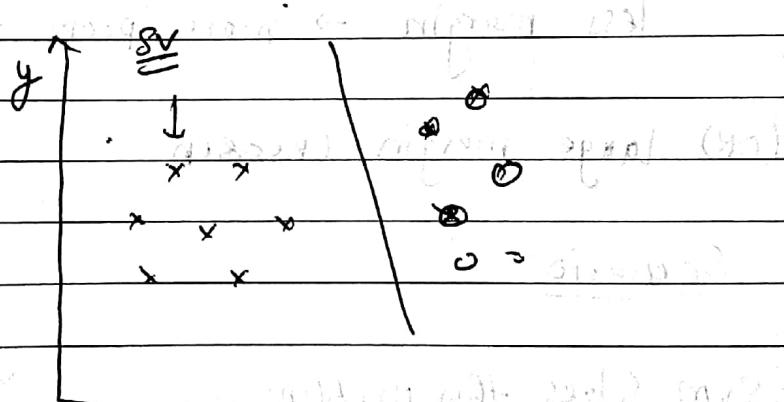


## Support Vector Machine (SVM)

classmate

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

- SVM is supervised ML Algo used for both classification & regression
- Mostly used for classification
- Support vector clustering (SVC) that also build on kernel functions but is used for unsupervised learning and Data mining.
- In this Algo, we plot each data item as a point in n-D space ( $n = \# \text{ features}$ ) with value of each feature being value of particular coordinate.
- Then perform classification by finding the hyperplane that differentiates classes very well.
- For multi-class classification for  $k$  class, we can perform  $k$ -time SVM with approach one-v-s rest classifier.



- Support vectors are simply the coordinates of individual observation.
- SVM is a frontier (line/border/boundary) which best segregates 2 classes (hyperplane/line).



### SVM properties :-

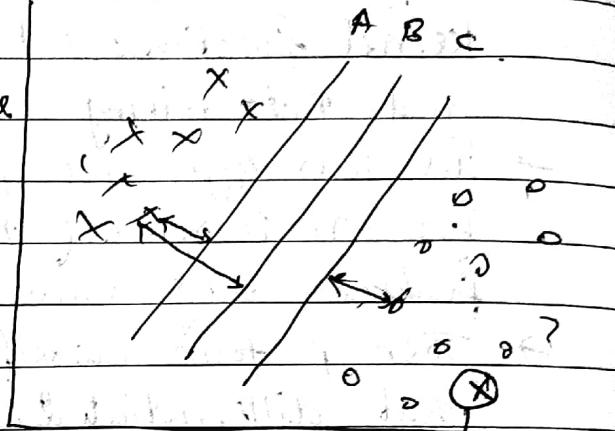
- Cost function - Convex ; Constraints - linear
- local minima is global minima -
- Solution of SVM is unique.

# Mathematics behind large margin (see previous)

## # How does SVM work?

How we can identify the right hyperplane?

→ Maximizing the distance b/w nearest data point and hyperplane help us to decide right hyperplane. This distance is called margin.



Q. → SVM has a feature to ignore outliers.

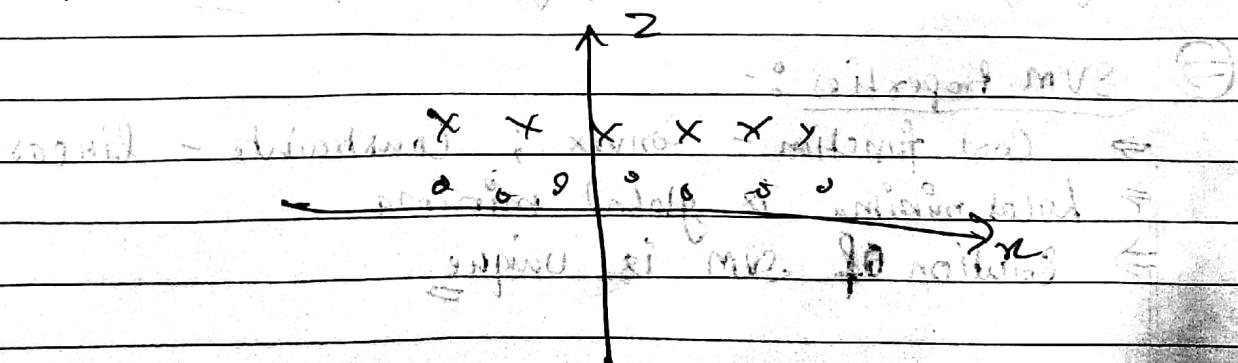
less margin  $\rightarrow$  more prone to overfitting

(OR) large margin chosen

- Scenario:

- SVM Solves this problem by introducing additional features,  
e.g.  $z = x^2 + y^2$ .

Now plot data point on axis x & z



- In SVM, easy to find linear hyperplane for 2 classes
- But for non-linear, we can't add additional feature manually.
- SVM have a technique called "Kernel" trick.
- Kernel :-  
These are the functions which takes low dimensional input space & transform it to a higher dimensional space. It converts non-separable problem to separable problem.  
- Useful in some complex data transformation

### ② Some standard Kernel :-

#### ① Polynomial (homogeneous) Kernel :-

$$k(\vec{x}_i, \vec{x}_j) = (\vec{x}_i \cdot \vec{x}_j)^d$$

↑                      ↑↑                      degree of  
 kernel function    vectors of              polynomial  
 feature space

#### ② Polynomial & non-homogeneous) Kernel :-

$$k(x, y) = (x^T y + c)^d$$

c = constant / free parameter

#### ③ Radial basis kernel function :- RBF kernel

- most popular kernel used to draw completely non-linear hyperplanes
- distance matrix, squared euclidean ~~tri~~ distance
- useful

$$M = \min_{i=1, \dots, m} y_i \left( \frac{w_0}{\|w\|} x + \frac{b}{\|w\|} \right) \rightarrow \text{margin for classmate}$$

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

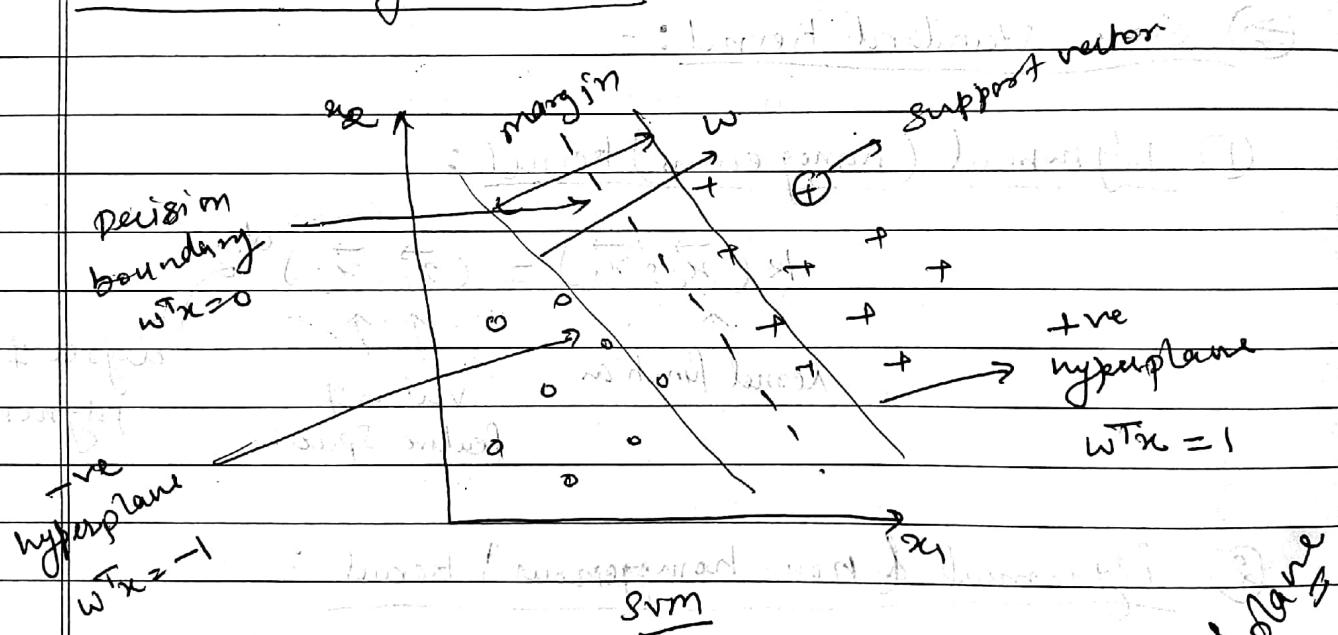
$$K(x, x') = \exp \left( -\frac{\|x - x'\|^2}{2\sigma^2} \right)$$

$x, x'$  → features of vector space.

$\sigma$  → free parameter.

- Selection of parameter is a critical choice.
- Using a typical value of the parameter can lead to overfitting our data.

## # Maximum margin intuition:



$$w_0 + w^T x_{\text{pos}} = 1 \quad \text{--- (1)}$$

$$w_0 + w^T x_{\text{neg}} = -1 \quad \text{--- (2)}$$

$$w^T (x_{\text{pos}} - x_{\text{neg}}) = 2$$

Some  
hyperplane

We can normalize this by length of vector  $w$ , which is

$$\|w\| = \sqrt{\sum_{j=1}^m w_j^2}$$

$$\frac{w^T (x_{\text{pos}} - x_{\text{neg}})}{\|w\|} = \frac{2}{\|w\|} \rightarrow \text{bc2 min } \frac{(w^T + b)}{\|w\|} \text{ for } \text{scale variant}$$

we are normalizing  
 $w^T + b$   
for var  
 $w^T w^T$

The left side of predicting  $y^{(i)}$  can be interpreted as distance b/w the +ve & -ve hyperplane which is also called margin that we want to maximize.

- Now objective function of SVM becomes maximization of this margin by maximizing  $\frac{2}{\|w\|}$  (width)
- Under the constraint that the samples are classified correctly, can be written as

$$w_0 + w^T x^{(i)} \geq 1 \text{ if } y^{(i)} = 1$$

$$w_0 + w^T x^{(i)} \leq -1 \text{ if } y^{(i)} = -1$$

This shows that -ve sample fall under -ve hyperplane

Compactly we can write as minimize  $J(w, w_0) = \frac{1}{2} \|w\|^2$

subject to

$$y^{(i)} (w_0 + w^T x^{(i)}) \geq 1 \quad \forall i = 1, 2, \dots, n$$

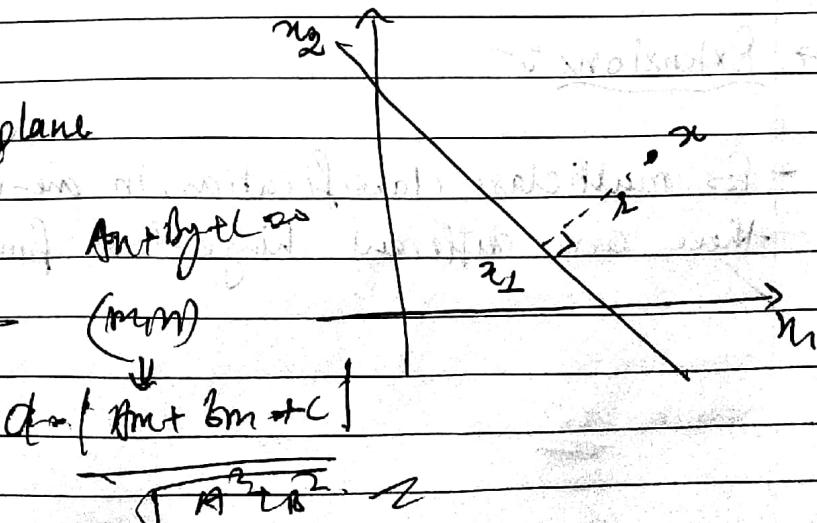
- In practice, it is easier to minimize the reciprocal term

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

Distance from  $x$  to plane

$$w^T x + b = 0$$

$$d = \frac{|w^T x + b|}{\|w\|}$$





## Hinge loss :-

- In ML, hinge loss is a function used to training classifiers.
- Hinge loss is used for "maximum-margin" classification most notably used for (SVMs).
- For out  $t = \pm 1$  & classifier score  $y$ ,  
hinge loss of prediction  $y$  is

$$l(y) = \max(0, 1 - t \cdot y)$$

Note:-  $y$  should be "raw" output of classifier's decision function, not predicted class label.

- For linear SVM,
- $y = w \cdot x + b$
- $w, b \rightarrow$  parameter of hyperplane
- $x \rightarrow$  point to classify.
- $t, y$  have same sign (meaning  $y$  predicts right class) and  $\|y\| \geq 1$ , the hinge loss  $l(y) = 0$  but when they have opposite sign,  $l(y)$  inc. linearly with  $y$  (one-sided-error).

## Extensions :-

- For multi-class classification. In one-vs-all or one-vs-one there are different hinge loss function proposed.

- Crammer & singh defined it for linear classification

$$\ell(y) = \max_{t \neq y} (0, 1 + w_t x - w_y y)$$

- Weston & Watkins provided a similar function, but with sum rather than a max.

$$\ell(y) = \sum_{t \neq y} (0, 1 + w_t x - w_y y)$$

- quadratically,

$$\ell(y) = \frac{1}{2Y} \max (0, 1 - t.y)^2$$

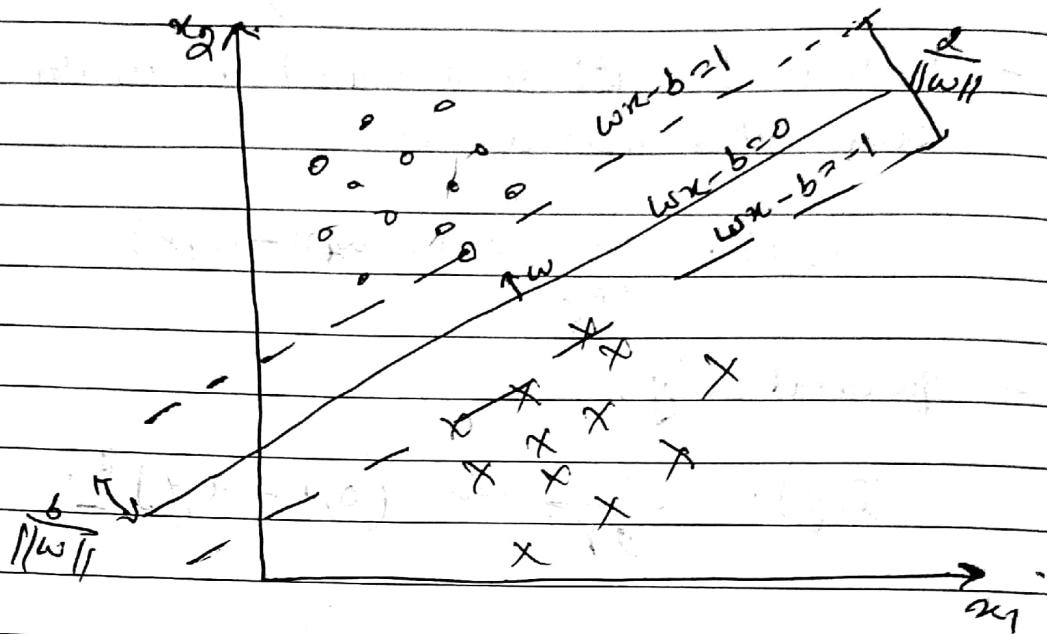


### How to choose a kernel:

- ① if # features are large as compared to # observations.  
↳ SVM with linear kernel of text classification with lots of words, small training example.
- ② if # features are small, # observation is intermediate  
↳ Gaussian Kernel
- ③ # features is small, # observations is small  
↳ linear kernel.

## → Selecting sum Hyperplanes :-

### ① Linearly separable :



For the data which can be separated linearly, we select 2 parallel hyperplane that separate 2 classes of data so that distance  $d_w$  between both the lines is max. The region b/w these 2 hyperplanes is known as "margin". & maximum margin hyperplane is one that lies in the middle of them.

$$w_i x_i - b \geq 1 \text{ if } y_i = 1$$

$$w_i x_i - b \leq -1 \text{ if } y_i = -1$$

$\vec{w}$  is normal vector to the hyperplane,  
 $y_i \rightarrow \text{classes}$ ,  $x_i \rightarrow \text{features}$

- The distance b/w 2 hyperplanes is  $\frac{2}{\|\vec{w}\|}$ , to

maximize this distance, ~~(or)~~  $\|\vec{w}\|$  should be minimized

$$\|\vec{w}\|_{\text{min}} \text{ for } y_i(\vec{w} \cdot \vec{x}_i - b) \geq 1, \forall i = 1, 2, \dots, n$$

(2) Non-linearly separable:- To build classifier for non-linear data, try to

minimize

Sub-gradient  
descent

$$f(\vec{w}, b) = \left[ \frac{1}{n} \sum_{i=1}^n \max(0, 1 - y_i(\vec{w} \cdot \vec{x}_i - b)) \right] + \lambda \|\vec{w}\|^2$$

Convex function Here,  $\max()$  method will be zero, if  $x_i$  is on the correct side of margin, for data that is on other side of margin, the function value is proportional to the distance from the margin.

$\lambda \rightarrow$  determine trade off b/w increasing the margin size

~~that~~

## # How to tune parameter of SVM :

→ Main 3 parameter that having higher impact on Model

Sklearn-SVM. SVC.C →  
① "kernel", "gamma", "C"

$\downarrow$   
rbf  
gamma → kernel coefficient for "rbf", "poly", "sigmoid",  
Higher the value of gamma, will try to exact fit the as per training data set i.e., generalization error and cause overfitting problem.

e.g.: gamma = 0, 10 (or) 100.

C → Penalty parameter, C of the error term. It also controls the tradeoff between smooth decision boundary & classifying training points correctly.

C=1, 100, 1000, ...

## # Pros & Cons of SVM :

→ Pros:-

- ① SVM is effective in high-dimensional spaces
- ② SVM is effective in case when # dimension > # Samples (as # features are very large)
- ③ Non-linear data is classified using customized hyperplane built by using kernel tricks.
- ④ SVM uses subset of training points in the decision function (called support vectors). so, it is also mem efficient.

→ Cons :-

- (1) biggest limitation is choice of kernel . The wrong choice of kernel can lead in error %.
- (2) SVM have good generalization performance but extremely slow in test phase.
- (3) with large no. of samples , SVM starts giving poor performance  $\rightarrow$  (large dataset).
- (4) SVMs don't directly provide probability estimate , these are calculated using an expensive five-fold cross validation .

## # Application of SVM :-

SVMs are by product of Neural Net.

- (1) facial expression classification :- Used to classify facial expression . It uses Statistical Model of shape & SVMs.
- (2) Speech recognition : Used to accept keywords & rejects non-keywords & build model to recognize speech .
- (3) Handwritten digit recognition
- (4) Text categorization

# SVC :-

Given training vector  $x_i \in \mathbb{R}^n$ ,  $i=1, \dots, n$  in 2 classes  
 $y \in \{-1, +1\}^n$

SVC solves following primal problem

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

$$\text{subject to } y_i (w^T \phi(x_i) + b) \geq 1 - \epsilon_i$$

$$\epsilon_i \geq 0, i=1, \dots, n.$$

its dual is

$$\min_{\alpha} \frac{1}{2} \alpha^T Q \alpha - e^T \alpha$$

$$\text{subject to } y^T \alpha = 0$$

$$0 \leq \alpha_i \leq C, i=1, \dots, n$$

$e$  = vector of all ones,  $C > 0$ , upper bound

$Q = n \times n$  semidefinite matrix

$$Q_{ij} \equiv y_i y_j K(x_i, x_j) \rightarrow \text{kernel} \quad K(x_i, x_j) = \phi(x_i)^T \phi(x_j)$$

training vectors are implicitly mapped to higher dimensional space mapped by function  $\phi$ .

The decision function:

$$\text{sign} \left( \sum_{i=1}^n y_i \alpha_i K(x_i, x) + p \right)$$

$$C = \frac{\# \text{Sample}}{\alpha}$$

↑  
Intercept

# SVR: Support Vector Regression

SVR solves following primal problem:

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

$$\text{subject to } y_i - w^T \phi(x_i) - b \leq \epsilon + \epsilon_i$$

$$w^T \phi(x_i) + b - y_i \leq \epsilon + \epsilon_i^*,$$

$$\epsilon, \epsilon^* \geq 0, i=1, \dots, n$$

Its dual is:

$$\min_{\alpha, \alpha^*} \frac{1}{2} (\alpha - \alpha^*)^T Q (\alpha - \alpha^*) + C e^T (\alpha + \alpha^*) - y^T (\alpha - \alpha^*)$$

$$\text{subject to } e^T ((\alpha - \alpha^*)) = 0$$

$$0 \leq \alpha_i, \alpha_i^* \leq C, i=1, \dots, n$$

Decision function

$$\sum_{i=1}^n (\alpha_i - \alpha_i^*) K(x_i, x) + p.$$

dual-coff: