T]$$

$$= E[((x^T x)^{-1} x^T y - w)((x^T x)^{-1} x^T y - w)^T]$$

$$\leftarrow y = xw + \varepsilon$$

$$= E_D[(x^T x)^{-1} x^T (xw + \varepsilon) - w)$$

$$\begin{aligned} & ((x^T x)^{-1} x^T (xw + \varepsilon) \\ & - w)^T \end{aligned}$$

$$= E_D[w + (x^T x)^{-1} x^T \varepsilon - w]$$

$$= E_D \left[ ((x^T x)^{-1} x^T \varepsilon) ((x^T x)^{-1} x^T \varepsilon)^T \right]$$

$$= E_D \left[ (x^T x)^{-1} x^T (\varepsilon \varepsilon^T) x (x^T x)^{-1} \right]$$

$$= (x^T x)^{-1} x^T \underbrace{E(\varepsilon \varepsilon^T)}_{\text{var-cov matrix}} x (x^T x)^{-1}$$

$$= (x^T x)^{-1} x^T \frac{1}{B} I_{n \times n} x (x^T x)^{-1} \quad (\varepsilon \sim N(0, \frac{1}{B}))$$

$$\text{var}(w_{MLE}) = \boxed{\frac{1}{B} (x^T x)^{-1}}$$

$$\left[ (x^T x)^{-1} x^T x = I \right]$$

var-cov matrix

$\Rightarrow$  Proving the 'Best' part in BLUE.

$$\underline{T.P.} : -\text{var}[w_{MLE}] + \text{var}[\tilde{w}] > 0.$$

$$\tilde{w} = A \cdot Y = A(xw + \varepsilon) = AYw + AE$$

$$E(Axw + Ae) = w \quad [\because \text{unbiased}]$$

$$\begin{matrix} \therefore Axw = w \\ \uparrow \qquad \uparrow \\ \text{matrix} \qquad \text{scalar} \end{matrix}$$

(characteristic eq<sup>b</sup>)

$\Rightarrow w$  is eigenvectors of  $Ax$  & corresponding eigenvalue is 1.

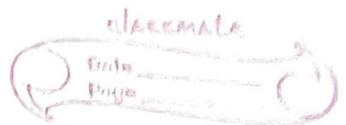
$$\text{var}[\tilde{w}] = E[(\tilde{w} - E(\tilde{w}))(\tilde{w} - E(\tilde{w}))^T]$$

$$\begin{matrix} \tilde{w} = Axw + Ae \\ \downarrow \\ \tilde{w} = E[(\tilde{w} - Axw)(\tilde{w} - Axw)^T] \\ = E[(A\varepsilon)(A\varepsilon)^T] \end{matrix}$$

$$= E[A\varepsilon\varepsilon^T A^T]$$

$$= \frac{1}{B} (AA^T)$$

$$w_{MLE} = \frac{(X^T X)^{-1} X^T Y}{M Y}$$



$$\text{var}(\tilde{w}) - \text{var}(w_{MLE})$$

(should  
be PSD)

$$= \frac{1}{p} (A A^T - (X^T X)^{-1}) \quad \textcircled{1}$$

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

multiply  $X$  on both sides,

$$AX = (X^T X)^{-1} X^T X + BX$$

$$\Rightarrow (AX - I) = BX$$

(prove that  $BX = 0$ )

$$(AX - I)w = 0$$

$$\Rightarrow BXw = 0$$

$$\text{var}[\tilde{w}] - \text{var}(w_{MLE})$$

$$= \frac{1}{p} [BB^T]$$

From ①, if we  
replace  $A$ , we get  
this)

( $BB^T$  is the semi-def.

$\downarrow$  true?

$$(x^T A x \geq 0 \forall x)$$

$$BB^T \Rightarrow$$

$$x^T B B^T x$$

$$(B^T x)^T B^T x$$

$$= \|B^T x\|_2^2$$

(which is  
always  $\geq 0$ )

$\therefore$  PSD.

④ Difference between  $w_{MLE}$  &  $w_{MAP}$ ?

Posterior  $\propto$  prior  $\times$  likelihood

$$P(w | (y_1, \dots, y_n) (x_1, \dots, x_n), \underbrace{\mu_0, \varepsilon_0, \frac{1}{\beta}}_{\text{prior}})$$

$$\propto P[(y_1, y_2, \dots, y_n) | w, \frac{1}{\beta}, (x_1, \dots, x_n)]$$

$$, P(w | \mu_0, \varepsilon_0)$$

↑  
prior.

16/3/2023

$$y_i = w^T x_i + \varepsilon_i$$

$$\begin{matrix} \uparrow & \uparrow \\ N(\mu_0, \varepsilon_0) & N(0, \frac{1}{\beta}) \end{matrix}$$

$$P(w | (y_1, \dots, y_n) (x_1, \dots, x_n)) \propto P(y_1, \dots, y_n | w, \frac{1}{\beta}, \underbrace{\mu_0, \varepsilon_0, \frac{1}{\beta}}_{(x_1, \dots, x_n)})$$

$$, P(w | \mu_0, \varepsilon_0)$$

$$P(w | (y_1, \dots, y_n) (x_1, \dots, x_n), \varepsilon_0, \mu_0, \frac{1}{\beta})$$

$$\propto \left[ \prod_{i=1}^n p(y_i | w, \frac{1}{\beta}, x_i) \right] \cdot P[w | \varepsilon_0, \mu_0]$$

$$\left[ \prod_{i=1}^n N(y_i | w^T x_i, \frac{1}{\beta}) \right] \cdot P[w | (\varepsilon_0, \mu_0)]$$

$$y_i^o = w^T x_i^o + \epsilon_i^o$$

$\uparrow$                        $\nwarrow N(0, \frac{1}{\beta})$   
 $N(\mu_0, \Sigma_0)$

multivariate gaussian

$$P(w | D, \Sigma_0, \mu_0, \frac{1}{\beta}) = \prod_{i=1}^n \frac{\sqrt{\beta}}{\sqrt{2\pi}}$$

$$= \exp \left( -\frac{1}{2} \frac{(y_i^o - w^T x_i^o)^2}{\beta} \right)$$

$$\cdot \frac{1}{(2\pi)^{d/2} |\Sigma_0|} \cdot \exp \left( -\frac{1}{2} \frac{(w - \mu_0)^T \Sigma_0^{-1} (w - \mu_0)}{\beta} \right)$$

(covariance matrix)

$$(w - \mu_0)^T \Sigma_0^{-1}$$

$$(w - \mu_0)$$

$$= -\log P(w | D, \Sigma_0, \mu_0, \frac{1}{\beta})$$

Posterior

$$= -n \log \sqrt{\beta} + n \log \sqrt{2\pi} + \frac{\beta}{2} \sum_{i=1}^n (y_i^o - w^T x_i^o)^2$$

$$+ \frac{1}{2} \log |\Sigma_0| + d \log \sqrt{2\pi} + \frac{1}{2} (w - \mu_0)^T \Sigma_0^{-1} (w - \mu_0)$$

$$= \log \left( \frac{1}{2\pi} \right)^{d/2} + \log \left( \frac{1}{|\Sigma_0|} \right)^{-1/2}$$

"Matrix cookbook" - book

Day  
Page

$$-\log P(w|D, \Sigma_0, \mu_0, \beta)$$

then

$$\frac{\partial}{\partial w} \left[ -\log P(w|D, \Sigma_0, \mu_0, \beta) \right] = 0$$

$$= \frac{\partial}{\partial w} \left[ \frac{\beta}{2} (xw - y)^T (xw - y) \right]$$

$$+ \frac{\partial}{\partial w} \left[ \frac{1}{2} (w - \mu_0)^T \Sigma^{-1} (w - \mu_0) \right]$$

$$= 0$$

$$\therefore \frac{\beta}{2} [2(x^T x)w - 2x^T y] + \frac{1}{2} [2\Sigma_0^{-1}w - 2\Sigma_0^{-1}\mu_0] = 0$$

$$(\beta(x^T x) + (\Sigma_0)^{-1})w - \beta x^T y - \Sigma_0^{-1}\mu_0 = 0$$

$$\downarrow$$

$$\therefore w_{MAP} = \beta x^T x + \Sigma_0^{-1}\mu_0$$

$$\therefore w_{MAP} = (\beta(x^T x) + \Sigma_0)^{-1} \cdot (\beta x^T y + \Sigma_0^{-1}\mu_0)$$

prior is infinitely broad, then we get exact MLE estimate

$$\text{If } \mu_0 = 0, \text{ spherical prior, } \Sigma_0^{-1} = \frac{1}{c} I$$

$$(\beta(x^T x) + \frac{1}{c} I)^{-1} \cdot (\beta x^T y)$$

$$(x^T x)^{-1} x^T y + \frac{\beta c}{\beta + c} x^T y$$

classmate

Date \_\_\_\_\_

Page \_\_\_\_\_

$$\left( (X^T X) + \frac{I}{\beta C} \right)^{-1} X^T Y$$

(like Ridge regressions)

Takeaway  $\Rightarrow$  (treated as  
L2 neg. problem)

20/3/23

## (1) Ridge Regression

$$\|Ax - b\|^2 + \lambda x^T x \Rightarrow (A^T A + \lambda I)^{-1} A^T y.$$

Spherical prior :-

$$\begin{aligned} E[w_{MAP}] &= E[(x^T x + \lambda I)^{-1} x^T y] \quad (y = xw + \varepsilon) \\ &= E[(x^T x + \lambda I)^{-1} (x^T w + x^T \varepsilon)] \\ &= E[(x^T x + \lambda I)^{-1} (x^T w) + (x^T x + \lambda I)^{-1} E(x^T \varepsilon)] \end{aligned}$$

$$\begin{aligned} &= \underbrace{(x^T x + \lambda I)^{-1} (x^T w)}_z \\ &= z x w \end{aligned}$$

$$\mu_0 = 0$$

Variance :-

$$\text{Var}[w_{MAP}] = E[w_{MAP} - E(w_{MAP})]$$

$$[w_{MAP} - E(w_{MAP})]^T$$

$$= E[(z y - z x w)(x y - z x w)^T]$$

$$= E[E(y - x w)(y - x w)^T z^T]$$

$$= z E[\varepsilon \varepsilon^T] z^T$$

$$= z \frac{1}{B} I_{n \times n} z^T$$

$$= \frac{1}{B} z z^T$$

$$= \frac{1}{B} (x^T x + \lambda I)^{-1} (x^T x) [x^T x + \lambda I]^{-1} \leftarrow \text{PD}$$

$$\underline{\text{SVD}} : - \underline{x^T x} = \underline{\sqrt{\sum^2} \sqrt{v^T}}$$

$$(x^T x)^{-1} = \sqrt{\text{diag}} \left( \frac{1}{\sigma_i^2} \right) v^T$$

$$x^T x + \lambda I = \sqrt{\sum^2} \sqrt{v^T} + \underline{\lambda v I v^T} \quad v \rightarrow \text{Rotating}$$

$$\begin{aligned} &= \sqrt{(\sum^2 + \lambda I)} v^T \\ &= \sqrt{\text{diag}} (\sigma_i^2 + \lambda) v^T \end{aligned}$$

$$(x^T x + \lambda I)^{-1} = \sqrt{\text{diag}} \left( \frac{1}{\lambda + \sigma_i^2} \right) v^T$$

$$\begin{aligned} (x^T x + \lambda I)^{-1} x^T x &= \sqrt{\text{diag}} \left( \frac{1}{\lambda + \sigma_i^2} \right) \frac{v^T (\sqrt{\sum^2} v^T)}{I} \\ &= \sqrt{\text{diag}} \left( \frac{1}{\lambda + \sigma_i^2} \right) \sum_i^2 v^T \quad (\text{as } \lambda \uparrow, \\ &\quad \text{Var} \downarrow \approx 0) \end{aligned}$$

Bias-Variance trade off

$n \rightarrow$  no. of points  
 $\lambda \rightarrow$  Regularized parameter

For  $n=0 \rightarrow$  w<sub>MLE</sub> doesn't exist.

$$n=0, \lambda=0$$

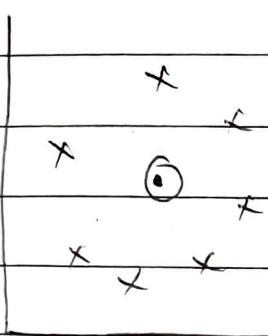
$$n=0, 0 < \lambda < \infty$$

$$\bullet u_0^T w$$

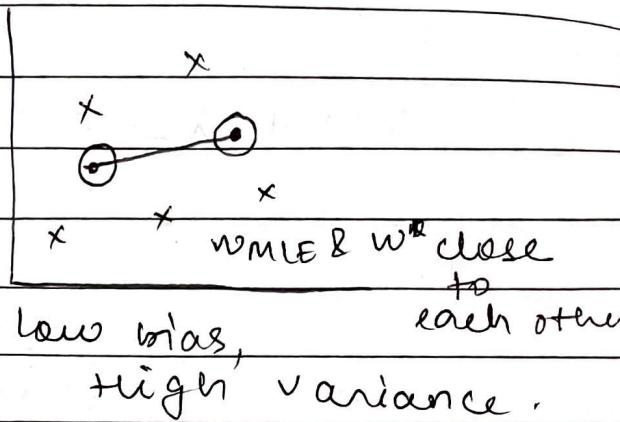
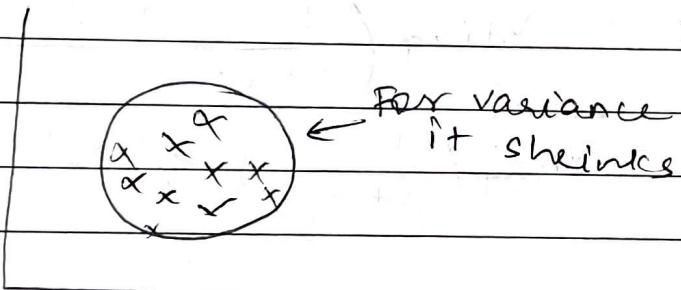
$$\bullet w^*$$

$$n=0, \lambda \rightarrow \infty$$

$$\bullet w$$

$0 < n < \infty, \lambda = 0$  $0 < n < \infty, D \propto \lambda \rightarrow \infty$ 

Reality scenario

 $0 < n < \infty, \lambda \rightarrow \infty$  $n \rightarrow \infty, \lambda = 0$  $n \rightarrow \infty, 0 < \lambda < \infty$  $n \rightarrow \infty, \lambda \rightarrow \infty$ 

- ④ Taxonomy of info - Intents are in terms of these categories.

IR Domain - IA Set

q - query

 $C_{(q)}$  - set of categories of q

d - document

 $C_{(d)}$  - " " " " dAssumptions :-

1) Taxonomy of info.

2)  $p(C|q)$  is drawn ( $C$  = classifier)

$$\sum_{C \in C_{(q)}} p(C|q) = 1 \rightarrow \text{Knowledge completeness}$$

23/3/23 ④ GMM:- Gaussian Mixture Models.

(Sir's notes → GMMS.pdf)

→ Gaussian distribution by itself may not be enough in modelling real datasets.

(Fourier transformation)

→ Sometimes a linear superposition of two or more Gaussians may capture data better. Such superpositions can be formulated as mixture distributions.

→ For a superposition of K Gaussians, we have,

$$p(x) = \sum_{k=1}^K w_k N(x | \mu_k, \Sigma_k)$$

Each  $N(x | \mu_k, \Sigma_k)$  is called a component of the mixture. The parameters  $w_k$  are called ~~mixto~~ mixing coefficients.

$$w = \{w_k\} \quad \mu = \{\mu_k\} \quad \Sigma = \{\Sigma_k\}$$

$$w_k \geq 0 \quad \forall k$$

$\therefore 0 \leq w_k \leq 1$  (can be thought of as prob.)

$$P(x) = \sum_{k=1}^K P(k) p(x|k) \quad (w_k = p(k))$$

Think of  $p(k)$  as prior of picking  $k^{th}$  component and  $p(x|k) = N(x|\mu_k, \Sigma_k)$

i.e. prob. of  $x$  conditioned on  $k$ .

→ The posterior prob.  $P(k|x)$  are given as

$$\gamma_k(x) = p(k|x)$$

$$= \frac{p(k) p(x|k)}{\sum_l p(l) p(x|l)} \quad (\text{Bayes' theorem})$$

$$= \frac{w_k N(x|\mu_k, \Sigma_k)}{\sum_l w_l N(x|\mu_l, \Sigma_l)}$$

→ Parameters of GMM are  $w$ ,  $\mu$  &  $\Sigma$  (rest part - refer Sir's notes)

→ Topic Modelling - e.g. of latent model

24/3/23

$$\begin{aligned}
 & \text{Now } \pi(z_k) \text{ i.e. } P(z_k=1|x) \\
 & = \frac{P(z_k=1) P(x|z_k=1)}{\sum_i P(z_i=1) P(x|z_i=1)} \\
 & = \frac{w_k N(x|\mu_k, \Sigma_k)}{\sum_i w_i N(x|\mu_i, \Sigma_i)}
 \end{aligned}$$

→ Maximizing log likelihood of GMM is more complex problem than the case of single Gaussian. Because of summation over  $k$  inside log, log func<sup>n</sup> does not act directly ...

### ④ EM for GMMS:- (Expectation Maximization)

↳ The log likelihood for GMM given data  $X$  is given by :-

$$\ln P(X|w, \mu, \Sigma) = \sum_{n=1}^N \ln \left\{ \sum_{k=1}^K w_k N(x_n | \mu_k, \Sigma_k) \right\}$$

Derivating w.r.t.  $\mu_k$  and equating to 0,

$$\begin{aligned}
 0 = - \sum_{n=1}^N \frac{w_k N(x_n | \mu_k, \Sigma_k)}{\sum_j w_j N(x_n | \mu_j, \Sigma_j)} \cdot \Sigma_k (x_n - \mu_k) \\
 \underbrace{\qquad\qquad\qquad}_{\pi(z_{nk})}
 \end{aligned}$$

Multiplying by  $\Sigma_k^{-1}$ ,

$$u_k = \frac{1}{N_k} \cdot \sum_{n=1}^N \gamma(z_{nk}) x_n$$

$$\text{where, } N_k = \sum_{n=1}^N \gamma(z_{nk})$$

$N_k$  = effective no. of points belonging to cluster  $k$ .

AA

$$\text{Similarly, } \Sigma_k = \frac{1}{N_k} \cdot \sum_{n=1}^N \gamma(z_{nk}) (x_n - u_k) (x_n - u_k)^T$$

Finally we need to maximize for  $P(x|w, u, \Sigma)$  w.r.t.  $w_k$

We consider,

$$\ln P(x|w, u, \Sigma) + \lambda \left( \sum_{k=1}^K w_k - 1 \right)$$

which gives,

$$0 = \sum_{n=1}^N \cdot \frac{N(x_n|u_k, \Sigma_k)}{\sum_j w_j N(x_n|u_j, \Sigma_j)} + \lambda$$

Multiplying both sides by  $w_k$  and sum over  $k$ , we get  $\lambda = -N$

$$\text{Finally, } w_k = \frac{N_k}{N}$$

④ EM for gaussian mixtures:-

1. Initialize  $\mu_k, \Sigma_k$  and initial value of  $w_k$  and evaluate log likelihood

2. E-step

$$\gamma_k(z_n) = \frac{w_k N(x_n | \mu_k, \Sigma_k)}{\sum_{j=1}^k w_j N(x_n | \mu_j, \Sigma_j)}$$

(Recalculates posterior)

$$3. \mu_k^{\text{new}} = \frac{1}{N_k} \sum_{n=1}^N \gamma_k(z_n) x_n$$

$$\Sigma_k^{\text{new}} = \frac{1}{N_k} \sum_{n=1}^N \gamma_k(z_n) (x_n - \mu_k^{\text{new}})^T (x_n - \mu_k^{\text{new}})$$

$$w_k^{\text{new}} = \frac{N_k}{N}$$

$$\text{where, } N_k = \sum_{n=1}^N \gamma_k(z_n)$$

4. Evaluate  $\sum_{n=1}^N \ln \left\{ \sum_{k=1}^K w_k N(x_n | \mu_k, \Sigma_k) \right\}$

Repeat 2-4 till convergence

Quiz next Monday  $\rightarrow$  Till GMMs.  $\rightarrow$  till here.