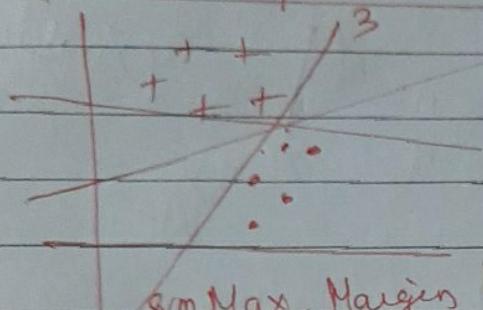


## Support Vector Machine (SVM)

- Vapnik in 1992
- work well on reasonably sized DS not on large DS

### I Optimal Separation



→ middle of separation  
b/w data points from 2 classes  
(equidistant)

Max. Margin Linear classifier

#### ① Margin & Support vectors.

- Measure the distance that we have to travel away from the line before we hit a point
- symmetric - The largest radius - Margin  $d_p(\text{datapt})$
- $\gamma$  → affects the speed at which it converges
- The datapoints in each class that lie closest to the classification line → support vectors

Best classifiers -

- 1) M should be as large as possible
- 2) SV are the most useful DP: after training we can throw all of the data except SV.

### Optimal decision boundary

weight vector, input vector  $x$

$$y = w \cdot x + b \rightarrow \text{bias weight} \rightarrow ax_1 + x_2 + b = 0 \quad \begin{matrix} x \rightarrow x_1 \\ y \rightarrow x_2 \end{matrix}$$

$$\rightarrow w \cdot x = \sum_i w_i x_i$$

the value  $\rightarrow$  '+' class

-ve value  $\rightarrow$  '0' class

for a given marginal value  $M$ ,

any point  $x$  where  $w^T x + b \geq M \rightarrow +$

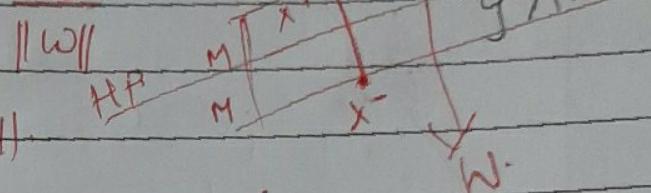
$$w^T x + b \leq -M \rightarrow 0$$

hyperplane  $\rightarrow w^T x + b = 0$

direction from  $x^+$  to  $x^-$  is along  $w$

$w$  is a unit vector

Margin is  $1/\|w\|$



To keep margin large ( $M$ ) length of vector  $x$   $\Rightarrow$  norm  
same as making  $w^T w$  small.

find boundary classifier  $\|x\| = \sqrt{x_1^2 + \dots + x_n^2}$   
well & also making  $w^T w$  a measure of  $w$  small.

$$\frac{1}{2} w^T w \text{ subject to } \|x\| \leq 1$$

$\Rightarrow$  absolute value of  $(w^T x + b)$

But we fail to distinguish the good

and bad hyperplane  $\Rightarrow$  adjust this metric

## ② Constrained Optimisation Problem:

- consider target answers  $+1, -1$

$$\min \frac{1}{2} w^T w \text{ subj to } t_i (w^T x_i + b) \geq 1 \text{ for all } i=1, \dots, n$$

problem of finding  $w$  &  $b$  is called Optimisation problem.

### SVM Lagrange problem

Lagrange state that if we want to find  $\min g$  under the equality constraint  $g$ , we just need to solve for

$$\nabla f(x) - \lambda \nabla g(x) = 0.$$

$$\frac{1}{2} w^T w$$

$\rightarrow$  lag. multiplier

$$t_p (w^T x_p + b) = 1$$

To find unique optimum

using Quad. prog  $\rightarrow$  corner

with linear contr.

unique minimum.

do this in poly time

when we find that ~~optimal~~ opt.  $\rightarrow$  KKT  
(Karush - Kuhn - Tucker) will be satisfied

$$\lambda_j^* (1 - t_p (w^T x_p + b^*)) = 0.$$

$$1 - t_p \leq 0.$$

$$\lambda_j^* \geq 0.$$

$\lambda_j^* \neq 0$ , so  $(1 - t_p) = 0 \Rightarrow$  true for SV.  
Lag. fn.

$$\begin{aligned} L(w, b, \lambda) &= \frac{1}{2} w^T w - \sum_{j=1}^n \lambda_j (t_p (w^T x_j + b) - 1), \\ &= \frac{1}{2} w^T w + \sum_{j=1}^n \lambda_j (1 - t_p (w^T x_j + b)), \end{aligned}$$

①

diff w.r.t  $w^T b$

$$\nabla_w L = w - \sum \lambda_i t_i x_i.$$

$$\frac{\partial L}{\partial b} \text{ or } \nabla_b L = \sum_{i=1}^n \lambda_i t_i$$

Set derivatives to zero to find maximum

$$w^* = \sum_{i=1}^n \lambda_i t_i x_i, \quad \sum_{i=1}^n \lambda_i t_i = 0 \quad (\text{saddle pts})$$

Substitute this to ①,

$$L(w^*, b^*, \lambda) = \sum_{i=1}^n \lambda_i t_i x_i - \sum_{j=1}^m \lambda_j t_j x_j$$

$$= \frac{1}{2} \sum_{i=1}^n \lambda_i t_i x_i \sum_{j=1}^m \lambda_j t_j x_j + \sum_{i=1}^n \lambda_i t_i \underbrace{\sum_{j=1}^m \lambda_j t_j}_{[w^T x_i]} - \sum_{i=1}^n \lambda_i t_i b.$$

$$= \sum_{i=1}^n \lambda_i - \frac{1}{2} \sum_{i=1}^n \sum_{j=1}^m \lambda_i \lambda_j t_i t_j x_i x_j$$

↓ This equation is dual problem  
max it w.r.t.  $\lambda_i$ . (Wolfe dual lag fn)

Constraints are  $\lambda_i \geq 0$  for all  $i$ .

$$\sum_{i=1}^n \lambda_i t_i = 0$$

then  $b^*$  is,

$$\text{WKT, } \nexists_{\mathbb{P}} (\omega \cdot x_p + b) = 1 > 0$$

$$\nexists_{\mathbb{P}} (\omega \cdot x_p + b) = 1$$

closest pt to HP will have functional margin 1.

$$\nexists_{\mathbb{P}} (\omega \cdot x_p + b) = 1$$

multiply both sides by  $t_p$ .

$$t_p^2 (\omega \cdot x_p + b) = t_p \quad \therefore t_p = 1.$$

$$\omega \cdot x_p + b = t_p$$

$$b = t_p - (\omega \cdot x_p)$$

$$b = \sum_{j \in SV} \left( t_j - \sum_{i=1}^n \lambda_i t_i x_i^T \cdot x_j \right)$$

$$= \frac{1}{N_S} \sum_{SVj} \left( \dots \right)$$

In case of error, it is not stable, aug it  $\oplus$

To classify new data point ( $z$ ),

$$\omega^{*T} z + b^* = \left( \sum_{i=1}^n \lambda_i t_i x_i \right)^T z + b^*$$

## Slack Variables for Non-linearly separable problem

- Non linear data, Slack var  $\eta_i \geq 0$ ,

so the constraint becomes,

$$t_p (w^T x_p + b) \geq 1 - \eta_p$$

for correct inputs, let  $\eta_i = 0$

SV  $\rightarrow$  tells one classifier  $\rightarrow$  mistake  $\rightarrow$  put on wrong side  
 another  $\rightarrow$  same but still longer side  
 includes this info in min. criteria  
 add a term  $w^T w + C \sum \eta_i$  (distance of misclassified point from BL)

if small  $C \rightarrow$  large margin with few errors to min.

$$L(w, \epsilon) = w^T w + C \sum_{i=1}^n \eta_i$$

dual problem  
inations, same except  $0 \leq \lambda_i \leq C$ .

$\forall i$  with  $\lambda_i > 0$ .

$\lambda_i$  condition,

$$\lambda_i^* (1 - t_p (w^T x_p + b^*) - \eta_i) = 0$$

$$(C - \lambda_i^*) \eta_i = 0 \rightarrow \eta_i \leq C, \text{ then}$$

$$\sum_i \lambda_i^* t_i = 0 \quad \eta_i = 0$$

$\lambda_i = C$ , if  $\eta_i \geq 1$ . then

Classifiers make a mistake  $\downarrow$  SV

## Slack Variables for Non-linearly separable problem

- Non linear data, Slack var  $\eta_i \geq 0$ ,

so the constraint becomes,

$$t_p (\mathbf{w}^T \mathbf{x}_p + b) \geq 1 - \eta_p$$

for correct inputs, set  $\eta_i = 0$

$\rightarrow$  SV  $\rightarrow$  tells one classifier  $\rightarrow$  mistake  $\rightarrow$  put on wrong side  
another  $\rightarrow$  same but still longer side

$\rightarrow$  include this info in min. criteria

$\rightarrow$  add a term  $\mathbf{w}^T \mathbf{w} + C \sum \eta_i$  (distance of misclassified point from BL)

if small  $C \rightarrow$  large margin with few errors  
fn to min is,

$$L(\mathbf{w}, \mathbf{e}) = \mathbf{w}^T \mathbf{w} + C \sum \eta_i$$

derivation, dual problem is same except  $0 \leq \lambda_i \leq C$ .

& SV with  $\lambda_i > 0$

KKT condition,

$$\lambda_i^* (1 - t_p (\mathbf{w}^* T \mathbf{x}_p + b^*) - \eta_i) = 0$$

$$(C - \lambda_i^*) \eta_i = 0 \rightarrow \lambda_i^* \leq C, \text{ then } \eta_i = 0$$

$$\sum_i x_i^* t_i = 0$$

$\downarrow$  if  $\lambda_i = C$ , if  $\eta_i \geq 1$ , then

Classifiers make a mistake SV

Kernels: new function  $\phi(x)$  of our input features

$$\omega^T x + b = \left( \sum_{i=1}^n \lambda_i t_i \phi(x_i) \right)^T \phi(z) + b,$$

Kernel function is defined as a function that corresponds to a dot product of two feature vectors in some expanded feature space

$$K(x_i, x_j) = \phi(x_i)^T \phi(x_j).$$

Polynomial of degree 2 then.

$$x_1 \dots x_n, x_1^2 \dots x_n^2, x_1 x_2, x_1 x_3 \dots x_{d-1} x_d.$$

$$\phi(x) \rightarrow d^2/2 \text{ elements. } \binom{n+n+\cancel{d(d-1)}}{2}$$

Using Kernel trick,  $O(d^2) \xrightarrow{\text{reduced}} O(d)$

↳ Gram Matrix / Kernel Matrix.

Kernels:

Linear kernel:  $K(x, y) = 1 + x^T y$ .

Poly. K:  $K(x, y) = (1 + x^T y)^k$ .  $k$ -degree

Sigmoid K:  $K(x, y) = \tanh(\beta_0 x^T y - \beta_1)$ .

RBF K:  $K(x, y) = \exp(-\|x-y\|^2/2\sigma^2)$ .

In general, fns that satisfy Mercer's theorem  
are Kernel fns.

## SVM Algo:

### Initialization

- for specified kernels & kernel param, compute kernel q distance between datapoints
- + compute  $K = X X^T$
- + for linear k, return  $K$ , for poly q degd return  $\frac{1}{C} K^q$ .
- + for RBF, compute  $K = \exp(-\|x - x'\|^2 / 2\sigma^2)$

### Training

- + assemble constraints as matrices

$$\min \frac{1}{2} x^T b_i b_j K x + q^T x \text{ subject to } b_i x \leq 1, Ax = b.$$

- + pass these matrices to solver.

- + identify SV and dispose rest of tr-data

- + Compute  $b^*$  using eqn

### Classification

- + for given test data z, use SV to classify data for relevant kernel using
  - compute inner prod. of test data & SV
  - perform classif as

$$\sum x_i t_p K(x, z) + b^* \text{ giving either belief or value (soft class).}$$

# 21

Friday • August  
Week 34 233-132

SUN

21 • 08 2014

$x_i$   $y_i$

9:00	1	1	-1				
9:30	2	2	-1	-④	④		
10:00	3	6	+1	1	2	3	4 5 6

$$\begin{aligned}
 L(w) &= \frac{w^2}{2} - \sum_{i=1}^3 \lambda_i (y_i (w \cdot x_i + b) - 1) \\
 &= \frac{w^2}{2} - \lambda_1 (-1 (w \cdot 1 + b) - 1) \\
 &\quad - \lambda_2 (-1 (w \cdot 2 + b) - 1) \\
 &\quad - \lambda_3 (+1 (w \cdot 6 + b) - 1).
 \end{aligned}$$

$$\begin{aligned}
 L(w) &= \frac{w^2}{2} - \lambda_1 (-w - b - 1) - \lambda_2 (-2w - b - 1) \\
 &\quad - \lambda_3 (6w + b - 1)
 \end{aligned}$$

$$\begin{aligned}
 \frac{\partial L}{\partial w} &= \frac{2w}{2} - \lambda_1 (-1) - \lambda_2 (-2) - \lambda_3 (6) \\
 &= w + \lambda_1 + 2\lambda_2 - 6\lambda_3.
 \end{aligned}$$

$$W = -\lambda_1 - 2\lambda_2 + 6\lambda_3 - \textcircled{1}$$



0

Weeks	27	28	29	30	31	August Weeks
Monday	6	13	20	27		Monday
Tuesday	7	14	21	28		Tuesday
Wednesday	1	8	15	22	29	Wednesday
Thursday	2	9	16	23	30	Thursday
Friday	3	10	17	24	31	Friday
Saturday	4	11	18	25		Saturday
Sunday	5	12	19	26		Sunday

September		October				
Weeks	16	17	18	19	20	21
Monday		7	14	21	28	
Tuesday	1	8	15	22	29	
Wednesday	2	9	16	23	30	
Thursday	3	10	17	24		
Friday	4	11	18	25		
Saturday	5	12	19	26		
Sunday	6	13	20	27		

October		Weeks				
Weeks	40	41	42	43	44	45
Monday		5	12	19	26	
Tuesday	6	13	20	27		
Wednesday	7	14	21	28		
Thursday	8	15	22	29		
Friday	9	16	23	30		
Saturday	10	17	24	31		
Sunday	11	18	25			

August + Saturday  
23-131 Week 34

22

22 \* 08 \* 2015

$$\frac{\partial L}{\partial b} = 0 - \lambda_1(-1) - \lambda_2(-1) - \lambda_3(1)$$

$$= +\lambda_1 + \lambda_2 - \lambda_3$$

$$\lambda_1 + \lambda_2 - \lambda_3 = 0, \quad \text{--- (2)}$$

$$\frac{\partial L}{\partial \alpha_1} = +w + b + 1 \quad \text{--- (3)} \quad \frac{\partial L}{\partial \alpha_3} = -6w - b + 1 \quad \text{--- (5)}$$

$$\frac{\partial L}{\partial \alpha_2} = 2w + b + 1 \quad \text{--- (4)}$$

$$\begin{array}{l|l} \text{--- (3)} & \text{--- (4) + (5)} \\ (3) = 0 & 2w + b + 1 = 0 \\ w + b + 1 = 0 & -6w - b + 1 = 0 \\ \hline w = 0 & -4w + 2 = 0 \\ & w = -2 \quad | -4 \\ & \boxed{w = 1/2} \end{array}$$

Sub in (3)

$$w + b + 1 = 0$$

$$0 + \frac{1}{2} + 1 = 0.$$

Sunday

23

$$b + \frac{3}{2} = 0.$$

$$\boxed{b = -3/2}$$

$$wx + b = 0$$

$$x - \frac{3}{2} = 0$$

$$x - \frac{3}{2} = 0$$

$$\boxed{x = 3}$$

Hggo!  
Time

24

4 • 08 • 2015

Monday • August  
Week 35 236-129

ω 4<sub>2</sub>

b  
-3<sub>2</sub>

o - note

July	27	28	29	30	31
Weeks					
Monday	6	13	20	27	
Tuesday	7	14	21	28	
Wednesday	1	8	15	22	29
Thursday	2	9	16	23	30
Friday	3	10	17	24	31
Saturday	4	11	18	25	
Sunday	5	12	19	26	

August	31	32	33	34
Weeks				
Monday	31	3	10	17
Tuesday	4	11	18	25
Wednesday	5	12	19	26
Thursday	6	13	20	27
Friday	7	14	21	28
Saturday	1	8	15	22
Sunday	2	9	16	23

$$g(x) = \omega x + b$$

$$x=1, g(1) = \frac{1}{2}(1) - \frac{3}{2} = \frac{1}{2} - \frac{3}{2} = -1 \quad (\text{we})$$

$$g(2) = \frac{1}{2}(2) - \frac{3}{2} = \frac{2}{2} - \frac{3}{2} = -\frac{1}{2} \quad (\text{we})$$

$$g(6) = \frac{6}{2} - \frac{3}{2} = \frac{6-3}{2} = \frac{3}{2} \quad (\text{the})$$

—d—