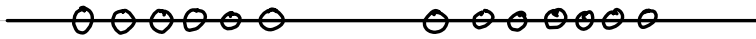


# SVM (Support Vector Machines)

- Used for both classification & regression
- Mainly used for classification of numeric data  
(Unlike logistic regression which acts on categorical data)
- Linear SVM → when data is linearly separable  
Nonlinear SVM → when data is not linearly separable
- Applications
  - Handwriting detection
  - Cancer diagnosis
  - Steganography detection
  - Face detection
  - Image classification
- Concept "Good fences makes good neighbors"

# Linear SVM

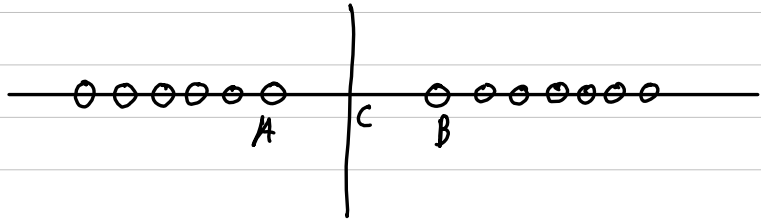
Consider classification of data points



In order to split these points, we use a threshold.

The threshold line can be the midpoint of the last points of the data

That is the point  $C$  is Midpoint of  $A$  &  $B$



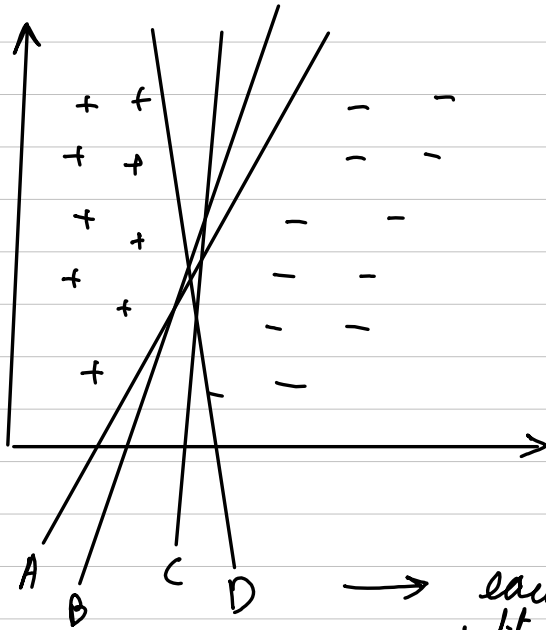
$A$  &  $B$  are the edge points of the data

They are called as support vectors

This form of classifier is called Maximal Margin classifier

For linearly separable data like that maximum margin hyperplane can be used.

For higher dimensions, an hyperplane is formed



Left of hyperplane is one class  
Right is another class

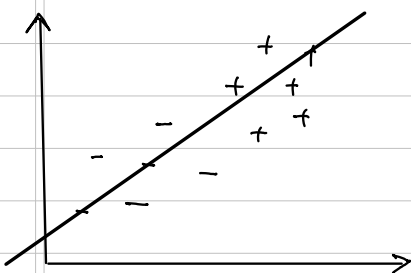
→ each hyperplane splits data into 2 parts

There are  $\infty$  such hyperplanes

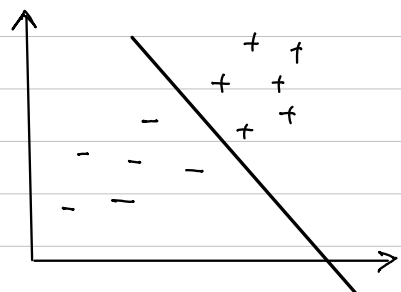
We need a hyperplane that not only performs well on the training data but also on unseen test data

So the plane which has maximum "margin" from the datapoints is required

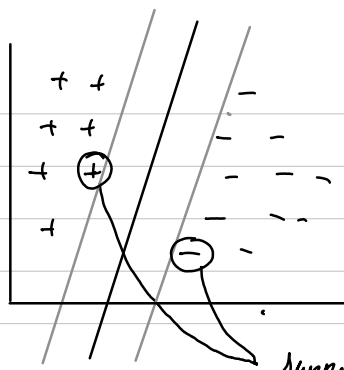
Note → This is different from linear regression



Linear regression



SVM



← Margin →  
support vectors  
ie values which lie on  
the boundaries

Classifier that has small margin is susceptible to overfitting

Hence we need to find plane with the Maximum margin

Optimize a hyperplane such that it has highest margin

Margin is distance of Margin with the nearest points of either class.

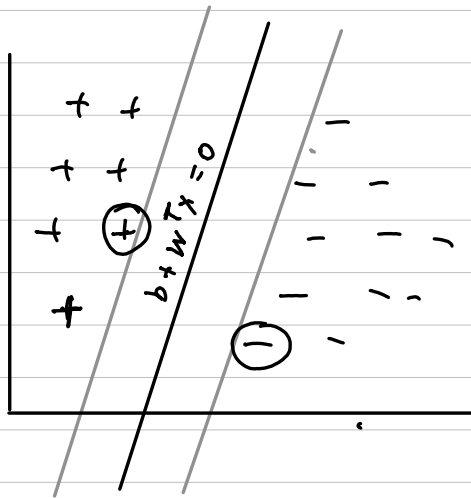
Intuitively classifier with small margin is more susceptible to model overfitting & classify with weak confidence on unseen data

Hence SVM is maximum margin classifier

Equation of hyperplane  $y = mx + c$  (2D)  
 $y = w_0 + w_1 x_1 + w_2 x_2 \dots$  (XD)

$$y = b + w^T x \quad (\text{matrix form})$$

Consider two Support Vectors  $x_+$  &  $x_-$



For all points above hyperplane

$$wx + b > 0 \quad \text{is class +}$$

For all points below hyperplane

$$wx + b < 0 \quad \text{is class -}$$

$$y = \begin{cases} + & \text{if } wx + b > 0 \\ - & \text{if } wx + b < 0 \end{cases}$$

We rotate the axes for hyperplane  
such that  $b + w^T x = 0$  for hyperplane  
( $b + w^T x \neq y$  now)

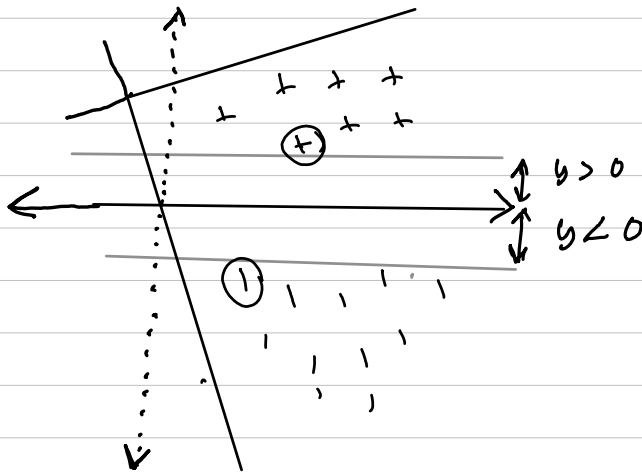
$$\begin{aligned} \text{For } x_+, \quad b + w^T x_+ &= k_1 \\ x_-, \quad b + w^T x_- &= k_2 \end{aligned}$$

Note that  $W$  &  $b$  are our numbers

We can scale them the way we want

They don't represent anything on the graph as long as the condition is met

From hyperplanes perspective the graph looks like this



In order to normalize things  
we will set scale as

$$WX + b = 1 \quad \text{at } x_+$$

$$WX + b = -1 \quad \text{at } x_-$$

$x_+ \rightarrow x_1$     $x_- \rightarrow x_2$  (new terminology)

Hence we get

$$w(x_1 - x_2) = 2$$

↑                      ↑  
matrix                  matrix of data  
of weights              points on boundary

Remember we want to maximize the margin that is  $(x_1 - x_2)$

$$\text{Hence we get Margin } d = (x_1 - x_2) = \frac{2}{\|w\|}$$

In order to Maximize Margin, we can minimize this function

$$\frac{\|w\|^2}{2} \quad \text{--- (1)}$$

We made it  $\|w\|^2$  by purpose

Also, we have condition

$$\begin{aligned} Wx_i + b &\geq 1 && \text{if class is } + \\ Wx_i + b &\leq -1 && \text{if class is } - \end{aligned}$$

Let us denote class + by  $y_i = 1$   
- by  $y_i = -1$

This  $y$  has no relation with any other  $y$   
(of axes)

Here  $y$  is the label (class) of the training dataset

$$\text{ie. } y_i (Wx_i + b) \geq 1 \quad - \quad (2)$$

for all datapoints in dataset

Note  $\rightarrow$  Here  $i$  are datapoints  $i = 1, 2, \dots, h$



So we want to minimize objective function<sup>①</sup>  
with respect to a condition<sup>②</sup>

This is called as convex optimization problem

Here the Objective function is Quadratic  $\|w\|^2$

Inequalities are linear.  $y_i(w \cdot x_i + b) \geq 1$

We can use the Lagrange multiplier method

LMM has got  $K \times K^T$  constraints

So Maximize  $\rightarrow$

$$L = \frac{\|w\|^2}{2} - \sum_{i=1}^n \lambda_i (y_i (w \cdot x_i + b) - 1)$$

$\lambda$  is Lagrangian multiplier for a datapoint

The KKT constraints are computationally expensive & can be solved using a linear / quadratic programming technique to get values of  $\lambda_i$

To classify a test tuple  $X$

$$\delta(X) = WX + b$$

$$= \sum_{i=1}^n \lambda_i y_i X_i + b$$

↑  
Matrix with  
length = length  $X$

$$b = \frac{b_1 + b_2}{2}$$

$$\begin{aligned} b_1 &= 1 - W \cdot X_1 \\ b_2 &= 1 - W \cdot X_2 \end{aligned} \quad \begin{array}{l} X_1 \text{ \& } X_2 \text{ are support} \\ \text{vectors} \end{array}$$

For non Support vectors  $\lambda = 0$

From  $k \times k^T$  condition we get

$$y_i (W X_i + b) - 1 \geq 0 \text{ for all datapoints } i$$

from this, we can solve and find values of  $\lambda$  as

$$y_i (W X_i + b) = 1 \quad \forall i$$

$$W X_i + b = y_i \quad \forall i \quad \text{since } \|y_i\| = 1 \text{ only sign is there}$$

$$(\sum \lambda_j y_j X_j) \cdot X_i + b = y_i$$

By augmenting  $b$  as additional dimension, we can remove  $b$  & get

$$(\sum \lambda_j y_j X_j) X_i = y_i$$

$$\sum \lambda_j y_j (X_j \cdot X_i) = y_i \quad \forall i$$

Unknown                      known   known   known

Since  $j = i$ , we have,  $i$  unknowns, &  $i$  equations

so solve to get  $\lambda$

for non support vectors  $\lambda = 0$  will be found

Q1)

Training data

	$A_1$	$A_2$	$y$	$\lambda_i$
1.	0.38	0.47	+	65.52
2.	0.49	0.61	-	65.52
3.	0.92	0.41	-	0
4.	0.74	0.89	-	0
	0.5	0.5	?	

→ 1. & 2. are support vectors

by quadratic programming we have

$$\begin{aligned} W_i &= \sum \lambda_i X_i y_i \\ &= 65.52 \times \begin{pmatrix} 0.38 \\ 0.47 \end{pmatrix} \times 1 \\ &\quad + 65.52 \times \begin{pmatrix} 0.49 \\ 0.61 \end{pmatrix} \times -1 \\ &\quad + 0 + 0 \end{aligned}$$

$$W = \begin{pmatrix} -7.20 \\ -9.17 \end{pmatrix}^T$$

$$\begin{aligned} b_1 &= 1 - W \cdot X_1 = 1 - \begin{pmatrix} -7.20 \\ -9.17 \end{pmatrix} \cdot \begin{pmatrix} 0.38 \\ 0.47 \end{pmatrix} \\ &= \frac{3.73}{5.30} = 0.7037 \end{aligned}$$

$$b_2 = 1 - W \cdot X_2 = \frac{4.528}{6.593} = 0.6867$$

$$b = \frac{b_1 + b_2}{2} = 10.0755$$

for 0.5, 0.5

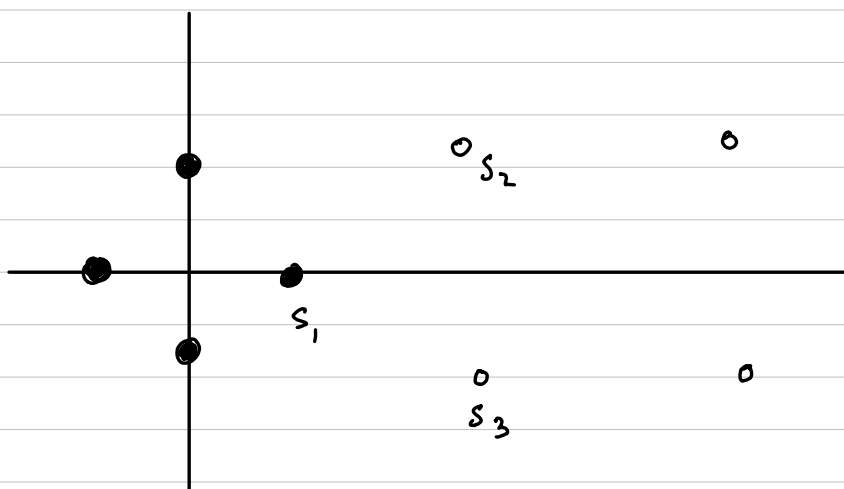
$$\begin{aligned} \delta(x) &= W \cdot X + b \\ &= \begin{pmatrix} -7.20 \\ -9.17 \end{pmatrix}^T \begin{pmatrix} 0.5 \\ 0.5 \end{pmatrix} + 10.0755 \end{aligned}$$

$$\begin{aligned} &= 1.8905 \\ \delta(x) &> 0 \quad \text{ie. class +} \end{aligned}$$

$$\textcircled{2} \quad \left. \begin{pmatrix} 3, 1 \\ 3, -1 \\ 6, 1 \\ 6, -1 \end{pmatrix} \right\} +ve \quad \left. \begin{pmatrix} 1, 0 \\ 0, 1 \\ 0, -1 \\ -1, 0 \end{pmatrix} \right\} -ve$$

classify (4, 1)

→



Support vectors are  $S_1 = (1, 0)$   
 $S_2 = (3, 1)$   
 $S_3 = (3, -1)$

for putting  $x_0 = 1$   
 & finding bias Add extra 1

$$S_1 = \begin{bmatrix} 1 \\ 0 \\ 1 \end{bmatrix} \quad S_2 = \begin{bmatrix} 3 \\ 1 \\ 1 \end{bmatrix} \quad S_3 = \begin{bmatrix} 3 \\ -1 \\ 1 \end{bmatrix}$$

$$[\alpha_1, \alpha_2, \alpha_3] \begin{bmatrix} S_1 \\ S_2 \\ S_3 \end{bmatrix} \cdot \begin{bmatrix} S_1 \\ S_2 \\ S_3 \end{bmatrix} = y$$

ie.

$$\begin{aligned} \alpha_1 S_1 S_1 + \alpha_2 S_1 S_2 + \alpha_3 S_1 S_3 &= y_1 \\ \alpha_1 S_2 S_1 + \alpha_2 S_2 S_2 + \alpha_3 S_2 S_3 &= y_2 \\ \alpha_1 S_3 S_1 + \alpha_2 S_3 S_2 + \alpha_3 S_3 S_3 &= y_3 \end{aligned}$$

$\alpha_1, \alpha_2, \alpha_3 \rightarrow$  constants

$S_1, S_2, S_3 \rightarrow$  Matrices

Multiplication is dot product

$$\therefore \alpha_1 \begin{bmatrix} 1 \\ 0 \\ 1 \end{bmatrix} \begin{bmatrix} 1 \\ 0 \\ 1 \end{bmatrix} + \alpha_2 \begin{bmatrix} 3 \\ 1 \\ 1 \end{bmatrix} \begin{bmatrix} 1 \\ 0 \\ 1 \end{bmatrix} + \alpha_3 \begin{bmatrix} 3 \\ -1 \\ 1 \end{bmatrix} \begin{bmatrix} 1 \\ 0 \\ 1 \end{bmatrix} = -1$$

$$\therefore \alpha_1 \begin{bmatrix} 1 \\ 0 \\ 1 \end{bmatrix} \begin{bmatrix} 3 \\ 1 \\ 1 \end{bmatrix} + \alpha_2 \dots \dots \dots = 1$$

$$\therefore \alpha_2 \begin{bmatrix} 1 \\ 0 \\ 1 \end{bmatrix} \begin{bmatrix} 3 \\ -1 \\ 1 \end{bmatrix} + \alpha_2 \dots \dots \dots = 1$$

$$\begin{aligned} \therefore \quad 2\alpha_1 + 4\alpha_2 + 4\alpha_3 &= -1 & \text{--- (1)} \\ 4\alpha_1 + 11\alpha_2 + 9\alpha_3 &= 1 & \text{--- (2)} \\ 4\alpha_1 + 9\alpha_2 + 11\alpha_3 &= 1 & \text{--- (3)} \end{aligned}$$

On solving (1) (2) & (3)

$$\alpha_1 = -3.5$$

$$\alpha_2 = 0.75$$

$$\alpha_3 = 0.75$$

$$W = \sum \alpha_i S_i$$

$$= -3.5 \begin{bmatrix} 1 \\ 0 \\ 1 \end{bmatrix} + 0.75 \begin{bmatrix} 3 \\ 1 \\ 1 \end{bmatrix} + 0.75 \begin{bmatrix} 3 \\ -1 \\ 1 \end{bmatrix}$$

$$= \begin{bmatrix} 1 \\ 0 \\ -2 \end{bmatrix} \leftarrow \text{Weights}$$

$$\leftarrow b = -2$$

for point (4, 1)

$$Wx + b = \begin{bmatrix} 1 \\ 0 \end{bmatrix}^T \begin{bmatrix} 4 \\ 1 \end{bmatrix} - 2$$

$$= 2$$

$2 > 0 \therefore$  point (4, 1) is positive class

!! Wait  $W = \sum \lambda_i y_i S_i$  right?

what is  $\alpha$ ?

also our equations were

$$\sum \lambda_i y_i (x_i \cdot x_j) = y_i$$

Also as per kkt,  $\lambda > 0$

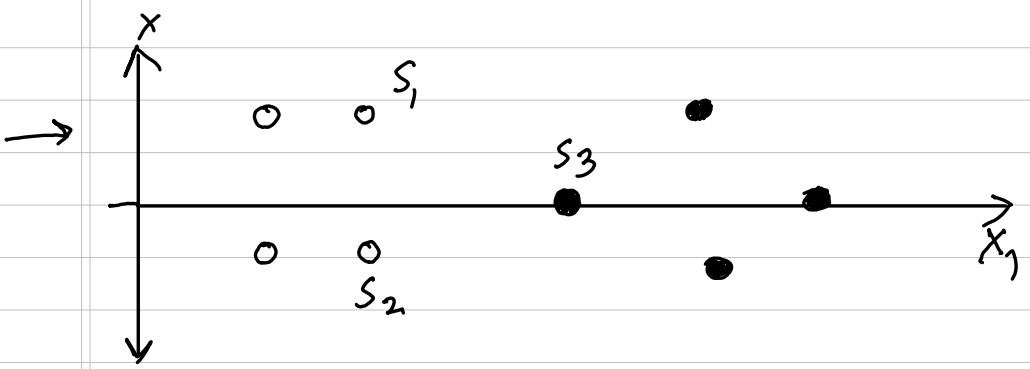
so basically we are clubbing

$$\alpha_i \rightarrow \lambda_i \cdot y_i$$

Classify the point (6,1) using SVM

③  $x_1$   $x_2$   $y$

1	1	-1
2	1	-1
1	-1	-1
2	-1	-1
4	0	1
5	1	1
6	0	1
5	-1	1



$$S_1 = \begin{pmatrix} 2 \\ 1 \\ 1 \end{pmatrix} \quad S_2 = \begin{pmatrix} 2 \\ -1 \\ 1 \end{pmatrix} \quad S_3 = \begin{pmatrix} 4 \\ 0 \\ 1 \end{pmatrix}$$

$$y_1 = -1 \quad y_2 = -1 \quad y_3 = 1$$

$$\alpha_1 S_1 S_1 + \alpha_2 S_1 S_2 + \alpha_3 S_1 S_3 = y_1$$

$$\alpha_1 S_2 S_1 + \alpha_2 S_2 S_2 + \alpha_3 S_2 S_3 = y_2$$

$$\alpha_1 S_3 S_1 + \alpha_2 S_3 S_2 + \alpha_3 S_3 S_3 = y_3$$

$$\alpha_1 \begin{pmatrix} 2 \\ 1 \\ 1 \end{pmatrix} \begin{pmatrix} 2 \\ 1 \\ 1 \end{pmatrix} + \alpha_2 \begin{pmatrix} 2 \\ 1 \\ 1 \end{pmatrix} \begin{pmatrix} 2 \\ -1 \\ 1 \end{pmatrix} + \alpha_3 \begin{pmatrix} 2 \\ 1 \\ 1 \end{pmatrix} \begin{pmatrix} 4 \\ 0 \\ 1 \end{pmatrix} = -1$$

$$\alpha_1 \begin{pmatrix} 2 \\ -1 \\ 1 \end{pmatrix} \begin{pmatrix} 2 \\ 1 \\ 1 \end{pmatrix} + \alpha_2 \begin{pmatrix} 2 \\ -1 \\ 1 \end{pmatrix} \begin{pmatrix} 2 \\ -1 \\ 1 \end{pmatrix} + \alpha_3 \begin{pmatrix} 2 \\ -1 \\ 1 \end{pmatrix} \begin{pmatrix} 4 \\ 0 \\ 1 \end{pmatrix} = -1$$

$$\alpha_1 \begin{pmatrix} 4 \\ 0 \\ 1 \end{pmatrix} \begin{pmatrix} 2 \\ 1 \\ 1 \end{pmatrix} + \alpha_2 \begin{pmatrix} 4 \\ 0 \\ 1 \end{pmatrix} \begin{pmatrix} 2 \\ -1 \\ 1 \end{pmatrix} + \alpha_3 \begin{pmatrix} 4 \\ 0 \\ 1 \end{pmatrix} \begin{pmatrix} 4 \\ 0 \\ 1 \end{pmatrix} = 1$$

$$6\alpha_1 + 4\alpha_2 + 9\alpha_3 = -1$$

$$4\alpha_1 + 6\alpha_2 + 9\alpha_3 = -1$$

$$9\alpha_1 + 9\alpha_2 + 17\alpha_3 = 1$$

$$\alpha_1 = -3.25$$

$$\alpha_2 = -3.25$$

$$\alpha_3 = 3.5$$

$$W_i = \sum \alpha_i S_i$$

$$= -3.25 \begin{pmatrix} 2 \\ 1 \\ 1 \end{pmatrix} + -3.25 \begin{pmatrix} 2 \\ -1 \\ 1 \end{pmatrix} + 3.5 \begin{pmatrix} 4 \\ 0 \\ 1 \end{pmatrix}$$

$$= \begin{pmatrix} 1 \\ 0 \\ -3 \end{pmatrix} \quad \} \quad w$$

$$b = -3$$

for point (6,1)

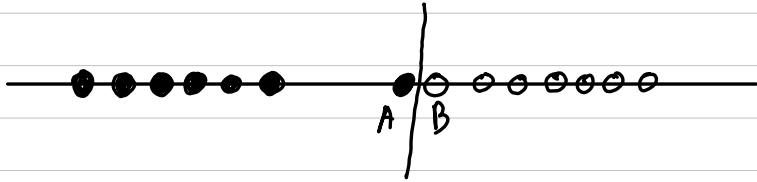
$w \cdot x$

( $w \cdot x + b$  here  $b$  is added dimension)

$$\begin{pmatrix} 6 \\ 1 \\ 1 \end{pmatrix} \begin{pmatrix} 1 \\ 0 \\ -3 \end{pmatrix} = 6 - 3 = 3$$

$3 > 0 \quad \therefore \text{class +ve}$

The maximum Margin classifier is the hard margin approach.

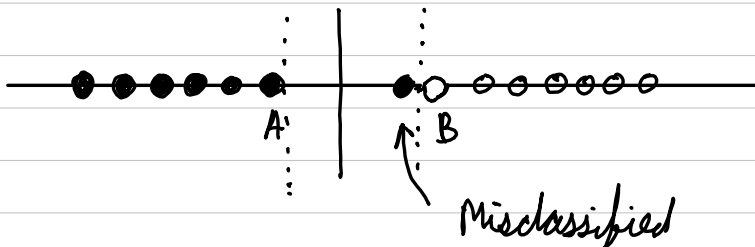


But that approach fails if we have an outlier

Hence a soft margin approach is used

Points are allowed to be Misclassified  
This is called support Vector classifier

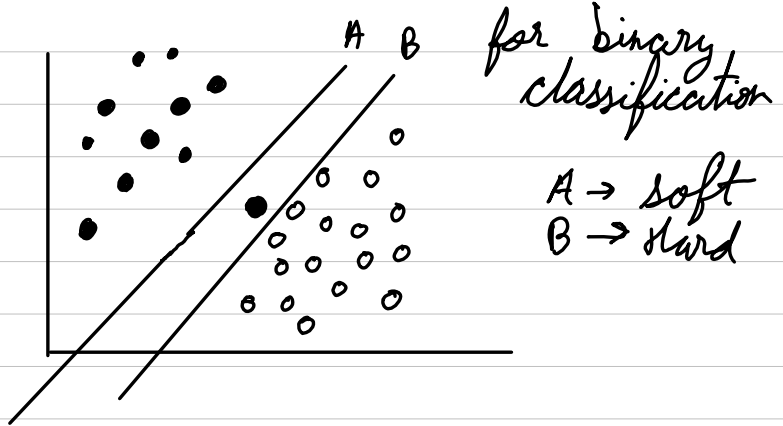
In order to provide variance to the dataset, the support vectors are not the edge points, but some points chosen by cross validation that are towards the edge



# Soft Margin SVM

Wider Margin is required, but also on training data

For 2D or  $n > 2$



For higher dimensional data, a hyperplane is formed

Misclassification must be allowed to reduce overfitting

Allow hyperplane to have errors



We need to penalize the misclassified points else model will break

Objective of soft Margin SVM

$$L = \frac{1}{2} \|w\|^2 + C (\dots \text{some math})$$

↑  
number of allowed mistakes

$C$  is a hyperparameter that decides tradeoff between maximizing margin and minimize mistake

$C$ : small  $\rightarrow$  focus more on avoiding misclassification  
large  $\rightarrow$  Maximize margin

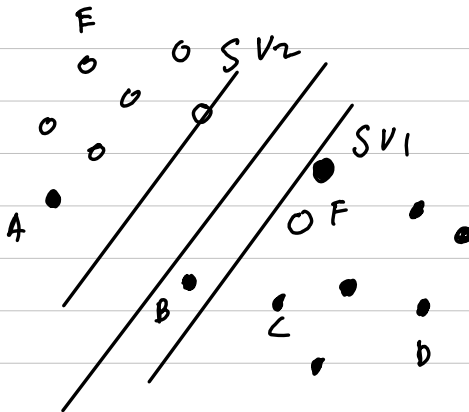
# Slack Variables

To check misclassified points, we put a variable  $\epsilon_i$

$\epsilon_i = 0$  if point correct

$0 \leq \epsilon_i \leq 1$  if point between Support Vectors

$\epsilon_i > 1$  if point on wrong side



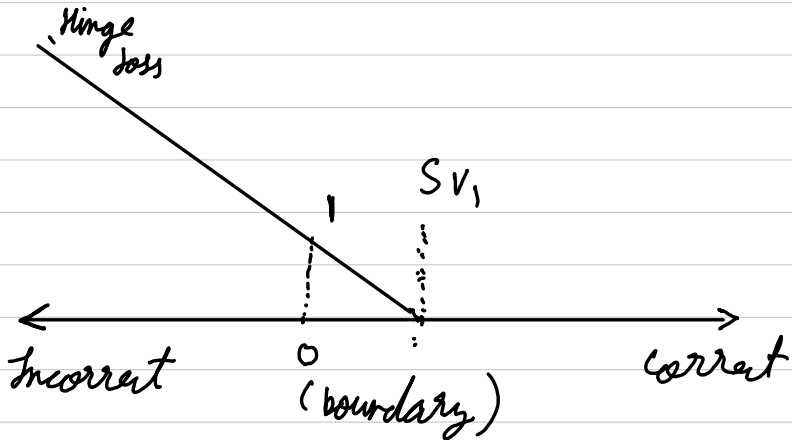
A, F  $\epsilon_i > 1$   
 B, G  $0 \leq \epsilon_i \leq 1$   
 C, D  $\epsilon_i = 0$

Hinge

$\max(0, 1 - \epsilon_i)$

# Hinge loss

$$\text{Hinge loss} = \max(0, 1 - y_i \bar{y}_i)$$



Example →

Actual $y_i$	Predicted $\bar{y}_i$	Hinge loss
+1	0.97	0.03
+1	1.20	0.00
+1	0	1.00
+1	-0.25	1.25
-1	-0.88	0.12
-1	-1.01	0.00

So we can write objective function as

$$L = \frac{1}{2} \|w\|^2 + C \sum_{i=1}^n (\xi_i)^\phi$$

$\phi \rightarrow$  parameter

& each  $x_i$  needs to be

$$y_i (w \cdot x_i + b) \geq 1 - \xi_i$$

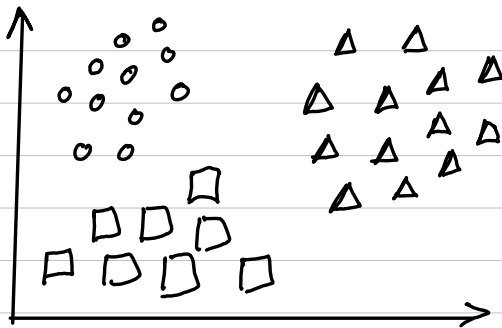
Now proceed with Lagrange multiplier

# Multiclass SVM

Till now we saw 2 class (binary) classification.

However even this approach has a drawback that we have restricted only to 2 classes.

What if we have 3 classes?



But SVM can only handle 2 classes at a time

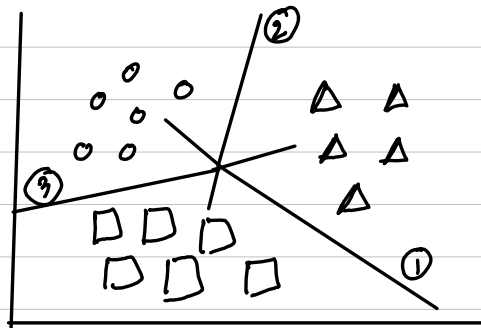
Then we need to use SVM multiple times

There are two strategies for this

OVO	One Vs One
OVA	One Vs All

## OVO (One vs One)

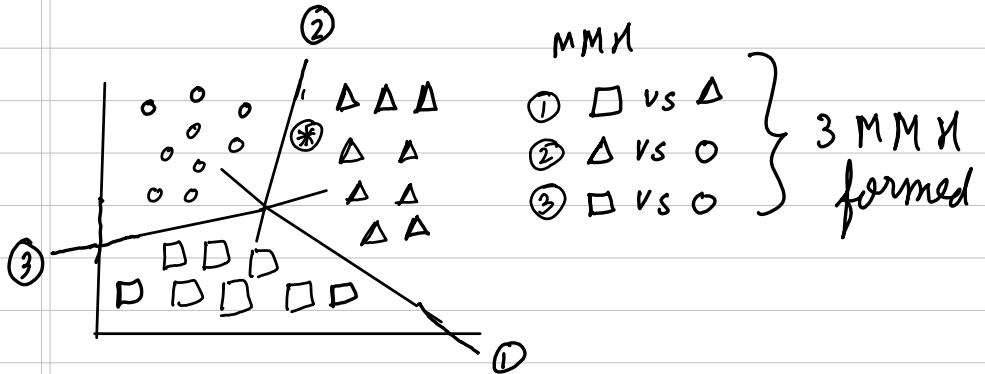
- Consider two classes at a time & ignore the rest. Run SVM only on the two classes. Do this for all pairs of classes
- Find MMH for each pair of classes.
- We can test OVO for each class and find the correct class.



- ① A vs B
- ② B vs C
- ③ C vs A

- But this is costly
- The amount of classification will increase exponentially with no of classes.

classify point marked (\*)



□ vs Δ → Δ  
 □ vs ○ → ○  
 Δ vs ○ → Δ

} Winner Δ

For 4 class

A B B  
 A C A  
 A D D  
 B C B  
 B D B  
 C D D

} B 6 MMH

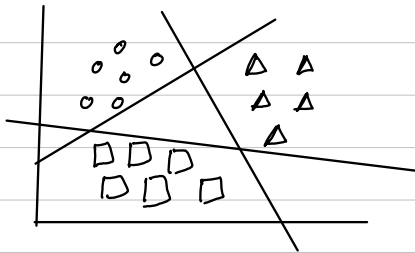
$$\text{No of MMH} = \frac{n(n-1)}{2} = O(h^2)$$

too large growth

## OVA (One Vs All)

Take one class, and group others into another. Do this for all classes.

Instead of making MMH for every pair of classes we will make MMH for every class.



Predicted label is the maximum of all labels.

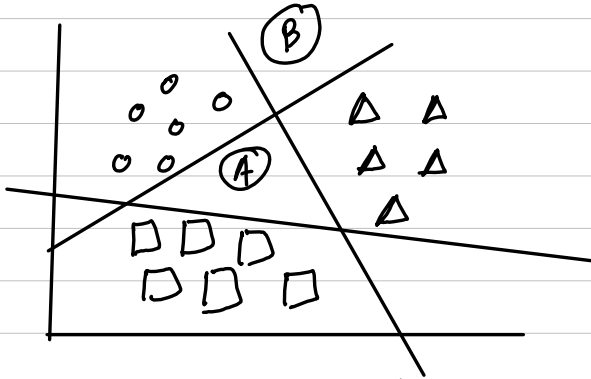
If classes are not linearly separable then this fails

This is a faster method of classification

Requires only  $n$  MMHs for  $n$  classes



Ambiguity in the classification



① → Unclassified tuples

② → In both classes ○ △

OVO & OVA can be thought as ensemble techniques where multiple models are combined and "winner take all" is used

# Error Correcting Output Codes

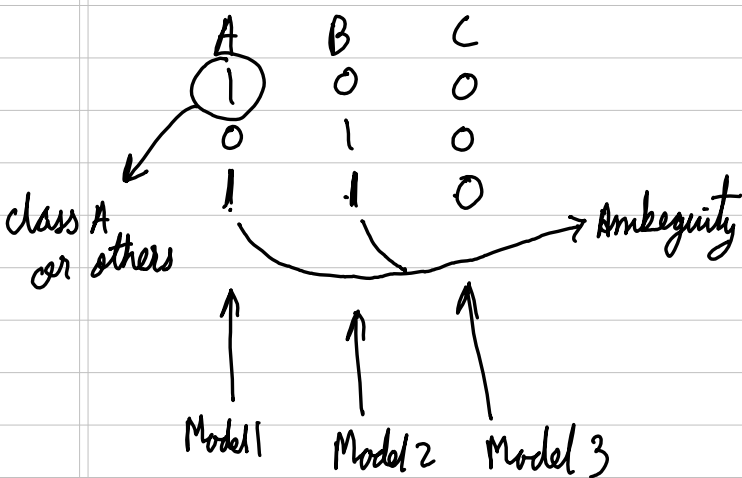
Technique that reframes multiclass classification as multiple binary classification.

A bit string (binary encoding) can be used to represent each class in the problem.

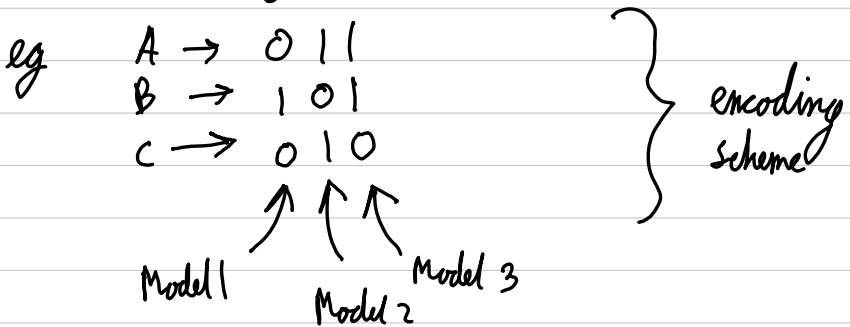
Each bit in the string can be predicted as 0 or 1 by a separate binary classification problem.

Each model receives the full input data and predicts one position in the output string.

Or A can be thought as one hot encoding



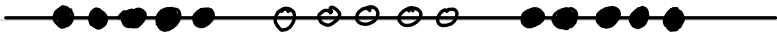
But EOC is more redundant and error correcting



Pred. 0 1 1  $\therefore$  class A  
0 0 0 closest to class  $\subseteq$

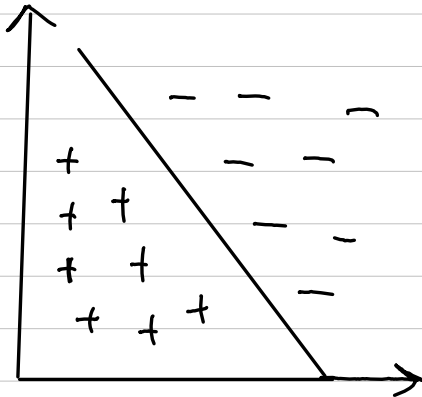
Care has to be taken that each encoded class has a very different encoded string

We saw  $MMH$  for linearly separable datapoints

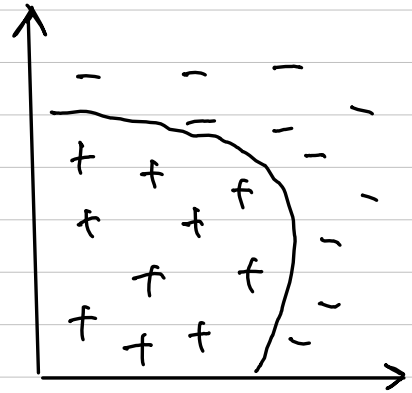


This type of data is not handled by the classifier because points are not linearly separable

In order to solve this problem, Non linear  $SVM_s$  are required

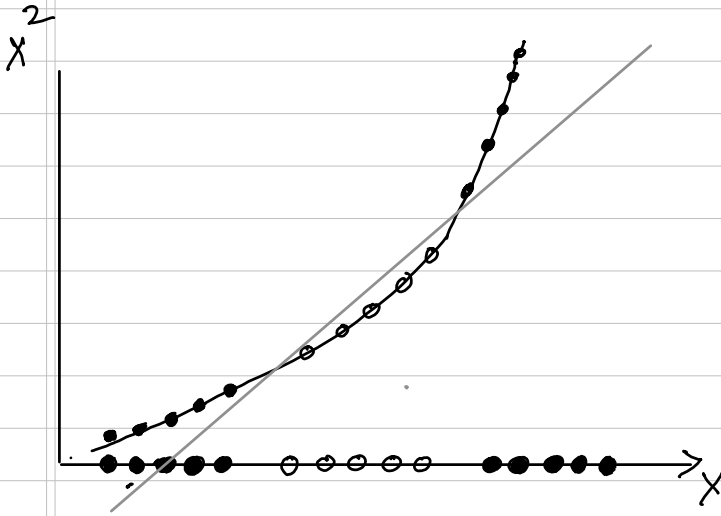


Linearly separable  
Data



Non linearly  
Data

what will happen if we square the data?



square the observations & make a 2d graph

x axis  $\rightarrow$  observations

y axis  $\rightarrow$  (observation)<sup>2</sup>

Now the data becomes linearly separable!

Now apply SVC to the graph

What we have done here is to apply some function on the input data to make it to a higher dimension

1D data which was not linearly separable  
 $\Updownarrow$  converted to  
2D data which is linearly separable

Here we can represent the function we applied as

$$\phi = \begin{pmatrix} x^2 \\ x \end{pmatrix} \text{ for input } (x)$$

$\uparrow$  2D plot  $x^2$  vs  $x$                        $\uparrow$  1D original  $x$

Here  $\phi$  is our transformation

This is a part of polynomial kernel

Why did we square the data and not cube it?

What if we can't draw a straight line even in 2D?

The value of the degree to which power is raised is  $d$

$d$  is decided by the SVM using cross validation

So basically we don't like the input space we are in and we want to transform the input space into a space where things are more convenient

Here we have raised the input space into higher dimension to make the transformation

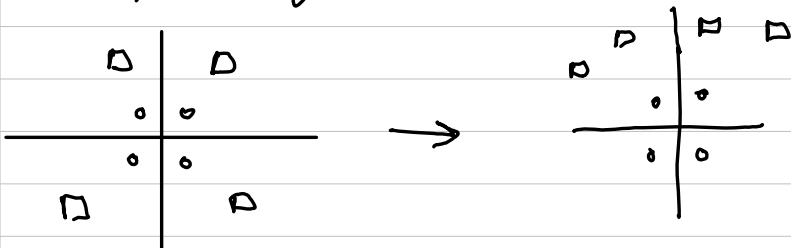
But there are other ways also

Different input spaces require different transformations

So we need a transformation  $\phi$

$\phi(\vec{x})$  will give transformed vector  $x$

$\phi$  may or may not increase the dimensions

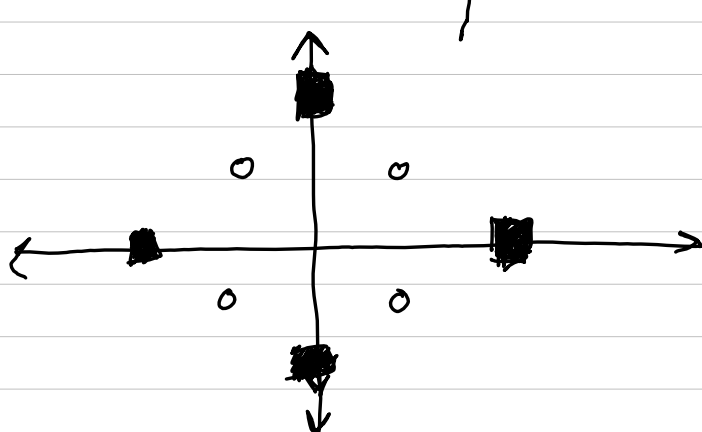


$$\textcircled{4} \quad \phi = \begin{cases} \begin{pmatrix} 6 - x_1 + (x_1 - x_2)^2 \\ 6 - x_2 + (x_1 - x_2)^2 \end{pmatrix} & \text{if } \sqrt{x_1^2 + x_2^2} \geq 2 \\ \begin{pmatrix} x_1 \\ x_2 \end{pmatrix} & \text{o/w} \end{cases}$$

Blue class vectors

Red class

$$\begin{pmatrix} 1 \\ 1 \end{pmatrix} \begin{pmatrix} -1 \\ 1 \end{pmatrix} \begin{pmatrix} -1 \\ -1 \end{pmatrix} \begin{pmatrix} 1 \\ -1 \end{pmatrix} \quad \left| \quad \begin{pmatrix} 2 \\ 0 \end{pmatrix} \begin{pmatrix} 0 \\ 2 \end{pmatrix} \begin{pmatrix} -2 \\ 0 \end{pmatrix} \begin{pmatrix} 0 \\ -2 \end{pmatrix} \right.$$



Apply transformation  $\phi$  to transform the space

$$\rightarrow \begin{pmatrix} 1 \\ 1 \end{pmatrix} \rightarrow \begin{pmatrix} 1 \\ 1 \end{pmatrix}$$

$$\begin{pmatrix} -1 \\ 1 \end{pmatrix} \rightarrow \begin{pmatrix} -1 \\ 1 \end{pmatrix}$$

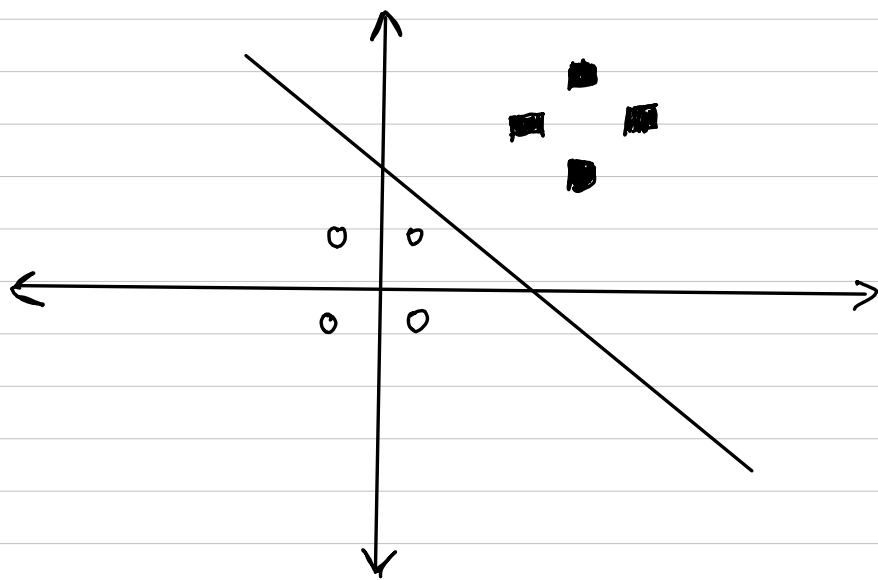
All blue will be same

$$\begin{pmatrix} 2 \\ 0 \end{pmatrix} \rightarrow \begin{pmatrix} 8 \\ 10 \end{pmatrix}$$

$$\begin{pmatrix} 0 \\ 2 \end{pmatrix} \rightarrow \begin{pmatrix} 10 \\ 8 \end{pmatrix}$$

$$\begin{pmatrix} -2 \\ 0 \end{pmatrix} \rightarrow \begin{pmatrix} 12 \\ 10 \end{pmatrix}$$

$$\begin{pmatrix} 0 \\ -2 \end{pmatrix} \rightarrow \begin{pmatrix} 10 \\ 12 \end{pmatrix}$$



Now solve as linear SVM

$$\bar{S}_1 = \begin{pmatrix} 8 \\ 10 \\ 1 \end{pmatrix} \quad \bar{S}_2 = \begin{pmatrix} 10 \\ 8 \\ 1 \end{pmatrix} \quad \bar{S}_3 = \begin{pmatrix} 1 \\ 1 \\ 1 \end{pmatrix}$$

$$\begin{aligned} \lambda_1 S_1 S_1 + \lambda_2 S_1 S_2 + \lambda_3 S_1 S_3 &= y_1 \\ \lambda_1 S_2 S_1 + \lambda_2 S_2 S_2 + \lambda_3 S_2 S_3 &= y_2 \\ \lambda_1 S_3 S_1 + \lambda_2 S_3 S_2 + \lambda_3 S_3 S_3 &= y_3 \end{aligned}$$

$$\begin{aligned} 165 \lambda_1 + 161 \lambda_2 + 19 \lambda_3 &= +1 \\ 161 \lambda_1 + 165 \lambda_2 + 19 \lambda_3 &= +1 \\ 19 \lambda_1 + 19 \lambda_2 + 3 \lambda_3 &= -1 \end{aligned}$$

$$\lambda_1 = 0.0859$$

$$\lambda_2 = 0.0859$$

$$\lambda_3 = -1.421$$

$$w = \sum \lambda \cdot S$$

$$= \begin{pmatrix} 0.1252 \\ 0.1252 \\ -1.2501 \end{pmatrix}$$

In order to reduce complexity, we can do

$$W = w / 0.1252$$

$$\therefore W = \begin{pmatrix} 1 \\ 1 \\ -9.98 \end{pmatrix}$$

$\therefore$  Hyperplane will be formed

The transformation here was 2D to 2D without adding any dimension

## Kernel Trick

So we are going to feed the transformed points  $\phi(x_i)$  to the SVM

It is found that every time applying  $\phi$  on every point is expensive

We need to optimize the process

But when we look at the internal structure of SVM, we notice that actually what is happening to the vectors we are providing is

for a test vector  $u$  ( $b$  is augmented)

$$y = W \cdot u$$

↑ Predicted class

Input

Pre calculated weights

Weights are calculated at training time

What if we don't do that?

$$W = \sum \lambda_i y_i x_i$$

step 1: found by solving equations

step 2: find actual value of  $W$

We skip step 2. & only calculate  $\lambda$

$$\therefore y_{\text{Pred}} = \left( \sum \lambda_i y_i x_i \right) \cdot u$$

$$y = \sum \lambda_i y_i (x_i \cdot u) \quad \text{--- (1)}$$

Now remember 1 equations?

$$\sum \lambda_i y_i (x_i \cdot x_i) = y_i \quad x_i \quad \text{--- (2)}$$

$$\left. \begin{aligned} \alpha_1 S_1 S_1 + \alpha_2 S_1 S_2 + \alpha_3 S_1 S_3 &= y_1 \\ \alpha_1 S_2 S_1 + \alpha_2 S_2 S_2 + \alpha_3 S_2 S_3 &= y_2 \\ \alpha_1 S_3 S_1 + \alpha_2 S_3 S_2 + \alpha_3 S_3 S_3 &= y_3 \end{aligned} \right\} \text{example}$$

from ① & ② we observe that everything depends only on dot products of the vectors

If we put transformed space vectors in the model, then in end we are going to get.

$$y = \sum \lambda_i y_i \phi(x_i) \phi(u)$$

$$\sum \lambda_i y_i \phi(x_i) \phi(x_j) = y_j$$

In place of  $x \cdot y \rightarrow \phi(x) \cdot \phi(y)$

So instead, can we replace the dot product with  $k(x, y)$ ?

$$k(x, y) = \phi(x) \cdot \phi(y)$$

We now need to know  $k$ . We don't need to know  $\phi$  anymore

$$\therefore \left. \begin{aligned} y &= \sum \lambda_i y_i k(x, y) \\ \sum \lambda_i y_i k(x, y) &= y_j \end{aligned} \right\} \text{New SVM}$$

So the function  $k$  is a function that provides the dot product of the vectors in another space

We don't need to know the transformation of vectors into another space

This is known as the kernel trick

Kernel Properties  $\rightarrow$

Mercer conditions must be satisfied

Addition of two kernels is a kernel



# Polynomial kernel

$$k(u, v) = (u \cdot v + r)^d$$

if  $r = 1/2$ ,  $d = 2$

$$(u \cdot v + 1/2)^2 = u \cdot v + u^2 \cdot v^2 + 1/4$$

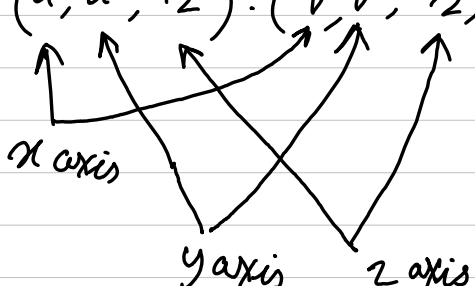
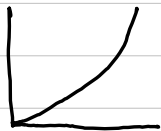
$$= (u, u^2, 1/2) \cdot (v, v^2, 1/2)$$


Diagram illustrating the mapping of  $u$  and  $v$  to their squared components and the constant term  $1/2$ . The  $u$  axis maps to  $u$  and  $u^2$ . The  $v$  axis maps to  $v$  and  $v^2$ . The  $z$  axis maps to the constant term  $1/2$  in both vectors.

Since  $z$  axis is same, ignore

$$(u, u^2) \cdot (v, v^2)$$



so this was what we did earlier

# $\phi$ vs $K$ Time complexity

Consider the last kernel

$$\begin{array}{cc} \phi & K \\ \begin{pmatrix} x^2 \\ x \end{pmatrix} \begin{pmatrix} u^2 \\ u \end{pmatrix} & (u \cdot v + 1/2)^2 \end{array}$$

$$= x^2 u^2 + x u$$

1  $\times$  multiplication  
1  $\times$  square

2  $\times$  squaring  
2  $\times$  multiplication

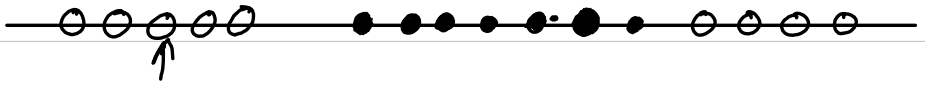
If the dimensionality was higher the growth would be massive for higher dimensions

Also kernels like radial basis operate in  $\infty$  dimensions, so they wouldn't be possible without 'kernel' trick

# Radial kernel

$$k(a, b) = e^{-\gamma(a-b)^2}$$

This kernel finds support vectors in  $\infty$  dimensions  
This can be found by expanding  $e$  into Taylor-like series



On data like these, the radial kernel acts as weighted nearest neighbor classifier to classify the points

$\gamma$  is found by crossvalidation

$\gamma$  controls the influence of the points

→ Advantages

Memory efficient

Works well for both low & high dimensional data

May Perform better than ANNs in some cases.

Works on text and images as well

→ Disadvantages

Long training time for large datasets

Doesn't work well on noise

less interpretability

Choosing right kernel is difficult

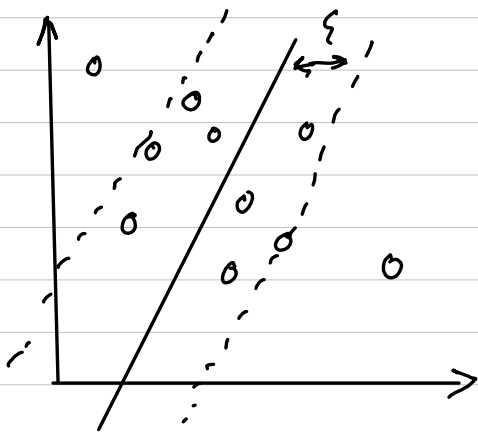
Doesn't give probability, directly the resultant class

# Support Vector Regression (SVR)

Used to handle complex non-linear data through kernels but hence slower

In SVM the goal is to find hyperplane that maximizes the margin between different classes.

But in SVR, the objective function is the one that deviates from the actual target values by a value no greater than  $\epsilon$



$\epsilon$ : tolerance value

The line must be as flat as possible

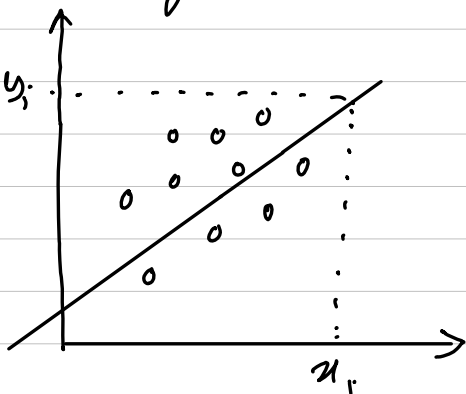
Epsilon insensitive loss  $\rightarrow$

$$\min_{w, b, \epsilon, \epsilon^*} \frac{1}{2} \|w\|^2 + C \sum_{i=1}^n (\epsilon_i + \epsilon_i^*)$$

Regularization parameter that determines the trade off between the flatness of the function & amount up to which deviations larger than  $\epsilon$  are tolerated.

$\epsilon_i \rightarrow$  upper boundary  
 $\epsilon_i^* \rightarrow$  lower boundary

Then the regression can be performed in same way as linear regression



find  $y_i$  gives  $x_i$

Kernelized SVR can be very accurate

