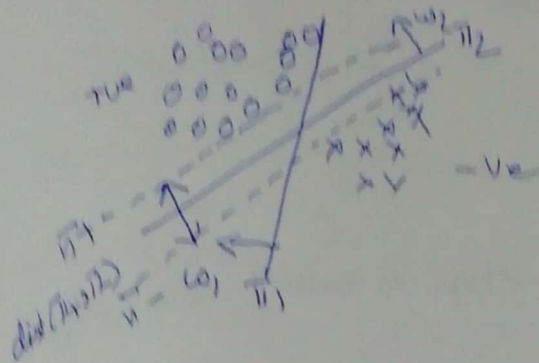


## Support Vector Machines (SVM)

→ SVM can do classification & regression.

→ we have bunch of negative & positive points



→ Key idea of SVM

we have to find the hyperplane ( $\pi$ ) that separates -ve & +ve pts as widely as possible.

$\pi_+$  &  $\pi_-$  are parallel to  $\pi$

$\pi_+$  &  $\pi_-$  are also parallel to each other

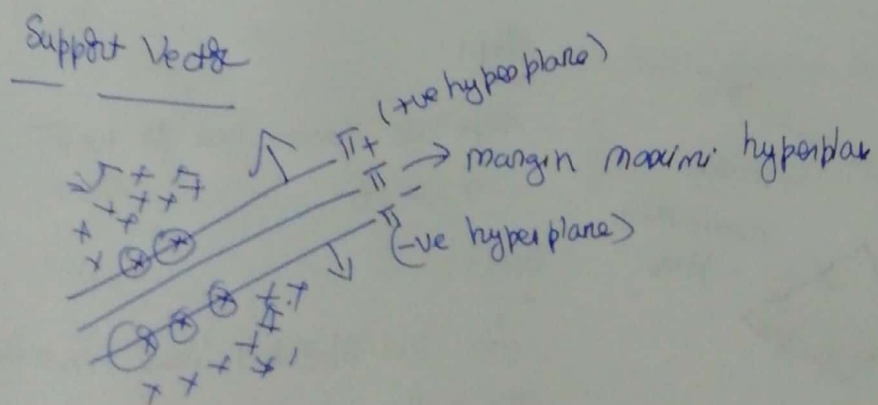
$\pi$ : margin maximizing hyper plane

SVM: Try to find a  $\pi$  that maximizes the margin =  $dist(\pi^+, \pi^-)$

margin  $\uparrow$ ; generalization acc  $\uparrow$

↳ Error on unseen data.

### Support Vector

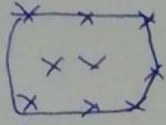


→ pts through which  $\pi^+$  &  $\pi^-$  pass through are support

vectors.

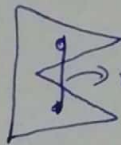
## Alternative geometry intuition of SVM

### Convex-hull



### Convex-polygon

A polygon which is passing through the shape is called convex polygon.



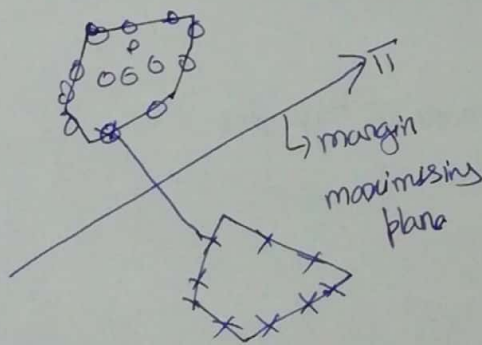
non convex polygon.

Convex polygon : If I want to connect the two points, the line connecting to both the points are inside the polygon is called

### Convex polygon

### Convex-hull

A convex hull is given a bunch of pts, if we can build the smallest convex polygon, that has every point inside the convex polygon or <sup>on</sup> inside the convex polygon.



#### Step 1

separate convex hull for the pts  
→ we pts

#### Step 2

find the shortest lines connecting the hulls

#### Steps

bisect the line  
(Equal parts)

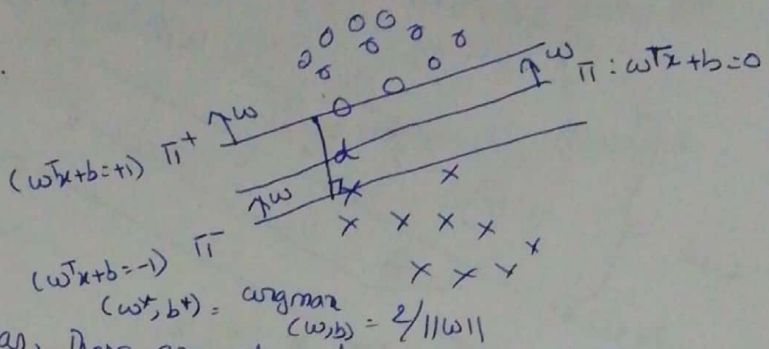
## 28.2 Mathematical derivation of SVM

$\Pi$  : margin-maximization.

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

$w$  is perpendicular to  $\Pi$

Since  $\Pi^+$  &  $\Pi^-$  are perpendicular, these are also perpendicular to  $\Pi$  so



$$\Pi^+ : w^T x + b = 1$$

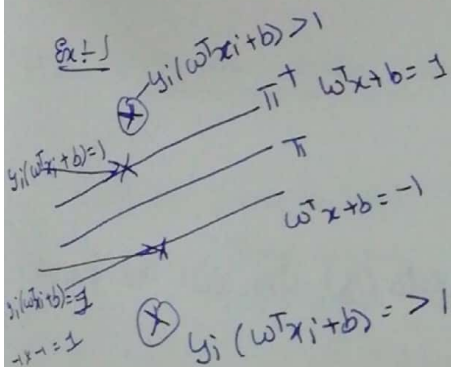
$$\Pi^- : w^T x + b = -1$$

Margin :

$$d = \frac{2}{\|w\|}$$

$$w^*, b^* = \text{argmax}_{(w,b)} \frac{2}{\|w\|} = \text{margin.}$$

hard margin SVM  
s.t.  $y_i (w^T x_i + b) \geq 1$  for all  $x_i$



Constraint

All the +ve points should be on one side & -ve points on the other side.

There should be any value in the margin area.

-> data should be linearly separable

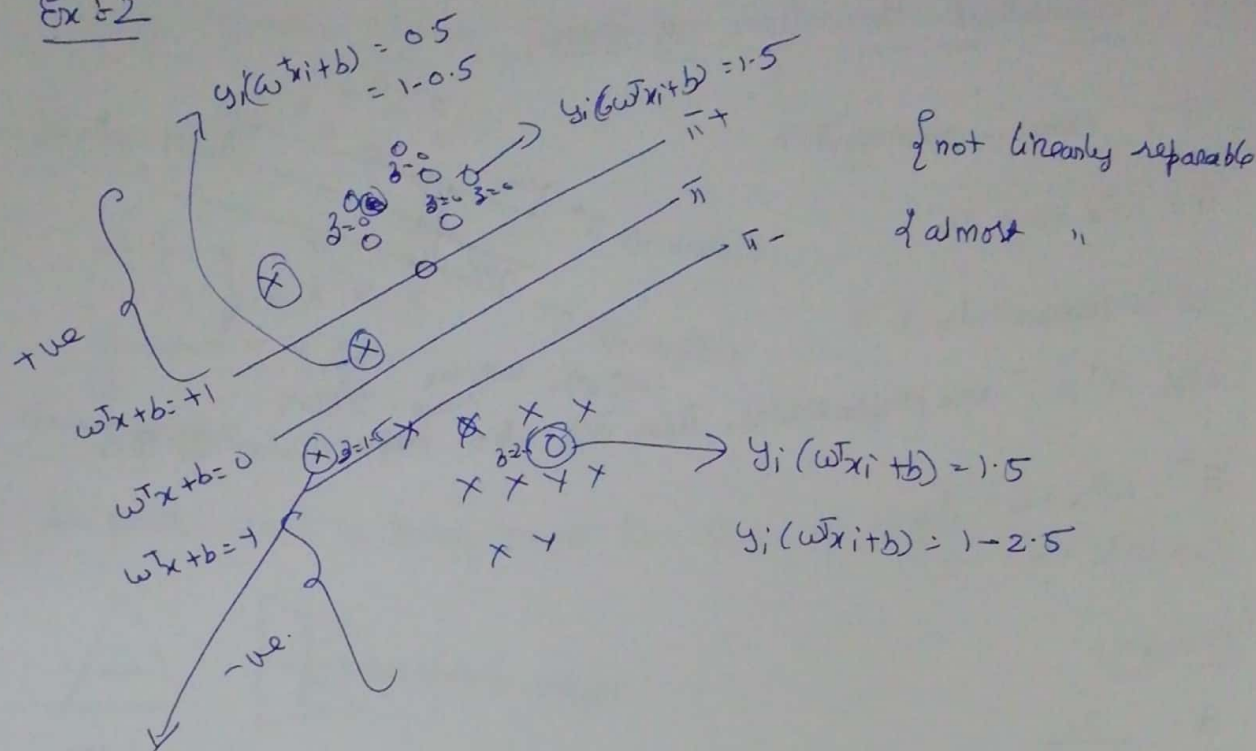
$$w^*, b^* = \text{argmax}_{w,b} \frac{2}{\|w\|}$$

s.t.  $\forall_i$   
for each

$$y_i (w^T x_i + b) \geq 1$$



Ex 2



half way b/w  $\pi^+$  &  $\pi^-$

$$\begin{aligned}
 y_i(w^T x_i + b) &= -0.5 \\
 y_i(w^T x_i + b) &= (1 - \underline{1.5}) \\
 &= \xi_i
 \end{aligned}$$

(zeta)

→ we will create a new variable called zeta ( $\xi$ ) for all the points such that

$\xi_i \geq 0$ , it is further away from correct  $\pi$  in the incorrect direction.

$$x_i \rightarrow \xi_i$$

$$\xi_i = 0 \quad \text{if } y_i(w^T x_i + b) \geq 1$$

↑ correctly classified  $\pi^+$  &  $\pi^-$

$\xi_i > 0$  & it is equal to the distance away from the correct hyperplane in the incorrect direction.

$$(w^*, b^*) = \underset{(w, b)}{\operatorname{argmax}} \frac{2}{\|w\|} = \underset{(w, b)}{\operatorname{argmin}} \frac{\|w\|}{2}$$

$w^*, b^* = \underset{(w,b)}{\text{argmin}} \left( \frac{\|w\|^2}{2} + C \cdot \frac{1}{n} \sum_{i=1}^n \xi_i \right)$

→ regularization  
 → hyperparameter  
 → average distance of misclassification pt for correct  $\pi_i$ 's

$$s.t. \begin{cases} y_i (w^T x_i + b) \geq 1 - \xi_i & \forall i \\ \text{for each } i \end{cases}$$

For correctly classified points  $\xi_i = 0$

$\boxed{\because \xi_i > 0}$

we have to minimize misclassification, means  $\sum_{i=1}^n \xi_i$  should be less.

→ As  $C \uparrow$ , tendency to make mistake on training set  
 $\Rightarrow$  overfit  $\Rightarrow$  high variance.

→ As  $C \downarrow$ ; tendency to underfit  $\Rightarrow$  high bias

→ This is called soft margin SVM.

28.3 why we take values +1 & -1 for support vector planes.

$\|w\| \neq 1$  (any vector)  
 need not be unit vector

$$\text{margin} = \frac{2}{\|w\|}$$

whole task is to maximize the margin.

$$\pi^+ : w^T x + b = K_1 ; \pi^- : w^T x + b = K_2 \quad \boxed{K > 0}$$

why we take only  $K_1$  &  $K_2$ , the  $\pi^+$  &  $\pi^-$  should be equally separable

$$\text{margin} = \frac{2K}{\|w\|} ; \underset{(w,b)}{\text{argmax}} \frac{2}{\|w\|} = \underset{w,b}{\text{argmax}} \frac{2K}{\|w\|} \xrightarrow{K=4} \frac{8}{\|w\|}$$

we took +1 & -1 for convenience.

Case 2  $\pi^\top \div w^\top x + b = k$

$$\left(\frac{w}{k}\right)^\top x + \left(\frac{b}{k}\right) = 1$$

$$w \perp \pi$$

1)  $w$  need not be  $\perp$

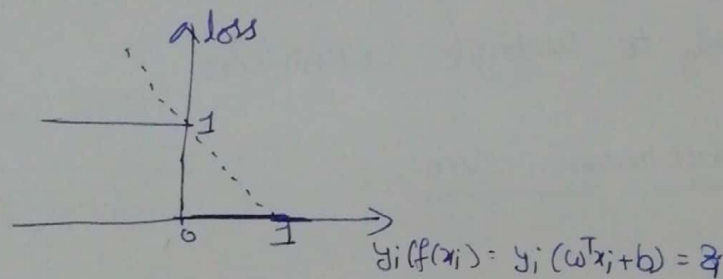
$$(w')^\top x + b' = 1$$

## 28.4 Loss function (Hinge Loss) based interpretation

logistic regression  $\rightarrow$  logistic loss + reg

lr. regression  $\rightarrow$  lr. loss + reg

SVM  $\div$  hinge loss + reg



when  $z_i > 0$  :  $x_i$  correctly classified

$z_i < 0 \div x_i$  is incorrectly classified

### Hinge Loss

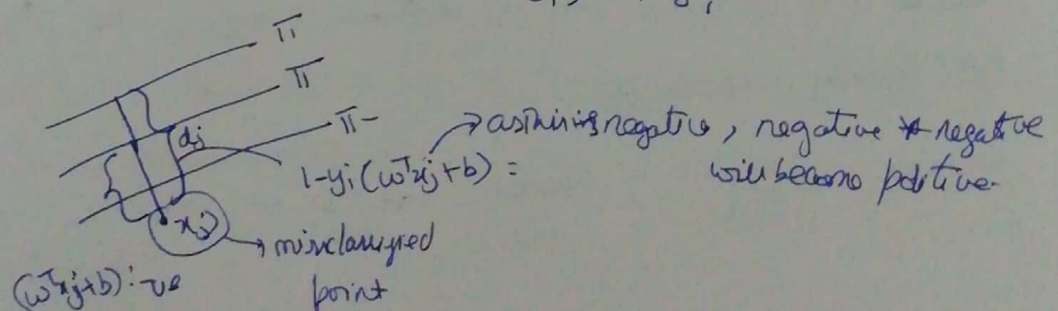
$z_i \geq 1$  ; hinge loss is 0

$z_i < 1$  ; hinge loss =  $1 - z_i$

$$\rightarrow \max(0, 1 - z_i)$$

Case 1  $\div z_i \geq 1$  ;  $1 - z_i$  is a negative value  $\Rightarrow \max(0, 1 - z_i) = 0$

Case 2  $\div z_i \leq 1$  ;  $1 - z_i > 0 \Rightarrow \max(0, 1 - z_i) = 1 - z_i$





$$d_j = 1 - y_j (\omega^T x_j + b) = z_j$$

$$\xi_j = \text{dist of } x_j \text{ to } \pi^+ = d_j = 1 - z_j$$

$$\xi_j = 1 - z_j \rightarrow \text{when } x_j \text{ is misclassified}$$

$$\max(0, 1 - z_i) = \xi_i$$

$$\text{Soft SVM} : \min_{\omega, b} \frac{\|\omega\|^2}{2} + \sum_{i=1}^n \xi_i$$

$$s.t. \begin{cases} (1 - y_i)(\omega^T x_i + b) \leq z_i \\ \xi_i \geq 0 \end{cases}$$

$\uparrow = \text{overfit}$

$\downarrow = \text{underfit}$

$$\text{loss min} : \min_{\omega, b} \sum_{i=1}^n \max(0, 1 - y_i(\omega^T x_i + b)) + \frac{\lambda}{2} \|\omega\|^2$$

$$\|\omega\| \geq 0 \Rightarrow \min \frac{\|\omega\|^2}{2} \text{ is same as } \min \|\omega\|$$

$\lambda \uparrow = \text{underfit}$

$\lambda \downarrow = \text{overfit}$

$\xi = 0$ ; for correctly classified points

$\xi > 0$ ; for incorrectly classified points.

## 28.5 Dual form of SVM formulation

$$\left\{ \begin{array}{l} \min_{w, b} \quad \frac{1}{2} \|w\|^2 + C \sum_{i=1}^n \xi_i \\ \text{s.t.} \quad y_i (w^T x_i + b) \geq 1 - \xi_i \quad \forall i \\ \xi_i \geq 0 \end{array} \right\} =$$

Primal of SVM

## dual form of SVM

$$\left\{ \begin{array}{l} \max_{\alpha_i} \quad \sum_{i=1}^n \alpha_i - \frac{1}{2} \sum_{i=1}^n \sum_{j=1}^n \alpha_i \alpha_j y_i y_j x_i^T x_j \\ \text{s.t.} \quad \alpha_i \geq 0 \end{array} \right\}$$

$\rightarrow \alpha_i = 0$  for non-SVs  
 $\alpha_i > 0$  for support vectors

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

## dual form

- ① for every  $x_i$  there is an  $\alpha_i$  corresponding
- ②  $x_i^T x_j \Rightarrow x_i$  is only occur in the form.

## End of the dual

$$x_q : (w^T x_q + b) = f(x_q)$$

- ③  $f(x_q) = \sum_{i=1}^n \alpha_i y_i x_i^T x_q + b$
- ④  $\alpha_i > 0$  only for support vectors  
 $\alpha_i = 0$  for non support vectors.

$f(x_q)$ : Only pts that matter are support vectors

$$x_i^T x_j = x_i \cdot x_j = \cos(\theta) = \text{sim}(x_i, x_j)$$

If  $\|x_i\| = 1 ; \|x_j\| = 1$

- $\rightarrow$  we can replace  $x_i^T x_j$  with any similarity function  $\text{sim}(x_i, x_j)$
- $\rightarrow$  This is what made SVM super popular
- $\rightarrow \text{sim}(x_i, x_j) = K(x_i, x_j)$



updated Eq

$$\max_{\alpha_i} \sum_{i=1}^n \alpha_i - \frac{1}{2} \sum_{i=1}^n \sum_{j=1}^n \alpha_i \alpha_j y_i y_j K(x_i, x_j)$$

$\uparrow$  Kernel-fn.

$$s.t. \sum_{i=1}^n \alpha_i y_i = 0 ; \alpha_i \geq 0$$

### 28.6 Kernel Trick

→ one class of similarity function is kernel function.

The most important idea in SVM is Kernel.

Soft svm hyperplane  $\approx$  log-reg

↳ margin-max.

If we ~~did~~ did not apply kernel trick and just leave it as  $x_i^T x_j$ .

Then it is called linear svm.

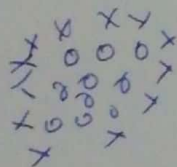
$$\text{lin. svm} : x_i^T x_j$$

$$\text{kernel svm} : K(x_i, x_j)$$

In a linear svm we are finding margin maximizing hyperplane in the space of  $x_i$ 's

log-reg we are ~~not~~ finding a hyperplane which minimizes logistic loss in the space of  $x_i$ 's

lin. svm & log-reg work in below case.



$$f_1, f_2 \rightarrow f_1^2, f_2^2$$

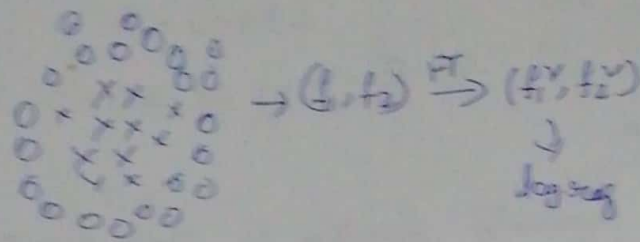
In the above case log-reg + feature transform is successful.

kernel svm will be successful.

↳ It will transform  $x_i \rightarrow x_i'$ , it tries to find hyperplane in the space of  $x_i'$

Kernel SVM  $\div$  non-lr separable data

25.7 Kernel Polynomial



Let's see how Kernelization solve this problem.

$$K(x_1, x_2) = (x_1^T x_2 + c)^d$$

$$(e.g) K(x_1, x_2) = \underbrace{(1 + x_1^T x_2)^2}_{\text{Quadratic Kernel}}$$

$$x_1 = \langle x_{11}, x_{12} \rangle$$

$$x_2 = \langle x_{21}, x_{22} \rangle$$

$$= (1 + x_{11}x_{21} + x_{12}x_{22})^2$$

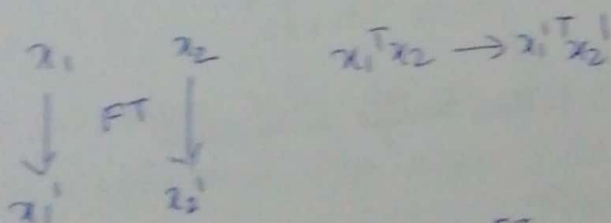
$$= 1 + x_{11}^2 x_{21}^2 + x_{12}^2 x_{22}^2 + 2x_{11}x_{21} + 2x_{12}x_{22} +$$

$$2x_{11}x_{21}x_{12}x_{22}$$

$$\Rightarrow [1, x_{11}^2, x_{12}^2, \sqrt{2}x_{11}, \sqrt{2}x_{12}, \sqrt{2}x_{11}x_{12}] = x_1'$$

$$[1, x_{21}^2, x_{22}^2, \sqrt{2}x_{21}, \sqrt{2}x_{22}, \sqrt{2}x_{21}x_{22}] = x_2'$$

$$= (x_1')^T (x_2')$$



Kernelization takes data  $(d)$   $\xrightarrow[\text{Internally Explicitly}]{FT}$   $d'$   $d' > d$

Mercer's Theorem :-

Shows what Kernel-trick is doing

$d \rightarrow d' \rightarrow$  Is-separable  
typically  $d' > d$   
not linearly separable.

SVM  $\rightarrow$  Explicit FT  $\rightarrow$  logit-reg

$\hookrightarrow$  finding the right kernel.

### 2.7 RBF - Kernel

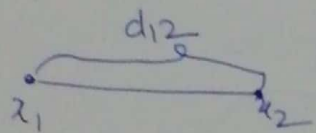
$\rightarrow$  Radial Basis Function (RBF)

$\rightarrow$  most popular / general purpose kernel : RBF

$(x_1, x_2)$

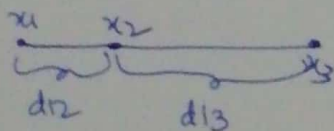
$$K_{RBF}(x_1, x_2) = \exp\left(\frac{-\|x_1 - x_2\|^2}{2\sigma^2}\right)$$

$\underbrace{2\sigma^2}_{\text{hyper parameter}}$   
 $\|x_1 - x_2\|^2 = d_{12}^2$



$$K_{RBF}(x_1, x_2) = \exp\left(\frac{-d_{12}^2}{2\sigma^2}\right) = \frac{1}{e^{d^2/2\sigma^2}}$$

$d_{12} \uparrow ; K(x_1, x_2) \downarrow$   
Similarity



$$K(x_1, x_2) > K(x_1, x_3)$$

② Impact of ' $\sigma$ '

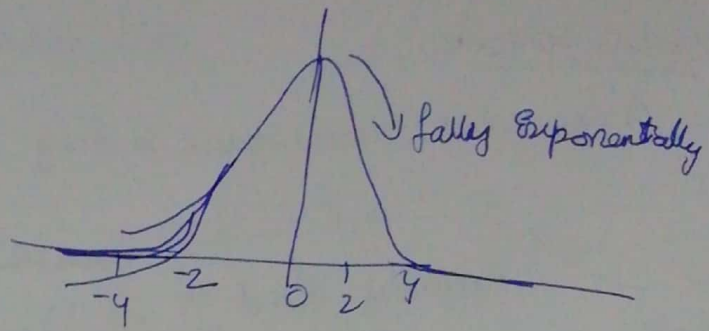
$$\sigma = 1, \sigma = 0.1, \sigma = 10$$



①  $d = 0$

$k = 1$

②  $d \uparrow ; k \downarrow$  Exponentially



kernel  $\rightarrow$  similar  
dist  $\rightarrow$  dissimilar.

RBF  $\approx$  gaussian PDF

Case

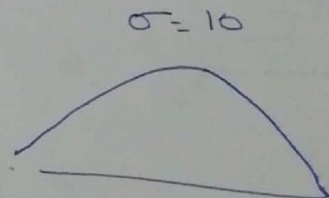
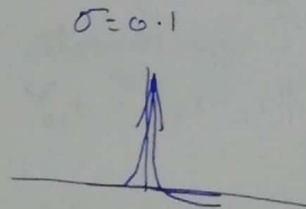
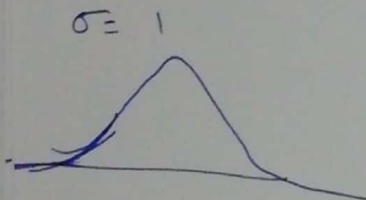
$\sigma = 0.1$

$\sigma^2 = 0.01$

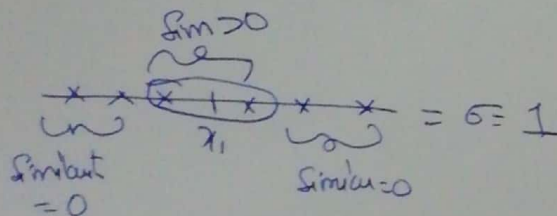
$d > 1 ; k = 0$

let  $\sigma$  increase

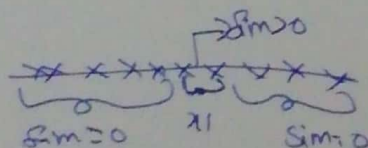
$\sigma = 10, \sigma^2 = 100$



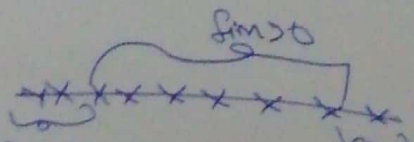
as  $\sigma \uparrow$



$\sigma = 0.1$



$\sigma = 10$



There is a relation b/w k-nn & RBF Kernel.

$\sigma \uparrow \Rightarrow k \uparrow$  in k-nn  
(in RBF)

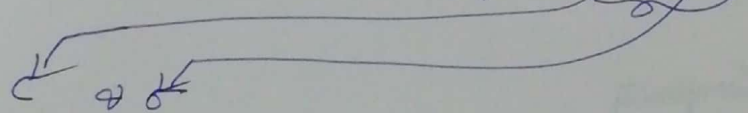
RBF-sum  $\approx$  k-nn

k-nn: store all the k pts. & LSH

RBF-sum: we can just have SV's & its  $\alpha_i$ 's

#SV's  $\ll n$   
(number of points)

If we don't have best Kernel simply use RBF-sum.

Soft margin: 

$(C, \sigma) \rightarrow$  grid search  
random search.

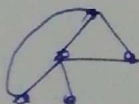
### 289 Domain Specific Kernels

RBF  $\rightarrow$  general purpose kernel

$\rightarrow$  String kernels (for text classification)

$\rightarrow$  genome kernels

$\rightarrow$  Graph kernels



$\rightarrow$  given the ~~data~~ problem, we have to select the appropriate kernels

For Amazon data set, String kernel will give better results than RBF

## 28-10 : Train & runtime Complexities

Train  $\rightarrow$  SGD  
 $\rightarrow$  Specialized algo (dual)  $\rightarrow$  Sequential minimal optimization (SMO)

libsvm: best libraries for training SVM's.

Training Time  $\approx O(n^2)$  for kernel SVM's

If  $n$  is large  $\rightarrow O(n^2)$  is very large.

Typically don't use SVM when  $n$  is large

applications

Internet

Runtime Complexity

$$f(x_q) = \sum_{i=1}^n \alpha_i y_i K(x_i, x_q) + b$$

$\uparrow$   
 $\alpha = 0$  for non SV's

$$\# \text{ SV's} = K \Rightarrow O(Kd)$$

$\uparrow$   
# of support vectors



20.11 nu-SVM: Control Errors & support vectors.

C-SVM  $\rightarrow$  Original formulation.  $(C \geq 0)$

alternative formulation of SVM:

$$\text{nu-SVM} \div 0 \leq \text{nu} \leq 1$$

nu  $\geq$  fraction of errors

nu: hyper parameter

SVM

Obtain

↓

I don't want 10% errors

↓

$$\text{nu} = 0.1$$

↓

I don't want 1% errors

↓

$$\text{nu} = 0.01$$

let assume

$$\text{nu} = 0.01 \Rightarrow \text{x. of errors} \leq 1\%$$

$$\underline{\underline{\# \text{ sv's} \geq 1\% \text{ of } n}}$$

$\rightarrow$  For runtime complexity : fewer sv's is ideal

## 28.13 Cases for SVM's

→ Feature Engineering & FT

↳ Finding the right kernel (SVM tends to work well)

→ Decision Surface

ls. SVM's: hyperplane

Kernel-SVM's:  $(d)$   $x_i \rightarrow$  non-linear surface

↓ ↓  
 $d'$   $x_i' \rightarrow$  linear surface.

$$\boxed{d' > d}$$

→ Imagine we are given distance or similarity fn

$$\rightarrow K(x_i, x_j)$$

→ Interpretability & Feature Importance

↳ no way for us to get feature importance directly

↳ we can use forward feature selection

→ Outlier: very little impact on the model

↳ SV's that matter

↳ RBF with a small  $\sigma \rightarrow$  knn with small ' $k$ '

↳ these two might get impacted with small  $\sigma$  or  $k$

→ Bias Variance

$C \uparrow \rightarrow$  overfit  $\uparrow$

$C \downarrow \rightarrow$  underfit no model

→ large  $d \rightarrow$  v. good for SVM ( $d \rightarrow d'$ )

↳ good kernels if we have otherwise RBF

→ Best cases :

↳ Sigmoid Kernel

→ Worst cases

↳ when 'n' is large → Training time is typically high  
Ex: internet based application.

↳ K is large =

↳ we cannot have low latency

when n is v.v large ppl end up using Logistic Regression.



## 28.14 Code Sample

→ on scikit learn

→ sklearn implements SVC, nu-SVC, SVR

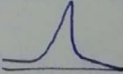
```
class sklearn.svm.SVC (C=1.0, kernel='rbf', degree=3, gamma='auto',  
coef0=0.0, shrinking=True, probability=False, tol=0.001,  
cache_size=200, class_weight=None, verbose=False, max_iter=1,  
decision_function_shape='ovr', random_state=None)
```

→ C: bias & Variance

degree is useful only when the kernel is polynomial

→ gamma: In RBF we have  $\sigma$  (sigma). In sklearn it is referred

to as gamma  $\gamma = \frac{1}{\sigma}$

$\sigma \downarrow \quad \gamma \uparrow$  

tol: 0.001

→ Tolerance says when we are moving from  $i^{\text{th}}$  iteration to  $i+1^{\text{th}}$  iteration  
( $w_i - w_{i+1}$ )

If difference b/w  $w_i$  &  $w_{i+1}$  is smaller than tolerance, terminate the loop.

→ class\_weight = None

Imbalance data if we want to do upsampling then we have to give class weight

→ max\_iter = -1

it will iterate until the tolerance is reached

→ decision\_function\_shape

If it is multiclass then we have to provide 'ovr' (one vs rest)

Import numpy as np

$X = \text{np.array}([[-1, -1], [-2, -1], [1, 1], [2, 1]])$

$y = \text{np.array}([1, 1, 2, 2])$

from sklearn.svm import SVC

$\text{clf} = \text{SVC}()$

$\text{clf.fit}(X, y)$

$\rightarrow \text{print}(\text{clf.predict}([[-0.8, -1]]))$

o/p  $\rightarrow [1]$

\*

gamma

The gamma parameters can be seen as the inverse of the radius of influence of samples selected by the model as support vectors.

$\text{gamma} = 10^{-1}, c = 10^{-2} \rightarrow g = 10^0, c = 10^{-2} \rightarrow g = 10^{-1}, c = 10^{-2}$

When  $\sigma \downarrow$ , we are overfitting  $\xrightarrow{\text{when moving}}$   
 $\sigma \uparrow$

When  $\sigma = 1$ ,  $\sigma \downarrow$  in RBF both are same, they will overfit

Similarly

When  $c \uparrow$  it will overfit

$\text{gamma} = 10^{-1}, c = 10^{-2}$

$g = 10^{-1}, c = 10^0$

$g = 10^{-1}, c = 10^{-2}$

## Sklearn. svm. NuSVC

only thing will change is  $\nu$  ( $\nu=0.5$ )

$\nu$ :

An upper bound on the fraction of training errors and a lower bound of the fraction of support vectors.

Should be in the interval of  $(0, 1)$

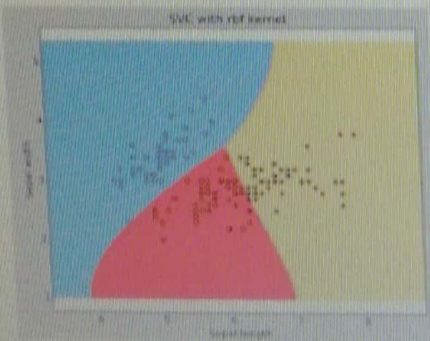


**gamma** Kernel coefficient for 'rbf', 'poly' and 'sigmoid'. Higher the value of gamma, will try to exact fit the as per training data set i.e. generalization error and cause over-fitting problem.

**Example:** Let's difference if we have gamma different gamma values like 0, 10 or 100.

```
svc = svm.SVC(kernel='rbf', C=1, gamma=0).fit(X, y)
```

**gamma = 0**



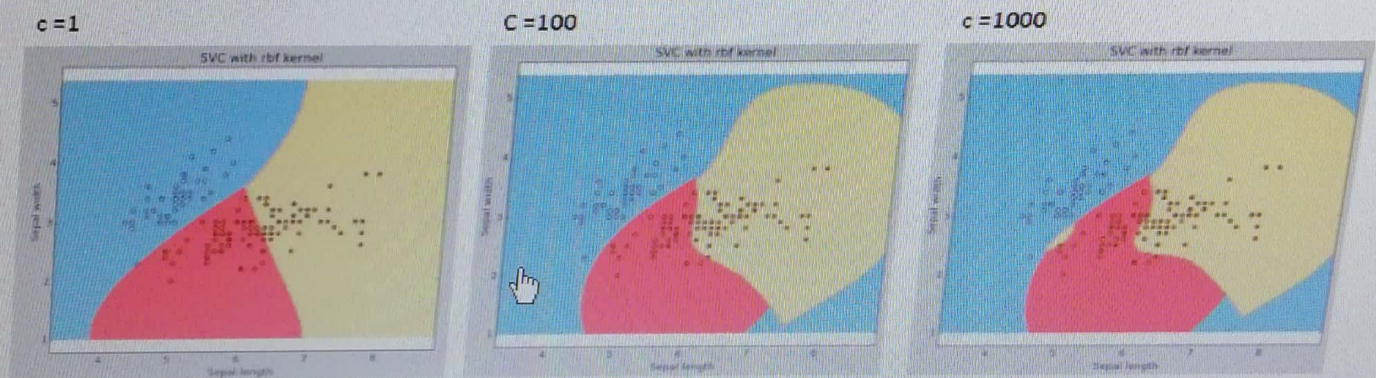
**gamma = 10**



**gamma = 100**



**C:** Penalty parameter  $C$  of the error term. It also controls the trade off between smooth decision boundary and classifying the training points correctly.



We should always look at the cross validation score to have effective combination of these parameters and avoid over-fitting.

## Pros and Cons associated with SVM

- **Pros:**

- It works really well with clear margin of separation
- It is effective in high dimensional spaces.
- It is effective in cases where number of dimensions is greater than the number of samples.
- It uses a subset of training points in the decision function (called support vectors), so it is also memory efficient.

- **Cons:**

- It doesn't perform well when we have large data set because the required training time is higher
- It also doesn't perform very well when the data set has more noise i.e. target classes are overlapping
- SVM doesn't directly provide probability estimates, these are calculated using an expensive five-fold cross-validation. It is related SVC method of Python scikit-learn library.



# Medium

Sign in

```
from sklearn import svm, grid_search

def svc_param_selection(X, y, nfolds):
    Cs = [0.001, 0.01, 0.1, 1, 10]
    gammas = [0.001, 0.01, 0.1, 1]
    param_grid = {'C': Cs, 'gamma' : gammas}
    grid_search = GridSearchCV(svm.SVC(kernel='rbf'), param_grid,
cv=nfolds)
    grid_search.fit(X, y)
    grid_search.best_params_
    return grid_search.best_params_
```