

ADVANCED MACHINE LEARNING

9/1/23

→ 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 func²: $f: x \rightarrow y$.

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) $L_s(h) = \frac{1}{m} \sum_{i=1}^m \delta_i$ if $i \in [m] : h(x_i) \neq y_i$

Search for a solⁿ that works well on available data

\rightarrow I.I.D assumption

\rightarrow Goal is that func^t performs well on test data.

Unknown target func^t

$$f: x \rightarrow y$$

Training exs

Unknown I.P distribution

Learning algo

Error measure

Hypothesis set \mathcal{H}

Final hypothesis (goal)

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 Marble 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 - proportion of red marbles out of N

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

If 'N' is large enough,
then μ is near to
empirical mean

② Hoeffding's Inequality

- type of Chernoff Bound.

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

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

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

ϵ - 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

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

(derivation)

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

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

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

\rightarrow when running the learning algo.

④ Define

Dichotomies

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

19/11/23

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

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

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

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 $\forall N$, $dvc(H) = \infty$

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

⇒ Sauer's Lemma

→ VC-Generalization bound

- more complex hypothesis - doesn't generalize well

Evaluation

① Reproducibility challenge - write reproducibility report

② Project

paperswithcode.com/

ICML 2021/22

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 implementation

- Replicate results (on same datasets)

- can you do something more?

◦ Tweak algo

◦ More dataset

◦ Hyperparameter tunings

- Replicate based on library changing

① Del. project empirical
② theory

repository ✓
for this

→ fix
distill. pub

(→

20/1/28

11

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 Worse:-

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

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

→ 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

11

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)

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

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

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

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

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

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

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

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

$$\begin{aligned} &= E_D[(g^{(D)}(x) - \bar{g}(x))^2] + (\bar{g}(x) - f(x))^2 \\ &= \text{Var}(x) + \text{bias} \end{aligned}$$

K
11

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 L₂ 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

27/1/23.

1/1

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

SVM
problem

Optimization for ML

continuous
Non convex

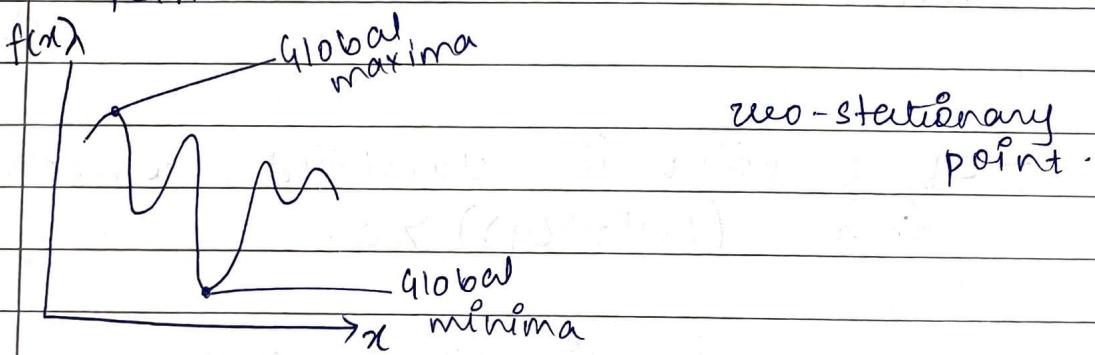
Convex

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 funcⁿ, which can also get stuck at local minima or saddle point.



— / —

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

J.

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:-

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

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

conds:-

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.

~~classmate~~

① _/_/_

- Eg. of unconstrained optimization problem → Linear regression
- Eg. of constrained → Lasso, Ridge
- constraints $\rightarrow \|x\|_1$ $\rightarrow \|x\|_2^2$

$H(x)$
↑
Hessian

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

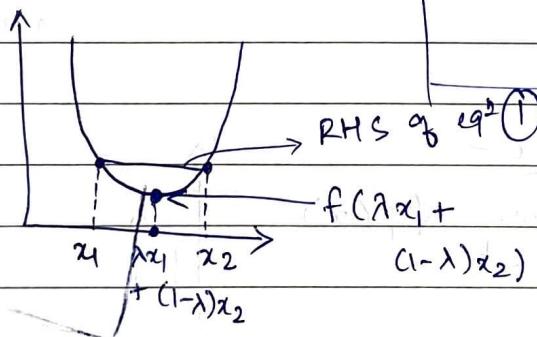
→ Convex function:- If $f: \mathbb{R}^n \rightarrow \mathbb{R}$

(def²) $\text{RHS of eq } ① \leq f(x_1) + (1-\lambda)f(x_2)$

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

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

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



$\therefore f$ is convex func

If in eq² ①, ' $<$ ' sign \rightarrow strictly convex

max (0, ∞)

$\|x\|_1$

→ 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 func² is convex, local minima = global minima

Eg:

① $w^T x + b$ on R^n - convex

② -ve log likelihood func² - "

③ quadratic func² $\Rightarrow f(x) = \frac{1}{2} x^T A x + b^T x + c$ en R^n
(A is PSD).

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

⑩ Operations preserving convexity

① $f(x)$ is convex

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

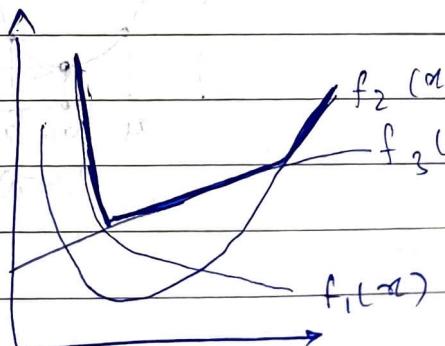
② Sum:- $f_1(x) + f_2(x) = f(x)$

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

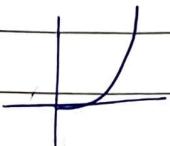
$$f(x) = \log \sum_{i=1}^n \exp(x_i)$$

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

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

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

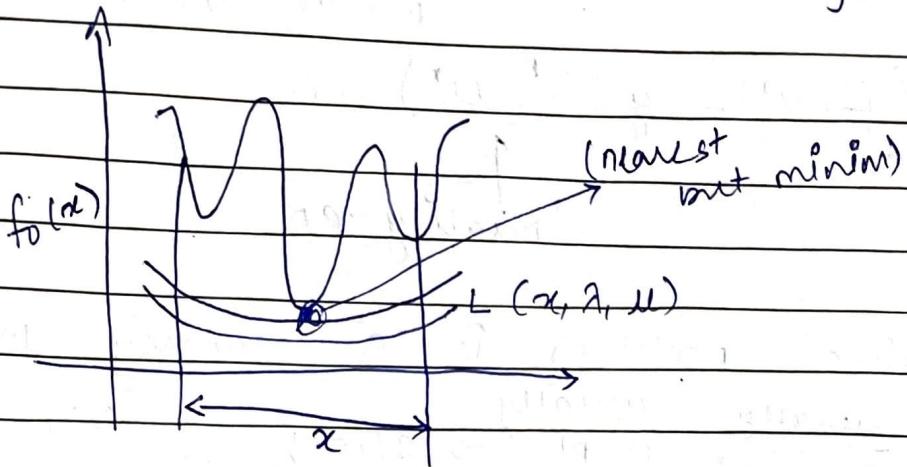
→ where λ_i, μ_j are called Lagrangian multipliers.

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

→ What if x is a feasible point?

↓
 means it
 satisfies
 constraints)

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



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

$$L(\lambda, \mu) = 2\lambda_1 \mu_1^2 + 3\lambda_2 \mu_2^2$$

$\min L(x, \lambda, \mu)$ such that $\lambda > 0$ } $\leq f_0(x^*)$
 lower bound on optimal solution

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

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

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

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 $\rightarrow g(\lambda, \mu) \leq [g(\lambda^*, \mu^*) \leq f_0(x^*)]$
 (λ^*, μ^*) are point of maxima.

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 λ^*, u^* - point of
 $\rightarrow x^*$ is feasible sol¹ for dual optimal
primal problem. 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(func² of x)

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

→ Any tuple (x, λ, μ)

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

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 func² is convex)

④ General stat.

If duality gap = 0

 x^* - Primal optimal λ^* - dual optimal μ^*

9/3/23

AML

Date
Pagescribble
notesApplication for SVM:-

Optimizations

book

⑩ 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 & +1) labels are used

margin

→ Primal SVM

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

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

⑩ Soft margin SVM → slack

$$(w^T x_i + w_0) \geq 1 - \xi_i \quad (\text{psi})$$

$$1 - \xi_i$$

allows some misclassifications

⇒ Hinge Loss ?

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

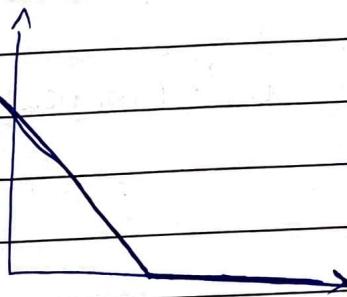
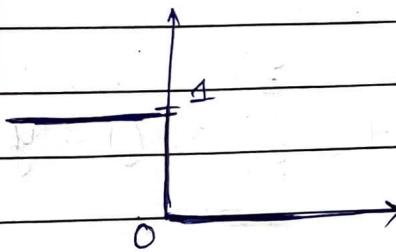
$$0 < 1 - y_i^o (w^T x_i^o + w_0)$$

Hinge

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

⇒ Zero-one loss

Hinge



⇒ Applying KKT conditions :-

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

(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

α_i^* for $i \in C^+$

α_j^* for $j \in C^-$

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

→ Any eg. for which dual var (α) 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:-

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

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

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

$$w^* = \sum_{i=1}^n \alpha_i^* y_i^* x_i^*, \quad w_0^* \text{ satisfies } \sum \alpha_i^* y_i^* = 0$$

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

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

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

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

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

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

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

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

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

Φ = feature func²

Date _____
Page _____

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

Mercer's theorem :-

Kernel func² :- $K: \mathcal{X} \times \mathcal{X} \rightarrow \mathbb{R}$

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 Φ .

If the matrix,

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

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.

④ 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

regularization parameter

① L1-SVM

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

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

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

② L2-SVM

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

③ L1C-SVM

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

④ L2C-SVM

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

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

Date _____
Page _____

Let $x = R$

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

$$= \exp\left(-\frac{x_i^{o2}}{2\sigma^2} - \frac{x_j^{o2}}{2\sigma^2} + \frac{2x_i^o x_j^o}{2\sigma^2}\right)$$

$$= \exp\left(-\frac{x_i^{o2}}{2\sigma^2}\right) \exp\left(-\frac{x_j^{o2}}{2\sigma^2}\right) * \exp\left(\frac{2x_i^o x_j^o}{2\sigma^2}\right)$$

$$= \exp\left(-\frac{x_i^{o2}}{2\sigma^2}\right) \exp\left(-\frac{x_j^{o2}}{2\sigma^2}\right) \left[1 + \frac{x_i^o x_j^o}{\sigma^2} + \frac{(x_i^o x_j^o)^2}{2!(\sigma^2)^2} + \dots \right]$$

$$\exp\left(-\frac{x_i^2}{2\sigma^2}\right) \exp\left(-\frac{x_j^2}{2\sigma^2}\right) \left[1, \frac{x_i^o}{\sigma}, \frac{x_i^{o2}}{\sqrt{2}\sigma^2}, \dots \right] \left[1, \frac{x_j^o}{\sigma}, \dots \right]$$

$$= \exp\left(-\frac{x_i^{o2}}{2\sigma^2}\right) \left[1, \frac{x_i^o}{\sigma}, \frac{x_i^{o2}}{\sqrt{2}\sigma^2}, \dots \right]$$

in terms of x_i^o
in infinite
dimensions

$$\phi(x_i^o)$$

13/2/23

separating & supporting
hyperplanes

SVC - default = RBF

Date _____

Page _____

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 in:-

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

'etc'

(SVD.pdf.) - Sir's notes.

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

(diagonal matrix)

(defn of
eigen vector
&
eigen vector)

$$A v = \lambda v$$

\downarrow vectors

λ = 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

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

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

- Linear comb^o of vectors
- Basis (size of basis?)
- Null space, column space (or range space)

Ax

$$\begin{bmatrix} | & & & \\ d & a^2 & \cdots & \\ | & | & & \end{bmatrix}$$

$$x = \begin{bmatrix} x_1 \\ x_2 \\ \vdots \\ x_d \end{bmatrix}$$

$$\sum_{j=1}^d x_j^0 (a^j)$$

↑
scalar

vectors

→ Rank of matrix

$$Ax = 0$$

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

$$\dim(\text{range space}) + \dim(\text{null space}) = \text{total dimension}$$

→ Every $n \times n$ real symmetric matrix

→ eigenvalues are real & eigen vectors can be chosen

orthonormal

$$Av = \lambda v$$

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

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

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

(linear algebra book) - ref,

Date _____
Page _____

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

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

→ (difference)

$\{u\}$ is 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

$n \times n$	real	symmetric

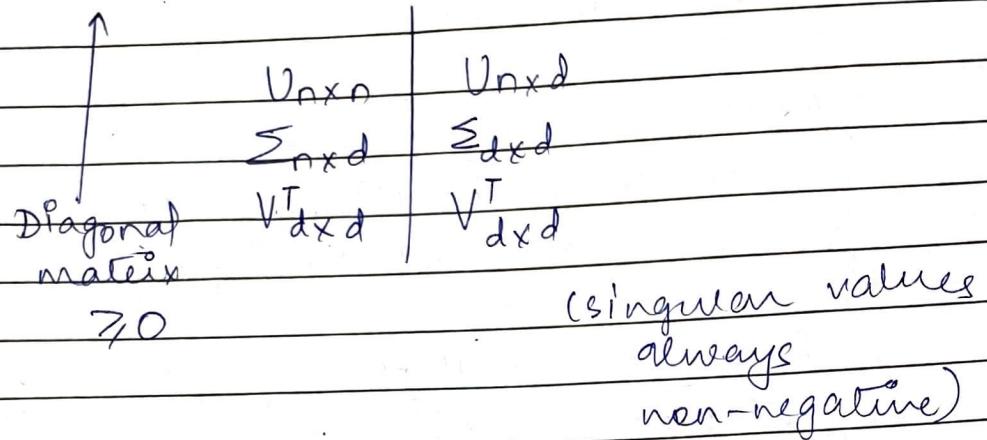
Matrices

$n \times n$ matrices

SVD

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

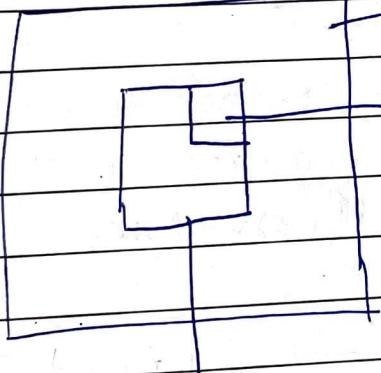
$$A = U \Sigma V^T$$



$$U^T U = I_n$$

$$V^T V = I_d$$

20/2/23



SVD (contd.)

Spectral theorem

→ Columns of U and V

are orthonormal

Σ 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 $-A^T A$.

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

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

Multiplying by A on both sides,

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

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

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

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

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

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

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

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

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

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

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

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

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 $|a_i \cdot v|$

i. Sum of length of squared projections is $\|Av\|_2^2$.

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

written in terms of matrix A

Maximizing $\|Av\|_2^2$ is best fit line.

First singular vector of A

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

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

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

$$\text{Why, } v_2 = \arg \max_{V \perp V_1, \|V\|_2=1} \|AV_2\|_2$$

Similarly find v_3, v_4, \dots, v_d

$$\sigma_2(A) = \|A v_2\|_2$$

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

Imp result:- Let A be $n \times d$ matrix where v_1, \dots, v_r 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