

(\*) constrained optimization problem:

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

$$\min f_0(x)$$

such that

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

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

classmate

Date \_\_\_\_\_

Page \_\_\_\_\_

8th. Feb 23

Lagrangian form of C.O.P

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

$$\|Ax - b\|_2^2 \rightarrow \|Ax - b\|_2^2 + \lambda \|x\|_2^2$$

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

$$\lambda > 0$$

$\lambda_i^*, \mu_j^*$  is Lagrangian multiplier.

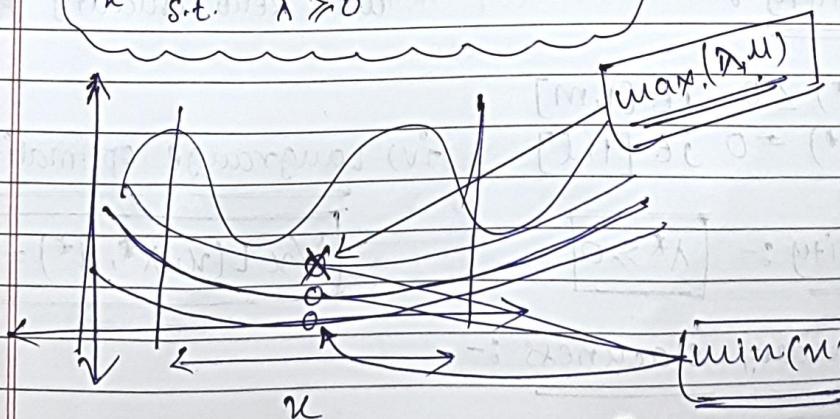
$$L(x, \lambda, \mu) = f_0(x) + \sum_{i=1}^m \lambda_i^* f_i(x) + \sum_{j=1}^l \mu_j^* f_j(x)$$

what if  $x$  is feasible point?

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

$$\min_x L(x, \lambda, \mu) \leq f_0(x^*)$$

s.t.  $\lambda \geq 0$



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

8th  
Feb  
2023

$$\min_{\lambda} L(\lambda, \lambda, u) = \underline{g(\lambda, u)}$$

4.

classmate

Date \_\_\_\_\_

Page \_\_\_\_\_

dual function.

$$\max_{\lambda, u} g(\lambda, u) \leq f_0(u^*)$$

$\lambda > 0$

$$g(\lambda, u) = \min_{\lambda} L(\lambda, \lambda, u)$$

$\Rightarrow (\lambda^*, u^*)$  are point of maximum, minimum

$$g(\lambda, u) \leq g(\lambda^*, u^*)$$

$$g(\lambda, u) \leq g(\lambda^*, u^*) \leq f_0(u^*)$$

maximum value of the dual function is equal to less than or equal to primal function.

$$f_0(u^*) - g(\lambda^*, u^*) \geq 0 \quad \left. \begin{array}{l} \text{weak duality} \\ \text{theorem} \end{array} \right\}$$

duality gap.

$\Rightarrow$  Primal Problem is convex problem.

(usually)

(not always)

strong duality hold

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

constraint should satisfy some condition

(\*) KKT conditions :-

$u^*$   $\leftarrow$  point of primal optimal

(i) Primal Feasibility :-

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

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

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

iv) Lagrange optimality :-

(ii) Dual feasibility :-

$$\lambda^* \geq 0$$

$$\nabla_{\lambda} L(\lambda, \lambda^*, u^*) = 0$$

(iii) Complementary slackness :-

$$\lambda_i^* f_i^*(u^*) = 0$$

8th Feb  
2023

Any tuple  $(x, \lambda, u)$  is called KKT point

$$R^n \uparrow \quad R^m \uparrow \quad R^l$$

classmate

Date \_\_\_\_\_

Page \_\_\_\_\_

if it satisfies the KKT conditions

when you have  $(x^*, \lambda^*, u^*) \Rightarrow$  KKT point.

if duality gap = 0

$x^*$  } primal optimal

$\lambda^*$  } dual optimal

$u^*$

$(x^*, \lambda^*, u^*)$  is KKT point

If Problem is convex

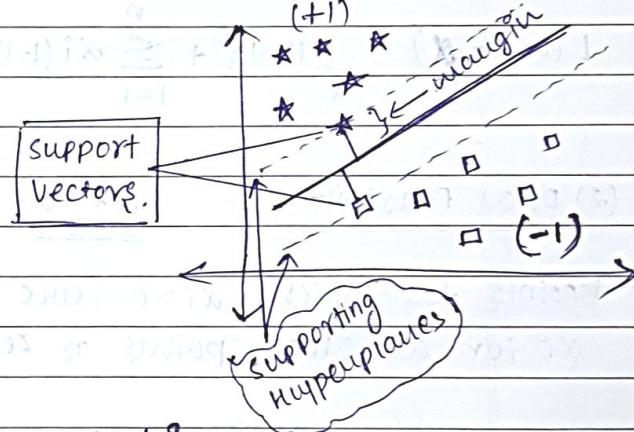
9th Feb. 2023

(\*) Hard Margin SVM :-

$$\text{Min}_{(w, w_0)} \frac{1}{2} \|w\|_2^2$$

s.t

$$y_i(w^T x_i + w_0) \geq 1 \quad \forall i \in [1, n]$$



How new data point will be classified?

$\Rightarrow x_{\text{new}} \leftarrow$  new data point.

$\rightarrow$  if  $(w^T x_{\text{new}} + w_0) \geq 1$  then positive class.  
else negative class ..

How to detect misclassification?

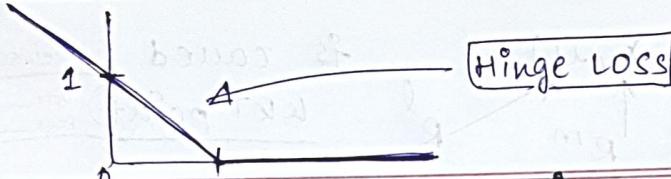
if  $y_i(w^T x_i + w_0) < 1$  then point is misclassified...

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

Hinge

$$\text{Hinge}[(x_i, y_i), (w, w_0)] = \max(0, 1 - y_i(w^T x_i + w_0))$$

9th Feb.  
2023



classmate

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

$$\text{Min}_{(w, w_0)} \frac{1}{2} \|w\|_2^2$$

s.t.  $y_i(w^T x_i + w_0) \geq 1 \quad \forall i \in [1, n]$

Primal form of  
Hard Margin.

= KKT conditions :-

(1) Primal Feasibility :-

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

Primal variables,

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

(2) Dual Feasibility :-  $\alpha_i \geq 0 \quad \forall i \in [1, n]$

→ points for which  $\alpha_i > 0$  are the support vectors.

$\alpha_i$  for all other points is zero.

(3) Complementary Slackness :-

$$\alpha_i (1 - y_i(w^T x_i + w_0)) = 0 \quad \forall i \in [1, n]$$

(4) Lagrange Optimality

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

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

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

9<sup>th</sup>  
Feb  
2023

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

classmate

Date \_\_\_\_\_

Page \_\_\_\_\_

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

$$\begin{aligned} \alpha_i &\text{ for } i \in C^+ \\ \alpha_j &\text{ for } j \in C^- \\ \sum \alpha_i &= \sum \alpha_j \end{aligned}$$

- Support Vectors :-

$$\{x_i | \alpha_i^* > 0\}$$

$$\min_{(w, w_0)}$$

$$\frac{1}{2} \|w\|_2^2$$

$$\text{such that } 1 - y_i(w^T x_i + w_0) \leq 0 \quad \forall i, i \in [1, n]$$

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

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

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

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

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

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

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

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

$$\text{Max } g(\alpha) = \sum_{i=1}^n \alpha_i - \frac{1}{2} \sum_{i=1}^n \sum_{j=1}^n \alpha_i \alpha_j y_i y_j x_i^T x_j$$

s.t.  $\alpha_i \geq 0$        $\uparrow$       dual form of Hard margin

$$\text{H} \cdot \sum_{i=1}^n \alpha_i y_i = 0$$

10th Feb 23

### (\*) Kernel SVM :-

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

$$x \rightarrow z \in \mathbb{R}^n$$

$$\text{s.t. } [n > d]$$

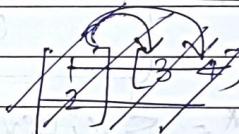
$$x_i \Rightarrow z_i \Rightarrow \phi(x_i)$$

$$\text{label}(x_{\text{test}}) = \text{sign} \left[ w^* \cdot \phi(x_{\text{test}}) + w_0^* \right]$$

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

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

- How to implement SVM in higher dimension without actually getting  $\phi(x)$  :-



Mercer's Theorem :-

Kernel fun :-

$$K: X \times X \rightarrow \mathbb{R}$$

A symmetric function  $K(x_i, x_j) = K(x_j, x_i)$

can be represented as  $K(x_i, x_j) = \phi(x_i)^T \phi(x_j)$   
for some ( $\phi$ ).

IFF,

$$K = \begin{pmatrix} K(x_1, x_1) & K(x_1, x_2) & \dots & K(x_1, x_n) \\ K(x_2, x_1) & K(x_2, x_2) & \dots & K(x_2, x_n) \\ \vdots & \vdots & \ddots & \vdots \\ K(x_m, x_1) & K(x_m, x_2) & \dots & K(x_m, x_n) \end{pmatrix}$$

Gram matrix.

$(m \times n)$

• Example of kernels :-

- Linear Kernel :-

$$K(x_i^o, x_j^o) = x_i^o \cdot x_j^o$$

classmate

Date \_\_\_\_\_

Page \_\_\_\_\_

10th Feb 2023

- Polynomial Kernel :-

$$K(x_i^o, x_j^o) = (1 + x_i^o \cdot x_j^o)^t ; t > 0$$

- Gaussian Kernel (RBF Kernel) :-

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

NOTE! A necessary and sufficient condition to make our custom kernel is that the Gram Matrix should be positive definite.

(\*) Soft SVM :-

$$\text{L1-SVM} :- \underset{(w, w_0, \xi)}{\text{Min}} \quad \lambda \frac{\|w\|_2^2}{2} + \frac{1}{2} \sum_{i=1}^n \xi_i$$

regularization term

$$\text{L2-SVM} :- \underset{(w, w_0, \xi)}{\text{Min}} \quad \frac{\lambda}{2} \|w\|_2^2 + \frac{1}{2} \sum_{i=1}^n \xi_i^2$$

$$\text{L2C-SVM} :- \underset{(w, w_0, \xi)}{\text{Min}} \quad \frac{C}{n} \sum_{i=1}^n \xi_i + \frac{\lambda}{2} \|w\|_2^2$$

$$\text{L2C-SVM} :- \underset{(w, w_0, \xi)}{\text{Min}} \quad \frac{C}{n} \sum_{i=1}^n \xi_i^2 + \frac{\lambda}{2} \|w\|_2^2$$

Example :- Let  $x \in \mathbb{R}$ .

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

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

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

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

10th Feb.  
2023

$$= \exp\left(-\frac{x_i^2}{\sigma^2}\right) \cdot \exp\left(-\frac{x_j^2}{\sigma^2}\right) \cdot \left[ 1 \quad \frac{x_i}{\sigma} \quad \frac{x_i^2}{\sqrt{2}\sigma^2} \dots \right] \left[ 1 \quad \frac{x_j}{\sigma} \quad \frac{x_j^2}{\sqrt{2}\sigma^2} \dots \right]$$

$$= \exp\left(-\frac{x_i^2}{\sigma^2}\right) \left[ 1 \quad \frac{x_i}{\sigma} \quad \frac{x_i^2}{\sqrt{2}\sigma^2} \dots \right] \cdot \exp\left(-\frac{x_j^2}{\sigma^2}\right) \left[ 1 \quad \frac{x_j}{\sigma} \quad \frac{x_j^2}{\sqrt{2}\sigma^2} \dots \right]$$

$\Phi(x_i)$

$\Phi(x_j)$

(\*) (\*)

(\*) Low Rank Approximation :-

$$\min_{A \in \mathbb{R}^{n \times d}} \|A - A'\|_F^2$$

subject to  $\text{rank}(A) = k$

$$\|Ux\|_2 = \|x\|_2$$

$$\|Vx\|_2 = \|x\|_2$$

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

solution

↑ Orthonormal.

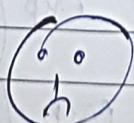
(\*) (\*) (\*) SMOTE :- Synthetic Minority Over-Sampling Technique.

- Avoid overfitting due to exact replicas of minority class samples.
- Subset of minority class is taken.
- New synthetic data samples ~~are~~ created from subsets and added.



No overfitting.

Reduced overfitting



Mainly introduce noise.

(\*) AdaBoost :- ~~A~~ classification  
 Boosting :- combining 'base' classifiers to form a "committee" whose performance is "better" than any of base classifiers.

classmate

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

Slide 8  
Deep learning  
RNN

It works even for "weak learners".

$\Rightarrow$  PCA, only considers features into account.

But in supervised learning we have label also.

20th April

(\*) Fisher's LDA :-

assumption:- Data coming from Gaussian distribn.

Fisher's

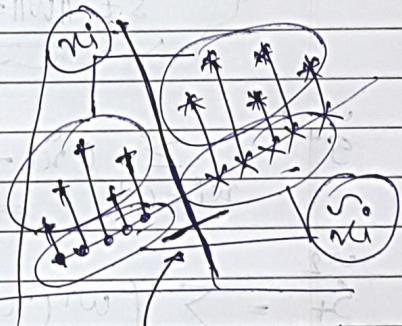
divide data into two classes

$$C^+ = \{x_i | y_i = +1\}, C^- = \{x_i | y_i = -1\}$$

$$|C^+| = n^+ \quad n^+ = n + n^-$$

$$|C^-| = n^-$$

Total no. of points



we need to find

$x_i$

projected

$\forall i \in C^+ \cup C^-$

mean of both classes.

$$m^+ = \frac{1}{n^+} \sum_{x_i \in C^+} x_i$$

$$m^- = \frac{1}{n^-} \sum_{x_i \in C^-} x_i$$

$$\tilde{m}^+ = \frac{1}{n^+} \sum_{x_i \in C^+} \tilde{x}_i$$

$$\tilde{m}^- = \frac{1}{n^-} \sum_{x_i \in C^-} \tilde{x}_i$$

$$\tilde{m}^+ = w^T m^+$$

$$\tilde{m}^- = w^T m^-$$

we need to maximize the distance between mean of two classes.

$$\therefore \text{maximize } \|m^+ - m^-\|$$

$$\therefore \text{maximize } |w^T(m^+ - m^-)|$$

$$\|m^+ - m^-\|^2 = \|w^T(m^+ - m^-)\|^2$$

$$= w^T(m^+ - m^-)(m^+ - m^-)^T w$$

$\uparrow$  matrix (rank = 1)

$S_B$  (Between class scatter matrix)

$$= w^T S_B w$$

{maximize  $w^T S_B w$ }

s.t.  $\|w\|=1$

problem

$$S^+ = \sum_{x_i \in C^+} (x_i - m^+)^2$$

$$S^- = \sum_{x_i \in C^-} (x_i - m^-)^2$$

$$S_t = \sum_{x_i \in C^+} (w^T(x_i - m^+))^2$$

scatter for projected classes.

$$S_t = \sum_{x_i \in C^+} w^T(x_i - m^+)(m^+ - m^+)^T w$$

$$= w^T \left( \sum_{x_i \in C^+} (x_i - m^+)(m^+ - m^+)^T \right) w$$

sum of rank 1 matrix.

$$S_t = w^T S^+ w$$

scatter matrix for +ve class.

$$\text{Similarly, } S^- = w^T S^- w$$

$$S^+ + S^- = w^T(S^+ + S^-)w$$

$$= w^T(S_w)w$$

within class scatter matrix.

# Fisher's Linear Discriminant criteria.

$$J(W) = \frac{W^T S_B W}{W^T S_W W}; R^d \rightarrow R$$

classmate

Date \_\_\_\_\_

Page \_\_\_\_\_

20/04/23

$$W_{\text{Fisher}} = \arg \max J(W)$$

$$\frac{\partial J(W)}{\partial W} = \frac{(W^T S_B W) \cdot (S_B W) - (W^T S_B W) (S_B W)}{(W^T S_W W)^2}$$

$$\left. \frac{\partial J(W)}{\partial W} \right|_{W^*} = W^T S_W W \cdot S_B W - W^T S_B W \cdot S_W W = 0$$

$$(W^T S_W W) S_B W = (W^* T S_B W) S_W W^* \quad J(W^*)$$

$$S_B W^* = \frac{(W^* T S_B W) S_W W^*}{(W^* S_W W^*)}$$

$$S_B W^* = \underbrace{J(W^*)}_{\text{scalar}} S_W W^*$$

$$\text{similar to } Av = \lambda Bv$$

$$S_B v = (m_t - m^-) (m_t - m^-)^T v$$

$$S_B v = \lambda (m_t - m^-) \quad \lambda$$

solving this eqn is  
generalised eigen value  
problem.

$$J(W^*) S_W W^* \propto (m_t - m^-)$$

$$S_W W^* \propto (m_t - m^-)$$

A and B are  
symmetric matrices

if  $S_W$  is invertible matrix, then

$$W^* \propto S_W^{-1} (m_t - m^-)$$

24<sup>th</sup> April

## \* Online learning \*

classmate

Date \_\_\_\_\_

Page \_\_\_\_\_

→ Traditional ML Methods assume that entire dataset  $D$  is available before training starts. In this case we perform "Batch learning!"

→ In many cases data arrives sequentially.

→ Let  $\hat{w}_{t-1}$  is our parameter for given data points from time  $1, 2, \dots, t-1$ . We want to update our parameter in constant time when we see the  $t^{\text{th}}$  point. We have to find rule.

$$y_t(-m + m)(-m + m) = y_t g_2 \\ \times (-m + m) f = y_t g_2$$

$$(-m + m) \times \hat{w}_{t-1} = y_t g_2 \\ (-m + m) \hat{w}_{t-1}$$