

# SUPPORT VECTOR MACHINES

[GEOMETRIC INTUITION](#)

[MATHEMATICAL DERIVATION](#)

[WHY WE TAKE +1 AND -1 FOR SUPPORT VECTOR PLANES](#)

[LOSS FUNCTION\(HINGE LOSS\) INTERPRETATION](#)

[DUAL FORM OF SVM](#)

[KERNEL TRICK FOR SVM](#)

[POLYNOMIAL KERNEL](#)

[RADIAL BASIS FUNCTION KERNEL \(RBF\)](#)

[TRAIN AND RUN-TIME COMPLEXITIES](#)

[SVM REGRESSION](#)

## GEOMETRIC INTUITION

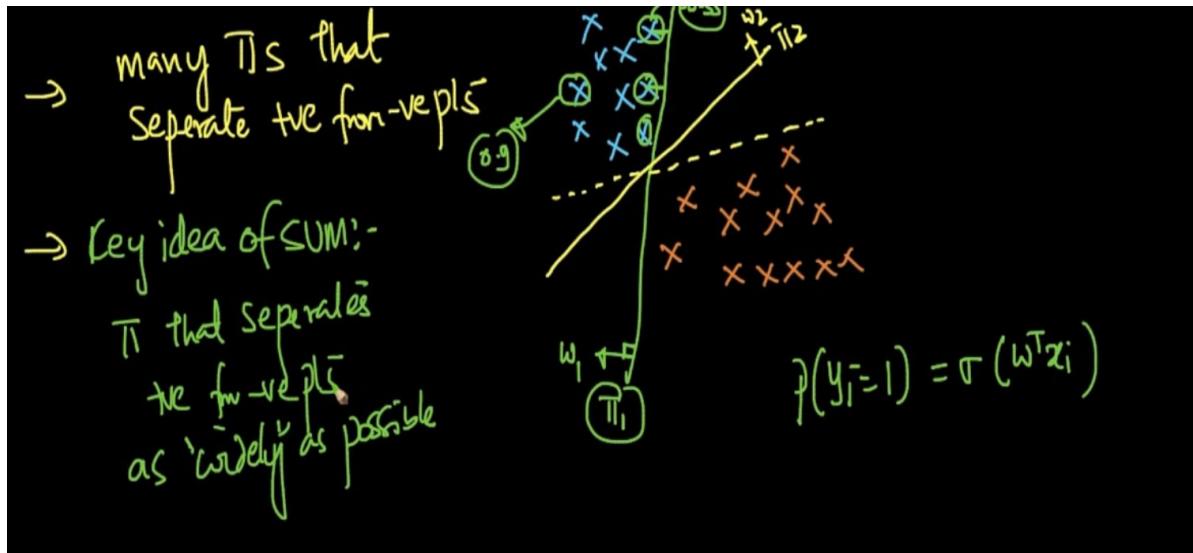
Support Vector Machines (SVM)

✓ (SVM) → popular ML  $\xrightarrow{\text{classifn}}$   $\xrightarrow{\text{regrsn}}$

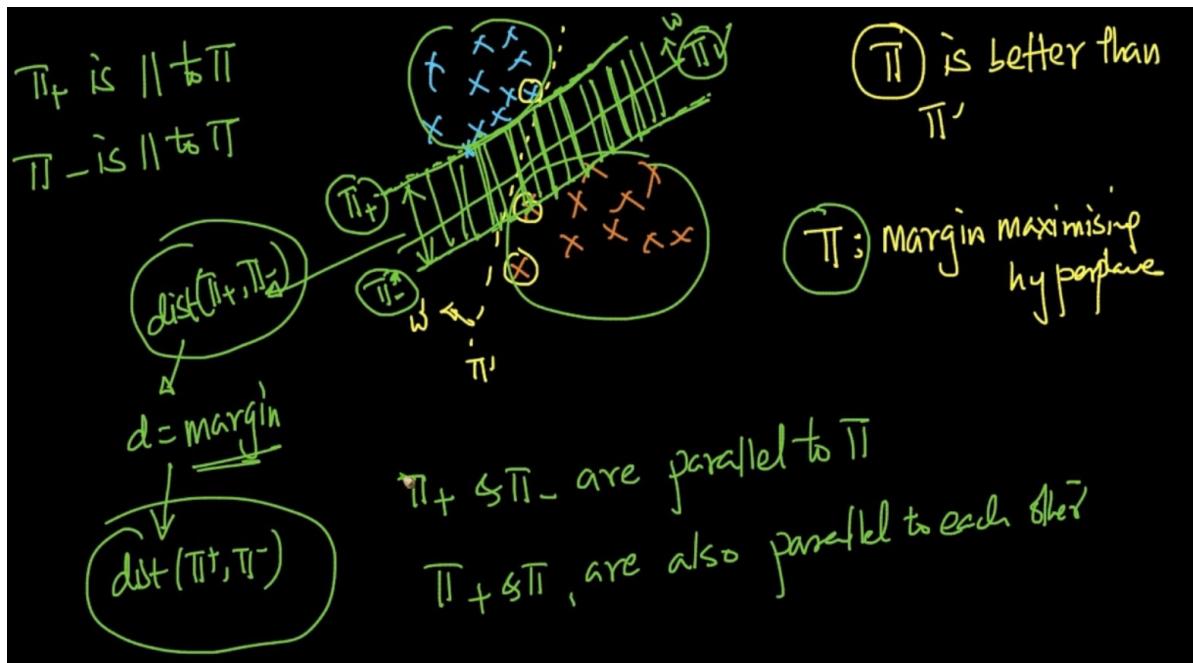
↳ 1990's

geom-intuition - Q

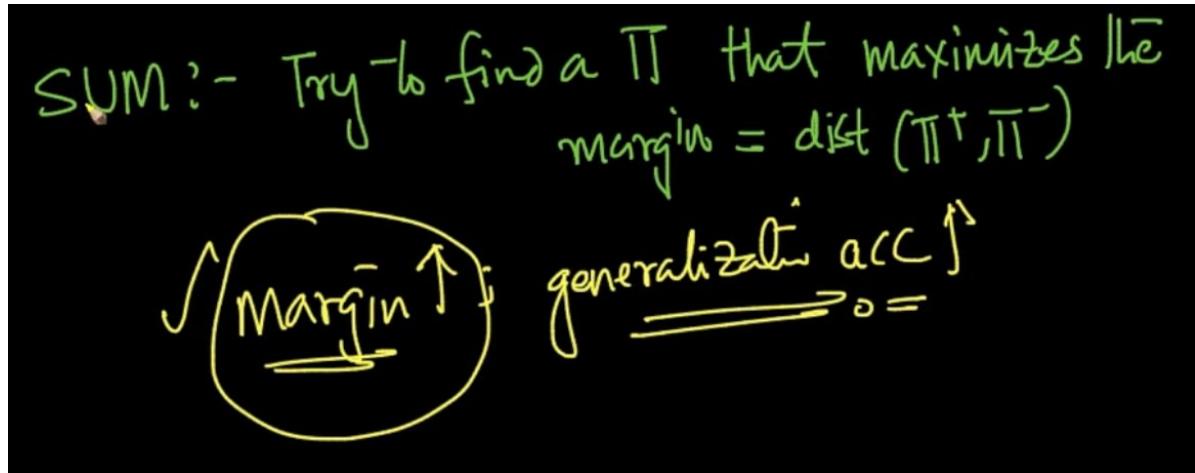
Support Vector Machines is a popular ML algorithm used for both Classification/ Regression



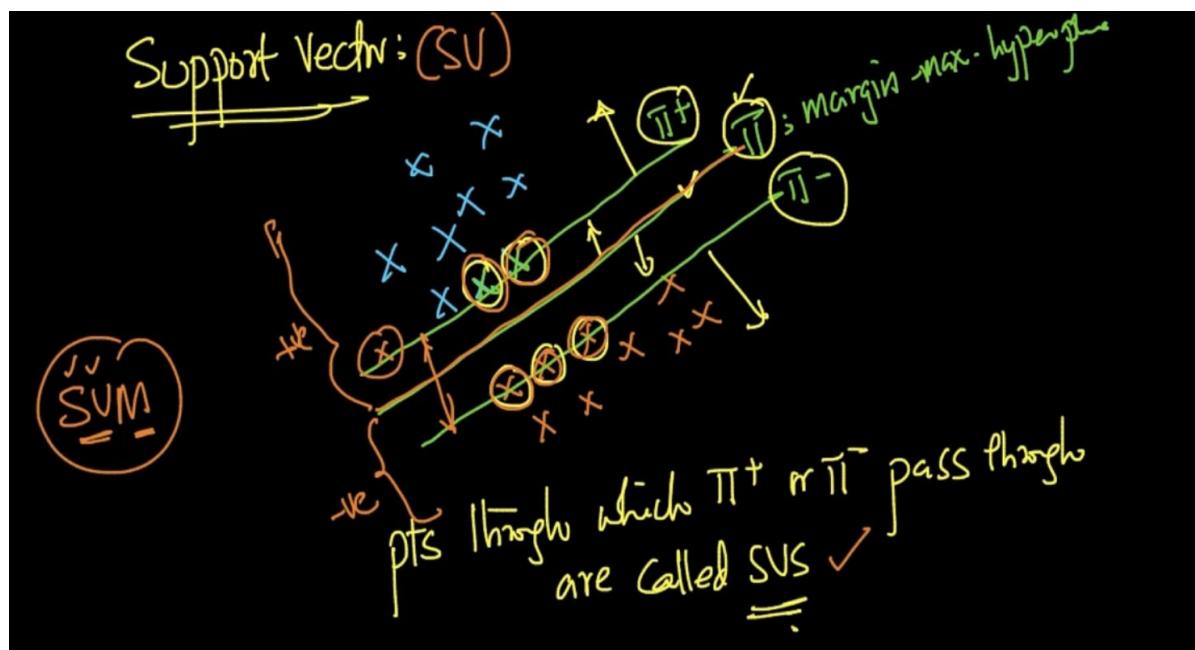
Dataset can be separated using many hyper-planes but if we take  $\pi_1$  the closest points will've lesser probability if we take something like sigmoid function which takes distance into account. SO our SVM wants a  $\pi$  that separates +ve and -ve as widely as possible



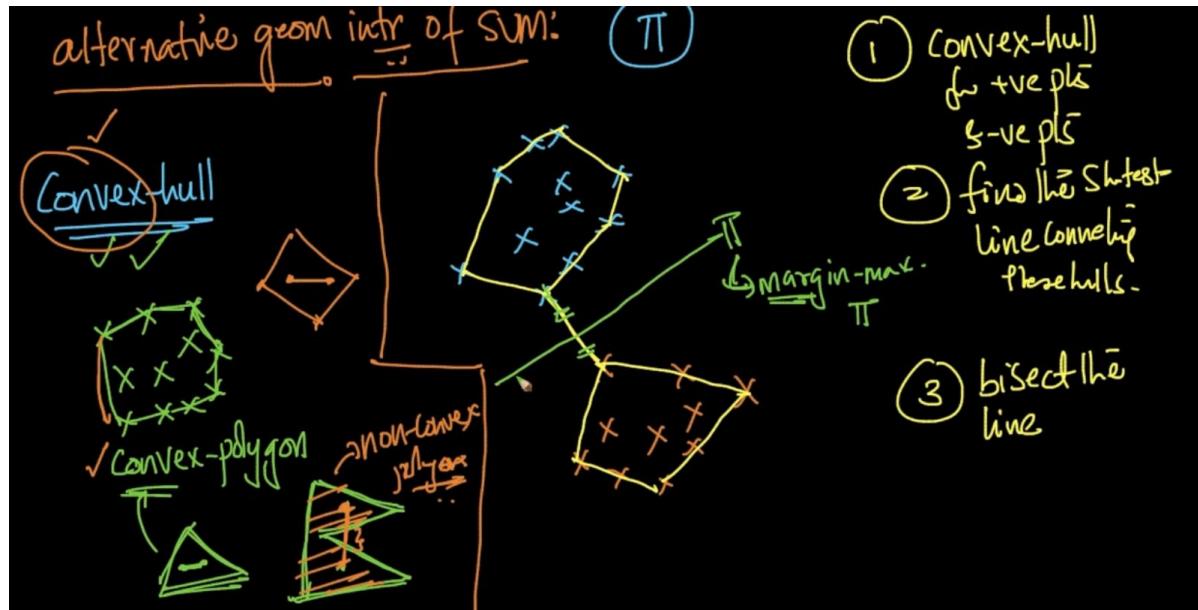
Our green line  $\pi$  is the margin maximizes hyperplane. If we take parallel lines which touch the *first points* in the +ve and -ve points. We calculate the distance which is known as margin between that two planes ( $\pi_+, \pi_-$ )



Svm tries to find the hyperplane  $\pi$  that maximizes the margin . If my margin is high then our chance of misclassifying decreases as margin increases our generalization accuracy i.e accuracy on unseen data increases



Support Vectors are the points through which  $\pi_+$  or  $\pi_-$  passes through. If a new point comes it's classified on the basis of  $\pi$ : Margin max hyperplane and that's it



A convex hull is a shape inside which all the data points are there or the points are on the convex hull .Follow all the steps to perform given above to perform SVM

## MATHEMATICAL DERIVATION

Mathematical formulation of SVM:

$\Pi^-$ : margin maximization

let  $\Pi^-: \underline{w}^T x + b = 0$

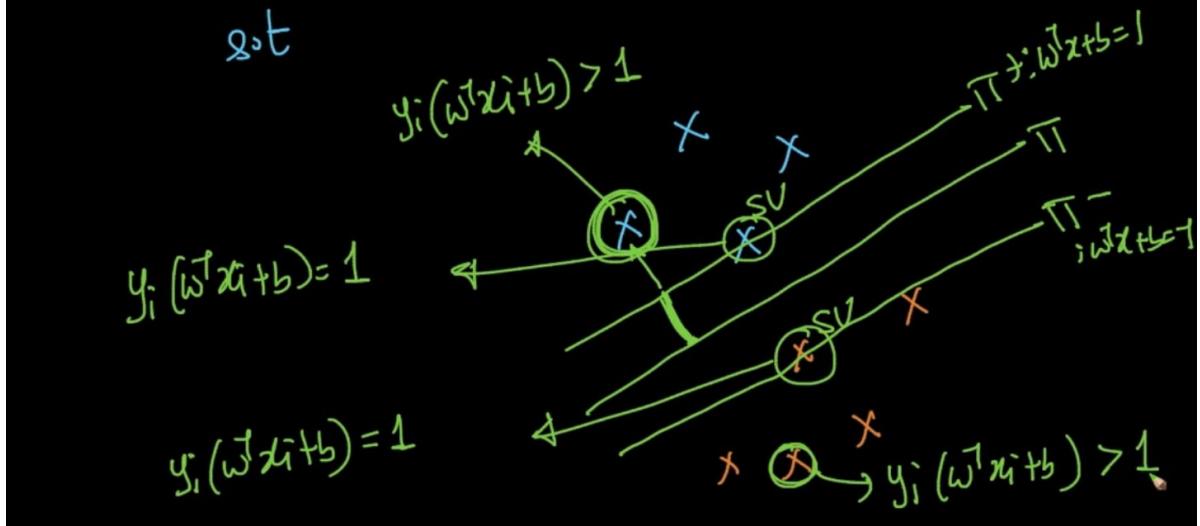
if  $\Pi^+: w^T x + b = 1$

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

Margin:  $d = \frac{2}{\|w\|}$  ✓

$\pi: w^T x + b = 0$  and our  $\pi^+: w^T x + b = 1$  and  $\pi^-: w^T x + b = -1$  so our margin i.e distance  $d = \frac{2}{\|w\|}$ . So we need to find  $(w^*, b^*)$  such that it maximizes  $\frac{2}{\|w\|}$

$$(\omega^T, b) = \underset{w, b}{\operatorname{argmax}} \frac{2}{\|w\|} = \text{Margin}$$



For support vectors :  $y_i(w^T x + b) = 1$

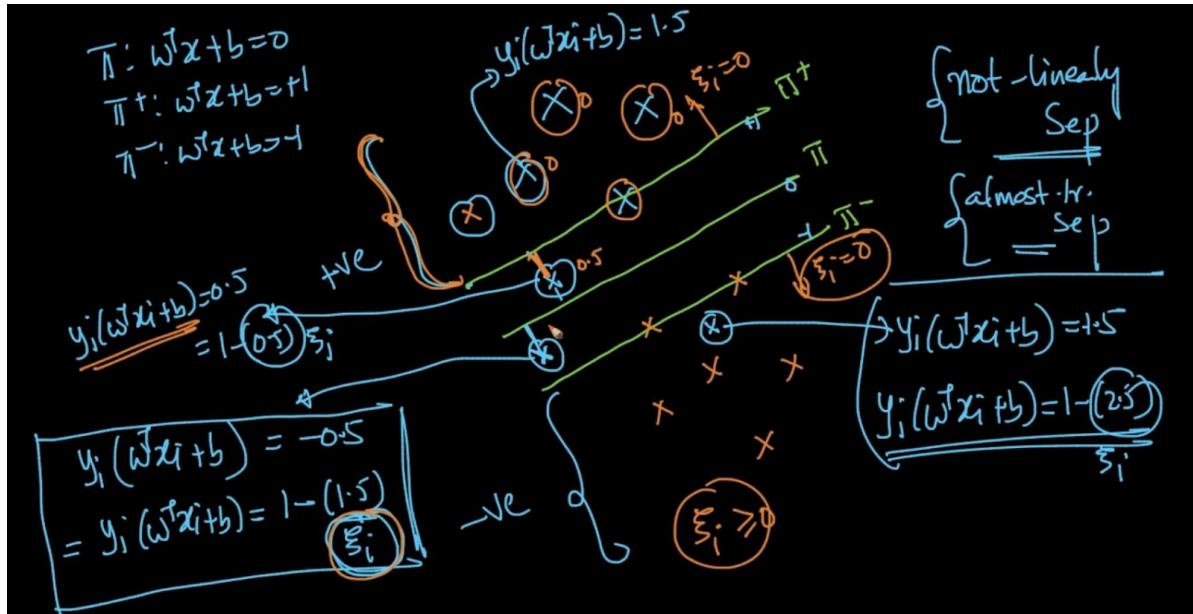
For data points greater than hyper-planes:  $y_i(w^T x + b) > 1$  as y will be -1 or +1 and  $(w^T x + b) > 1$  or  $(w^T x + b) < -1$

CONSTN,  
optimzn. of  
prob.  
of SVM

$$\left\{ \begin{array}{l} \omega^T, b^* = \underset{w, b}{\operatorname{argmax}} \frac{2}{\|w\|} \\ \text{s.t } \forall i, y_i(w^T x_i + b) \geq 1 \end{array} \right.$$

So we want to maximize the function with a constraint that for all i,  $y_i(w^T x + b) \geq 1$

To satisfy this constraint our data should be linearly separable



There can be cases like the above where a negative data is there in the positive part or vice versa. Here we'll never find the optimal  $w, b$  since it's not linearly separable. What do we do in these situations?

In the above, we can see that there's one +ve(blue) point amongst the negative data points. It's  $y_i(w^T x + b) = -1.5$  as  $y_i = 1$  and  $(w^T x + b) = -1.5$ . It can also be written as

$y_i(w^T x + b) = 1 - (2.5)$ . Same logic can be applied for the points in between the hyper-planes as well. So that 2.5 can be our  $\zeta$ . We use  $\zeta$  for all the points that are not in the region of their correct class. It works like a loss  $\zeta$ . In the region above the positive plane  $\pi^+$ ,  $\zeta = 0$ . Same for the correctly classified -ve points in the negative hyperplane

So, as  $\zeta \uparrow$ , the point is farther away from the correct hyper-plane. For ex: for the point which has  $\zeta = 2.5$ . 2.5 is the distance from the correct plane

$$x_i \rightarrow \xi_i$$

✓ if  $y_i(\mathbf{w}^T x_i + b) \geq 1$   
 ↗ correctly classified  
 $\pi^+ & \pi^-$   
 $\xi_i > 0$  & it is equal to the same units of  
 dist away from the correct hyperplane  
 in the incorrect dir

For every point  $x_i$  we've  $\xi_i$  and  $\xi_i = 0$  if the point is correctly classified and  $\xi_i > 0$  and it's equal to distance from the correct hyperplane and in the incorrect direction

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

margin  
 $C \cdot \frac{1}{n} \sum \xi_i$   
 avg. dist of  
 miss. pts for correct  
 H's

s.t.  $\left\{ \begin{array}{l} y_i(\mathbf{w}^T x_i + b) \geq 1 - \xi_i \\ \xi_i \geq 0 \end{array} \right.$ 
 }  
 corr. classif. pts  $\xi_i = 0$   
 incorr. classif. pts  $\xi_i > 0$

minimize errors = minimize misclassifications  
 $\min \|\sum \xi_i\|$

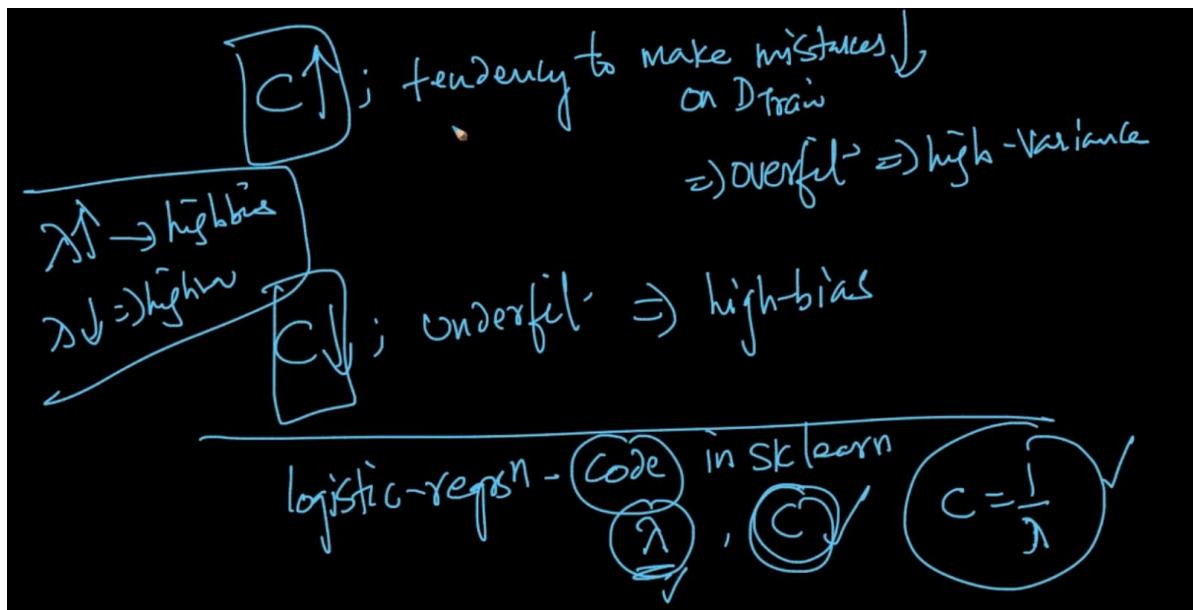
Our margin was  $\|\mathbf{w}\|/2$  and we want to maximize it/ So  $2/\|\mathbf{w}\|$  is the inverse which we want to minimize.  $\sum_{i=1}^n \xi_i$  : is what we want to minimize as well

$$(\hat{w}, \hat{b}) = \arg \min_{w, b} \left( \frac{\|w\|}{2} + C \cdot \sum_{i=1}^n \xi_i \right) \rightarrow \text{avg. dist for miss pl.}$$

min<sub>w,b</sub> C · loss + reg  
 s.t.  $y_i(w^T x_i + b) \geq 1 - \xi_i$   $\xi_i \geq 0$   
 margin  $\geq \xi_i$   $\xi_i \geq 0$   
 loss  $\leq \xi_i$   $\xi_i \geq 0$   
 hyp.  $y_i(w^T x_i + b) \geq 1 - \xi_i$   $\xi_i \geq 0$

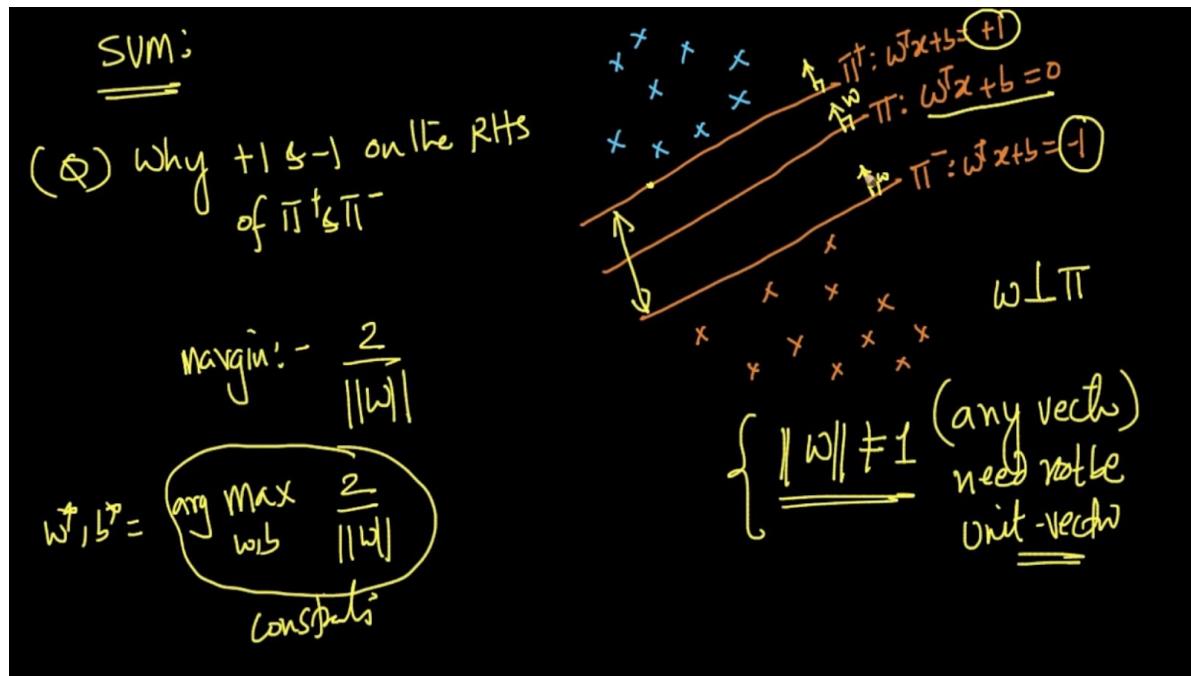
$$\boxed{\min_w (\text{logistic-loss}) + \lambda (\text{reg})}$$

$\|w\|/2$  is our regularization and  $\sum_{i=1}^n \xi_i$  is our Hinge-loss and  $C$  is hyperparameter



$C = \frac{1}{n}$  so it works exactly opposite to  $\lambda$  as seen above

## WHY WE TAKE +1 AND -1 FOR SUPPORT VECTOR PLANES



See, we said that  $\|w\| \neq 1$  it just should be  $w \perp \pi$ . Remember our objective was to maximize the margin  $(2/\|w\|)$  with some constraints which is basically that all the +ve points should be at one side and -ve points at other.

$$\textcircled{1} \quad \pi^+ : w^T x + b = K \quad K > 0 \quad K: \text{constant}$$

$$\pi^- : w^T x + b = -K$$

Margin:  $\frac{2K}{\|w\|}$

$$\arg \max_{w,b} \frac{2}{\|w\|} = \arg \max_{w,b} \frac{2K}{\|w\|} = \frac{8}{\|w\|}$$

$k$  can be anything there not just 1 but in the end it doesn't matter. Why can't we take  $k_1, k_2$  separately because our whole idea of SVM was to separate the +ve and -ve points well. It would be possible only if the hyperplanes  $\pi^+$  and  $\pi^-$  are equidistant from the optimal hyperplane  $\pi$  so it should be just  $k$ .  $\text{Max } 2/\|w\| = \text{max } (8/\|w\|)$  so it's the same.

$$\textcircled{2} \quad \text{Interpretation: } w^T x + b = k$$

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

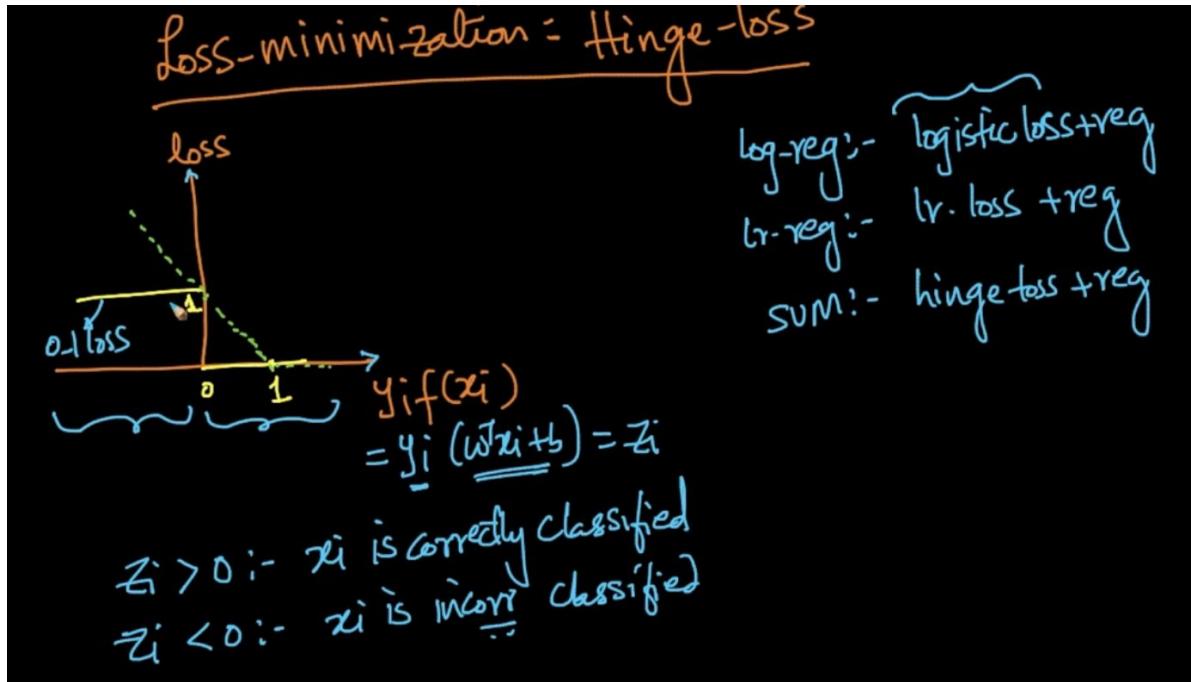
$$w \perp \pi \quad \|w\| \text{ need not be 1}$$

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

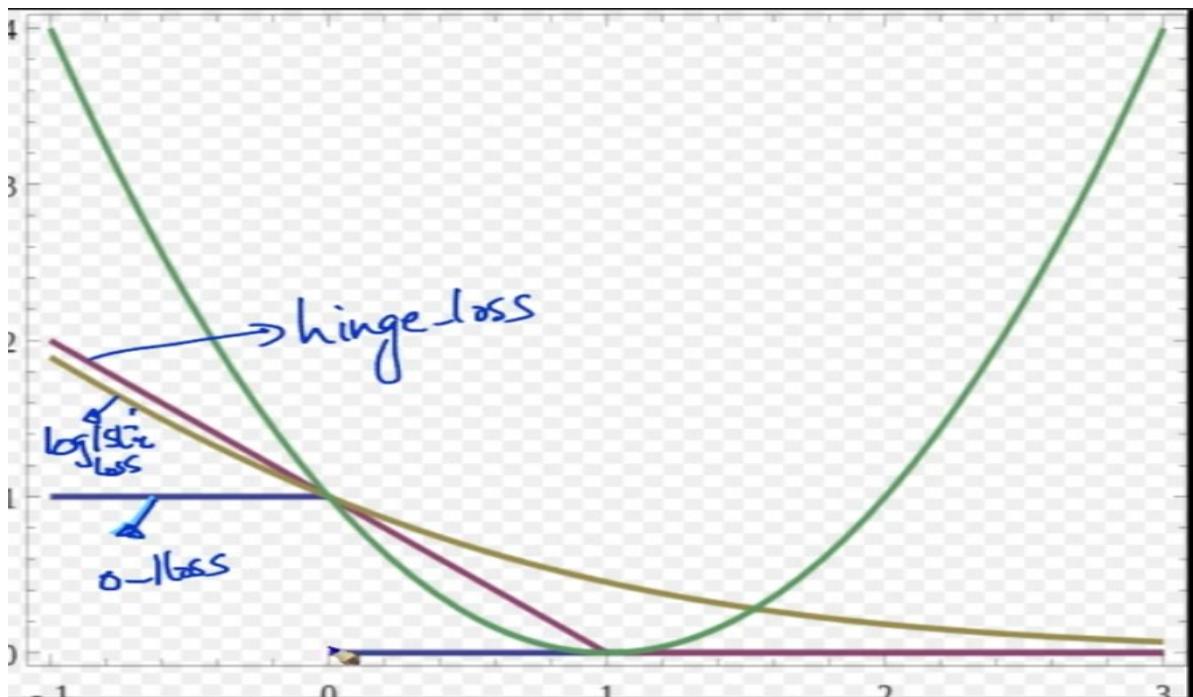
2nd interpretation is if we take  $w^T x + b = k$  as  $(\frac{w}{k})^T x + (\frac{b}{k}) = 1$  then also nothing would change since  $\|w\|$  need not be 1 it can be anything. So,  $(\frac{w}{k})^T x + (\frac{b}{k}) \Rightarrow (w')^T x + b' = 1$

So the reason we are taking +1 and -1 is for simplifying the maths.

## LOSS FUNCTION(HINGE LOSS) INTERPRETATION



$z_i = y_i(w^T x_i + b)$ . With 0-1 loss, when  $z_i > 0$ :  $x_i$  is correctly classified since  $y$  is +ve and  $w$  and  $x$  are at same direction whereas when  $z_i < 0$ :  $x_i$  is incorrectly classified since  $w$  and  $x$  will be at opposite direction.



hinge-loss :-  $\begin{cases} z_i \geq 1; \text{ hinge-loss} = 0 \\ z_i < 1; \text{ hinge-loss} = 1 - z_i \end{cases}$  ✓✓

$\rightarrow \underline{\max(0, 1 - z_i)}$

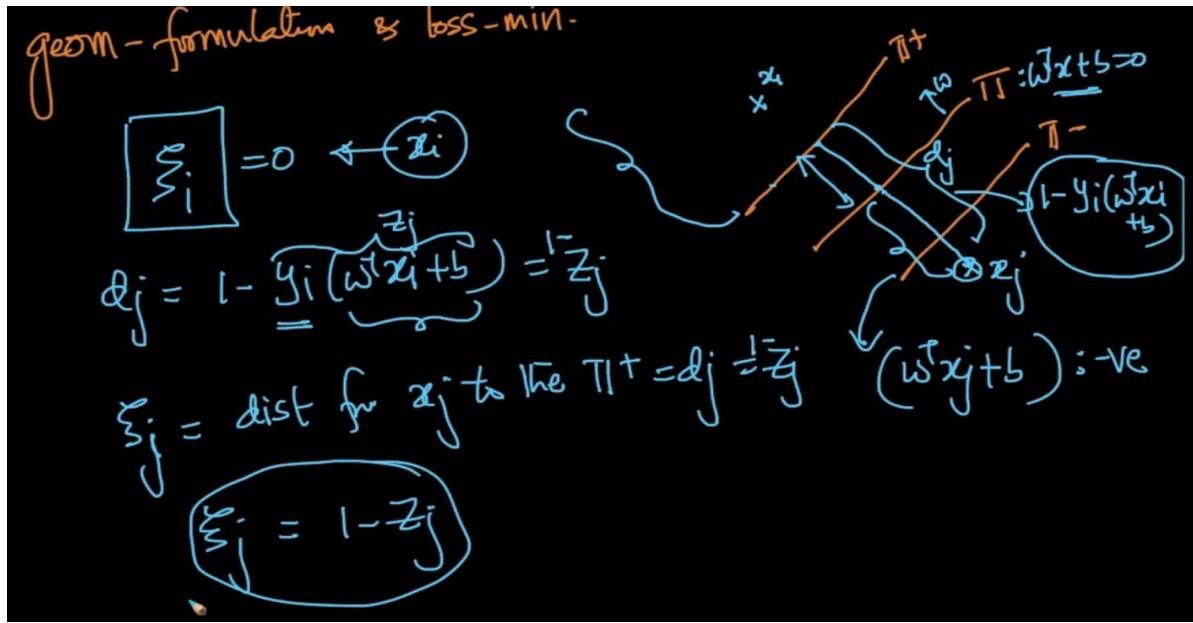
$\left\{ \begin{array}{l} \text{Case 1: } z_i \geq 1 \Rightarrow 1 - z_i \text{ is -ve value} \Rightarrow \max(0, 1 - z_i) = 0 \\ \text{Case 2: } z_i < 1; \quad 1 - z_i > 0 \Rightarrow \max(0, 1 - z_i) = 1 - z_i \end{array} \right.$

In Hinge Loss if  $z_i \geq 1$  then it's 0 or if  $z_i < 1$  it gets then loss increases

Hinge Loss :  $\max(0, 1 - z_i)$

Case 1:  $z_i \geq 1 \Rightarrow 1 - z_i$  is -ve value  $\Rightarrow \max(0, 1 - z_i) = 0$

Case 2:  $z_i < 1 \Rightarrow 1 - z_i > 0 \Rightarrow \max(0, 1 - z_i) = 1 - z_i$



For  $x_i$ ,  $\zeta = 0$  as it's correctly classified

For  $x_j$ , dist from  $\pi: w^T x + b \Rightarrow y_j \cdot w^T x_j + b$  but it's -ve since w and x are opposites

And from  $\pi^+$ , distance of  $x_j (d_j) = 1 - y_j \cdot (w^T x_j + b)$  since  $y = +ve$  and  $w^T x = -ve$  and 1 is distance from hyperplane  $\pi$  for  $\pi^+$  and  $y_j \cdot w^T x_j + b$  for  $x_j$ . So,  $d_j = 1 - z_j = \zeta_j$

soft sum:

$$\min_{w, b} \left( \frac{\|w\|}{2} \right) + C \sum_{i=1}^n \xi_i \quad \text{loss}$$

s.t.  $(1 - y_i(w^T x_i + b)) \geq \xi_i \quad \forall i$

$\xi_i > 0$

$\min_{w, b} \left[ \frac{\|w\|^2}{2} + \lambda \|w\|^2 \right]$

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

We are saying that Soft-SVM is equivalent to loss minimization formula it's just that 'C' is multiplied to loss and  $\lambda$  is multiplied to regularizer

## DUAL FORM OF SVM

Dual form of SVM

Primal of SVM

$$\begin{aligned} \text{soft-margin SVM} &= \min_{w, b} \frac{1}{2} \|w\|^2 + C \sum_{i=1}^n \xi_i \\ \text{s.t. } y_i(w^T x_i + b) &\geq 1 - \xi_i, \quad \xi_i \geq 0 \end{aligned}$$

Dual

$$\max_{\alpha} \sum_{i=1}^n \alpha_i - \frac{1}{2} \sum_{i=1}^n \sum_{j=1}^n \alpha_i \alpha_j y_i y_j x_i^T x_j$$

Equivalent

$$\begin{aligned} \text{s.t. } \alpha_i > 0 &\quad (4) \alpha_i > 0 \text{ only for SVs} \\ \sum_{i=1}^n \alpha_i y_i = 0 &\quad \alpha_i = 0 \text{ for non-SVs} \end{aligned}$$

$\boxed{\alpha_i}$ 
 $\boxed{x_i^T x_j}$ 
 $\boxed{y_i}$ 
 $\boxed{f(x_q)}$ 
 $\boxed{\sum_{i=1}^n \alpha_i y_i x_i^T x_q + b}$

Our soft margin SVM can also be written as like mentioned in dual form . The rules are also written there.

PS - In dual form we are not dealing with  $x_i$  separately. It's  $x_i^T x_j$

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

$\boxed{\alpha_i}$   $\boxed{y_i}$   $\boxed{x_i^T x_q}$   $\boxed{b}$

SVs :-  $\alpha_i > 0$   
 non SVs :-  $\alpha_i = 0$   
 $f(x_q)$  :- only pts that matter are SVs

SVM

As the 4th pt say  $\alpha_i > 0$  only for support vectors for other points the whole product will be 0  
 So, the only points that matter are Support Vectors

$$\max_{\alpha_i} \sum_{i=1}^n \alpha_i - \frac{1}{2} \sum_{i=1}^n \sum_{j=1}^n \alpha_i \alpha_j y_i y_j (\cancel{x_i^T x_j}) \rightarrow \underline{\text{Sim}(x_i, x_j)}$$

s.t.  $\alpha_i > 0 \rightarrow \begin{cases} \alpha_i = 0 \text{ for SVS} \\ \alpha_i > 0 \text{ for non-SVS} \end{cases} \checkmark$

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

$$\left\{ x_i^T x_j = x_i \cdot x_j = \underbrace{\text{Cosine sim}(x_i, x_j)}_{\text{if } \|x_i\| = 1, \|x_j\| = 1}$$

So you can replace  $x_i^T x_j$  with Cosine -similarity ( $x_i, x_j$ ). We can't do that with Soft -SVM

This part made SVM super-popular. We can represent Cosine -similarity ( $x_i, x_j$ )  $\Rightarrow K(x_i, x_j)$  which is also known as Kernel Function. We don't deal with  $x_i, x_j$  explicitly but with their dot product.

Run-time

$$f(\alpha_i) = \sum_{i=1}^n \alpha_i y_i \cancel{\frac{x_i^T x_j}{\|x_i\| \|x_j\|}} + b \rightarrow K(x_i, x_j)$$

$y_i$

This is same in Run-Time for SVM

## KERNEL TRICK FOR SVM

Kernel Trick:

$$\left\{ \begin{array}{l} \max_{\alpha_i} \sum_{i=1}^n \alpha_i - \frac{1}{2} \sum_{i=1}^n \sum_{j=1}^n \alpha_i \alpha_j y_i y_j x_i^T x_j \\ \text{s.t. } \sum_{i=1}^n \alpha_i y_i = 0, \alpha_i \geq 0 \end{array} \right. \rightarrow \begin{array}{l} \text{Sim}(x_i, x_j) \\ K(x_i, x_j) \\ \hookrightarrow \text{kernel fn} \end{array}$$

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

We'd seen that  $x_i^T x_j$  can be replaced with  $K(x_i, x_j)$ . It is called Kernel Trick

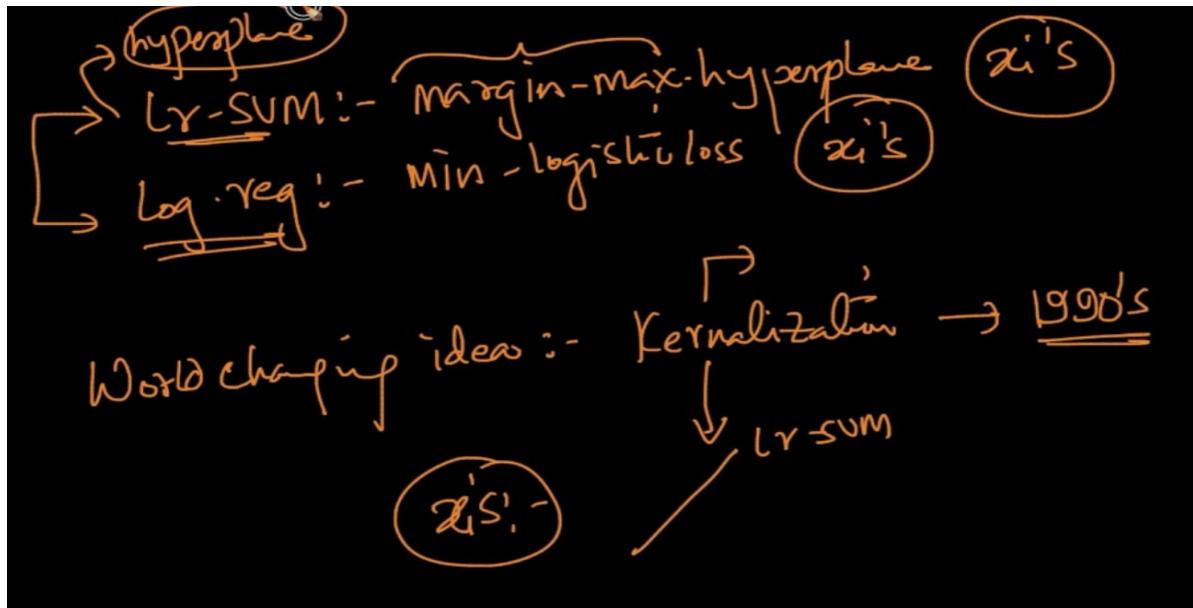
→ The most imp. idea in SVM is Kernel Trick

soft-SVM-hyperplanes  $\approx$  log-reg  
 $\hookrightarrow$  margin-max.

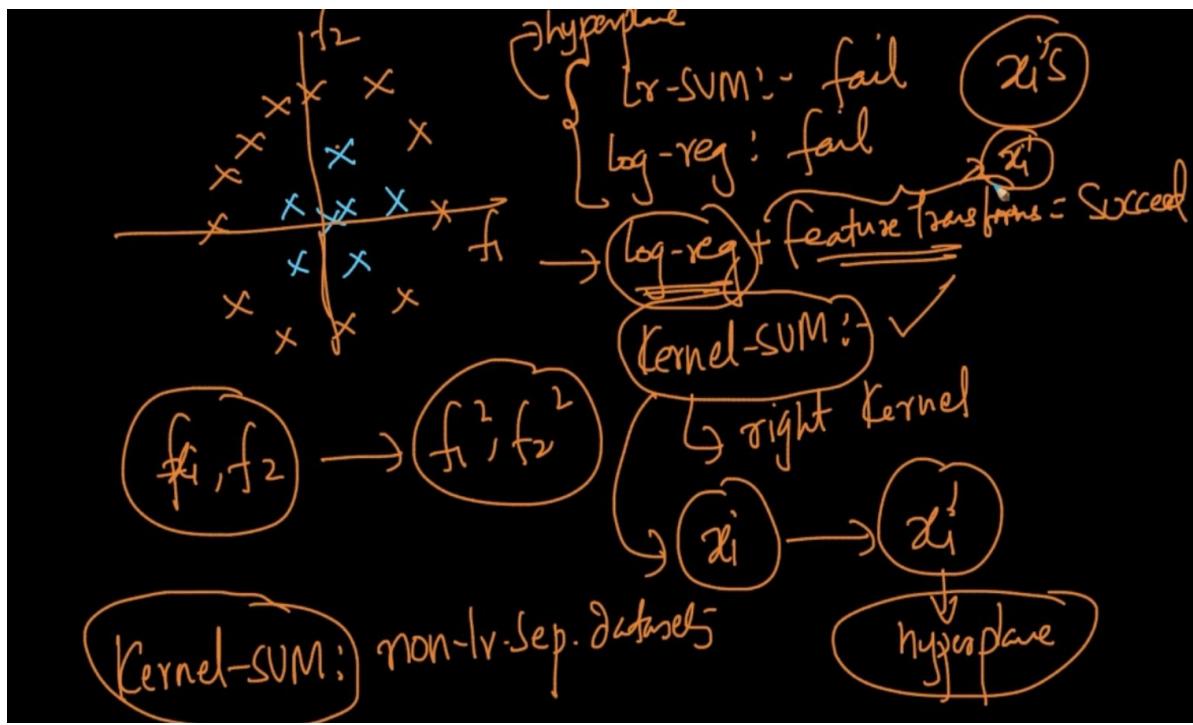
L.R.-SVM:- $x_i^T x_j$ Kernel-SVM:- <u><math>K(x_i, x_j)</math></u>	$\checkmark K(x_i, x_j) = x_i^T x_j$
--	--------------------------------------

Kernel-Trick is one of the most important idea in SVM.

If  $K(x_i, x_j) = x_i^T x_j$  then it is simply linear SVM but if we are using a function  $K(x_i, x_j)$  then Kernel-SVM



In Linear SVM we try to find the hyperplane which maximizes margin but it wasn't a world changing idea but Kernalization is



As we can see the above data isn't separable so our Linear SVM won't be able to find the margin. So it fails.. If we apply Kernel- SVM with the right Kernel then we can find the right hyperplane. Hence we can infer Kernel SVM works on non-linearly separated data

## POLYNOMIAL KERNEL

Polynomial Kernel

→ Kernelization

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

(e.g.)  $K(x_1, x_2) = (1 + x_1^T x_2)^2$ , Quadratic Kernel

We'd seen that with features  $f_1, f_2$  our logistic regression couldn't find a hyperplane to separate it so we apply feature transformation to get to a feature space where a hyperplane can divide the hyperplane

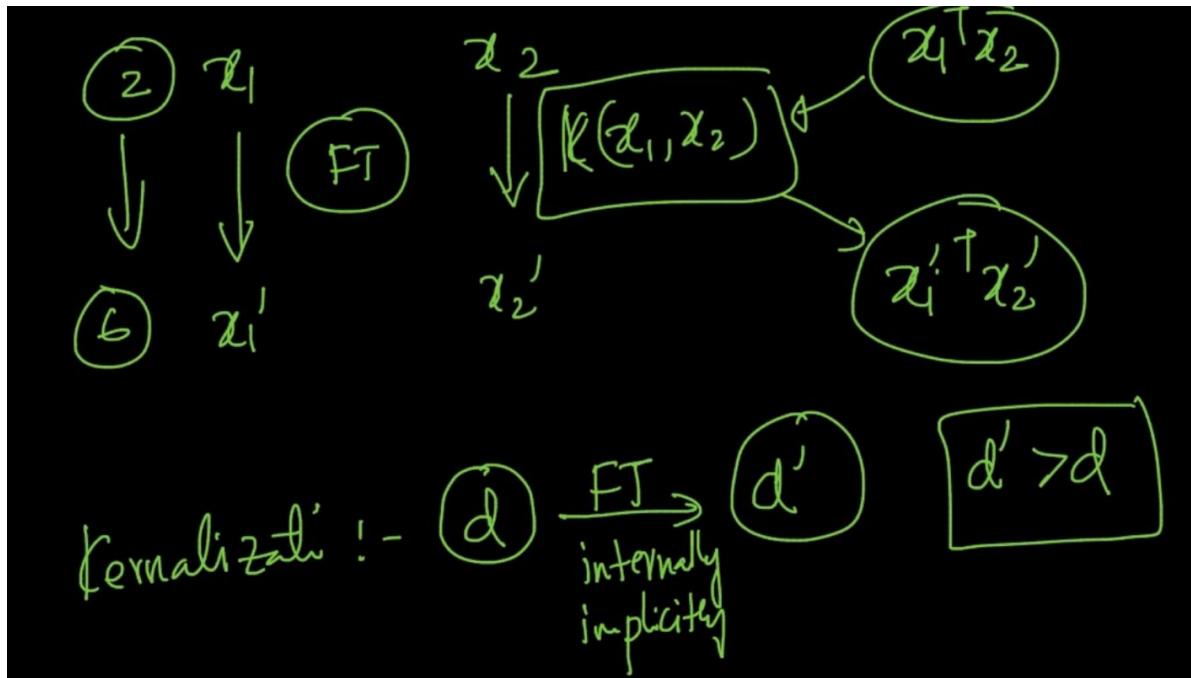
In SVM we can do that with kernelization.  $K(x_1, x_2) = (x_1^T x_2 + c)^d$  where d - degree of polynomial and c- constant

$K(x_1, x_2) = (x_1^T x_2 + 1)^2$  is called a quadratic kernel . It's a polynomial kernel with d = 2, c = 1

$$\begin{aligned}
 & \checkmark \left\{ \begin{aligned} K(x_1, x_2) &= (1 + x_1^T x_2)^2 \\ &= (1 + x_{11} x_{21} + x_{12} x_{22})^2 \end{aligned} \right. \quad \begin{aligned} x_1 &= \langle x_{11}, x_{12} \rangle \text{ (2D)} \\ x_2 &= \langle x_{21}, x_{22} \rangle \text{ (2D)} \end{aligned} \\
 &= 1 + x_{11}^2 x_{21}^2 + x_{12}^2 x_{22}^2 + 2x_{11} x_{21} + 2x_{12} x_{22} + 2x_{11} x_{21} x_{12} x_{22} \\
 &\text{Let } [1, \underline{x_{11}^2}, \underline{x_{12}^2}, \sqrt{2}x_{11}, \sqrt{2}x_{12}, \sqrt{2}x_{11} x_{12}] : \quad x_1' \\
 &\quad [1, \underline{x_{21}^2}, \underline{x_{22}^2}, \sqrt{2}x_{21}, \sqrt{2}x_{22}, \sqrt{2}x_{11} x_{22}] : \quad x_2' \\
 &= [(x_1')^T (x_2')] \quad \text{2D}
 \end{aligned}$$

We are just expanding the formula . Our  $x_1, x_2$  has 2 dimensions as seen. So after expanding it we'll be getting  $x'_1, x'_2$  so we can  $(x'_1)^T (x'_2)$

Note: Our 2D data is converted to 6D



With Kernelization ,  $x_1^T x_2 \Rightarrow K(x_1, x_2) \Rightarrow (x'_1)^T (x'_2)$ . With it our  $d$ -dimensional space we get to a new  $d'$ -dimensional space where  $d' > d$

PS - So in SVM our most important work is to find the right kernel in dataset. Ex: For the above concentric circle dataset a polynomial kernel with degree 2 is the right kernel

## RADIAL BASIS FUNCTION KERNEL (RBF)

Radial Basis Function (RBF)

SVM:- most popular / general-purpose : RBF

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

↑ hyperparameter

Self-marginal SVM       $\xrightarrow{x_1 \xrightarrow{d_{12}} x_2}$        $\|x_1 - x_2\|^2 = d_{12}^2$

( RBF Kernel )  
→ C : hyperparameter

RBF Kernel is the most popular kernel in SVM.  $K_{RBF}(x_1, x_2) = \exp\left(-\frac{\|x_1 - x_2\|^2}{2\sigma^2}\right)$  where

$\|x_1 - x_2\|$  = Distance between 2 points =  $d_{12}$  and  $\sigma$  => hyperparameter

$$K(x_1, x_2) = \exp\left(-\frac{d_{12}^2}{2\sigma^2}\right) \quad d_{12} = \|x_1 - x_2\|_2$$

①  $d_{12} \uparrow ; K(x_1, x_2) \downarrow$  Similarity

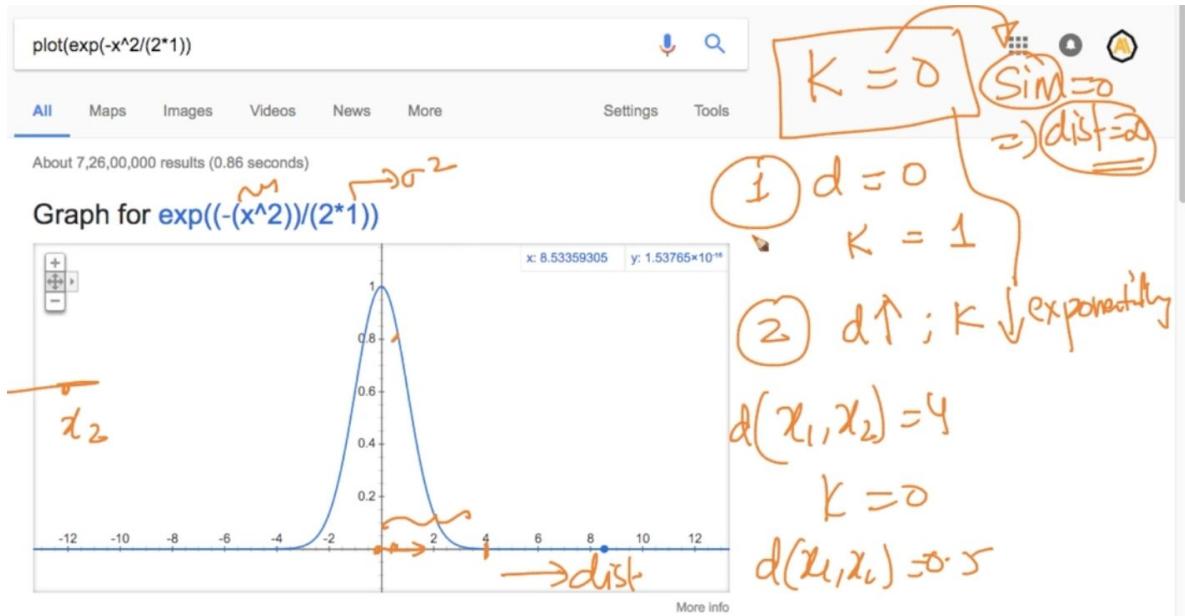
$\xrightarrow{x_1 \xrightarrow{d_{12}} x_2 \xrightarrow{d_{13}} x_3}$

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

$\frac{1}{e^{d^2/2\sigma^2}}$

$d \uparrow ; d^2 \uparrow$   
 $e^{d^2/2\sigma^2} \uparrow$   
 $\frac{1}{e^{d^2/2\sigma^2}} \downarrow$

As  $d_{12} \uparrow$ ;  $K_{RBF}(x_1, x_2) \downarrow$  i.e similarity  $\downarrow$ . In fig,  $K(x_1, x_2) > K(x_1, x_3)$  as  $d_{12} < d_{13}$



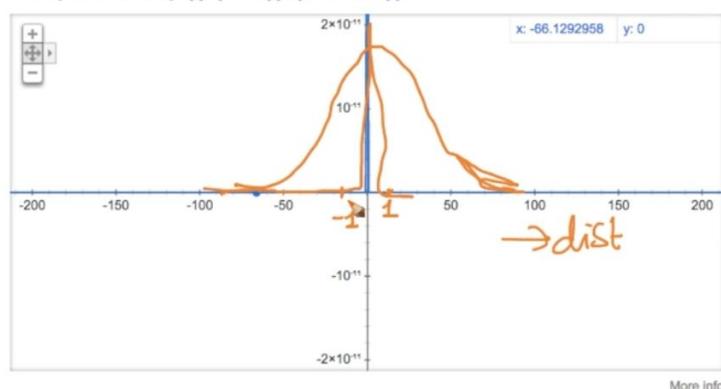
Here,  $\sigma = 1$ . So when  $d = 0$ ,  $K = 1$  but when  $d \uparrow$ ;  $K \downarrow$  exponentially. As we can see as the distance increases the Kernel function's value gets decreased.

NOTE - We can intuitively say that Kernel is like Similarity function i.e when two points are similar their value is greater and distance is like dissimilarity i.e if distance is greater the points are dissimilar

NOTE 2 - RBF ~ Gaussian PDF

About 2,30,00,000 results (0.51 seconds)

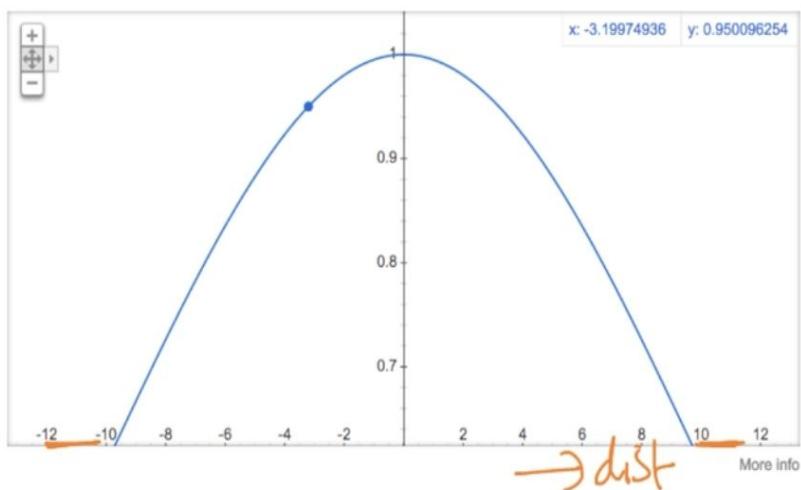
Graph for  $\exp\left(-\frac{(x^2)}{2 \cdot 0.01}\right)$



When we change  $\sigma = 1$  to 0.1 the graph gets more steep. Even when  $\text{dist} > 1$ ;  $K = 0$

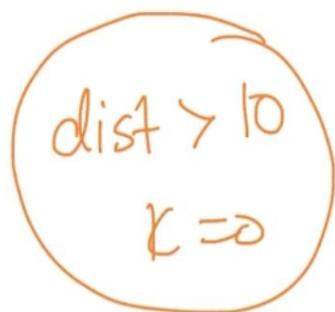
About 2,37,00,000 results (0.60 seconds)

Graph for  $\exp(-(x^2)/(2*100))$

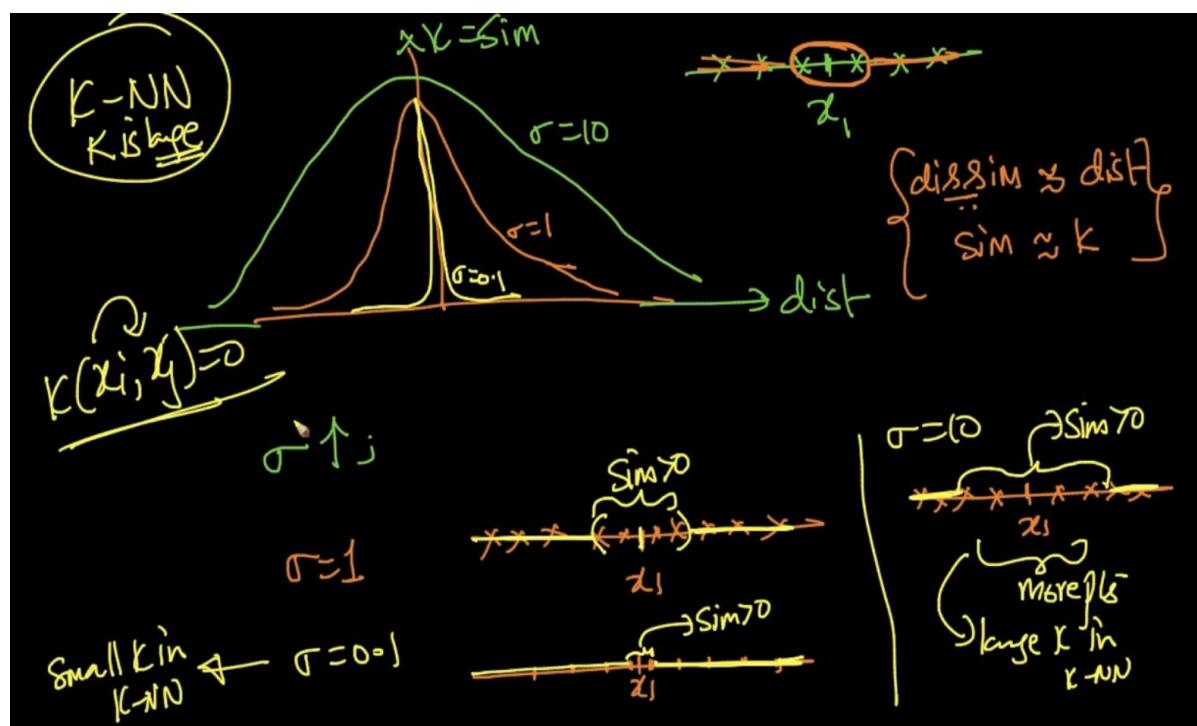


$$\sigma = 10$$

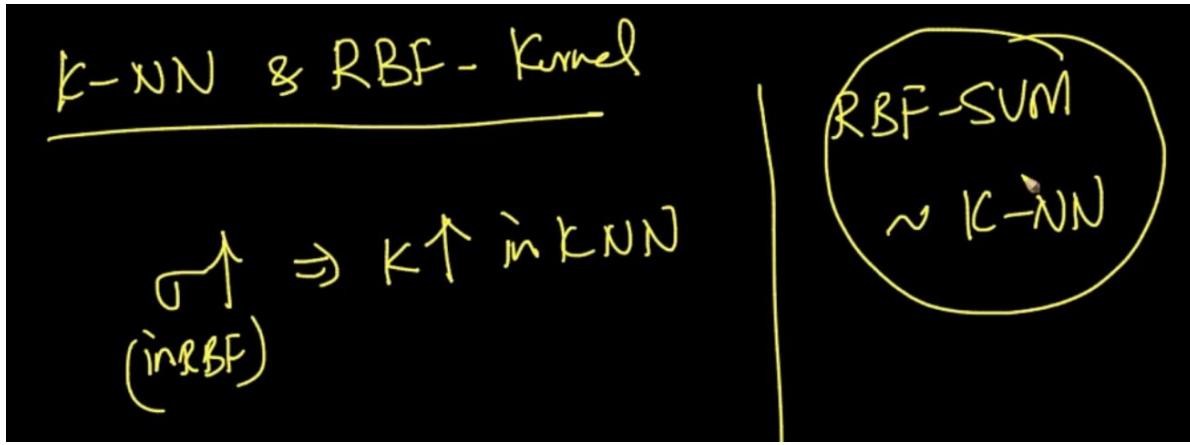
$$\sigma^2 = 100$$



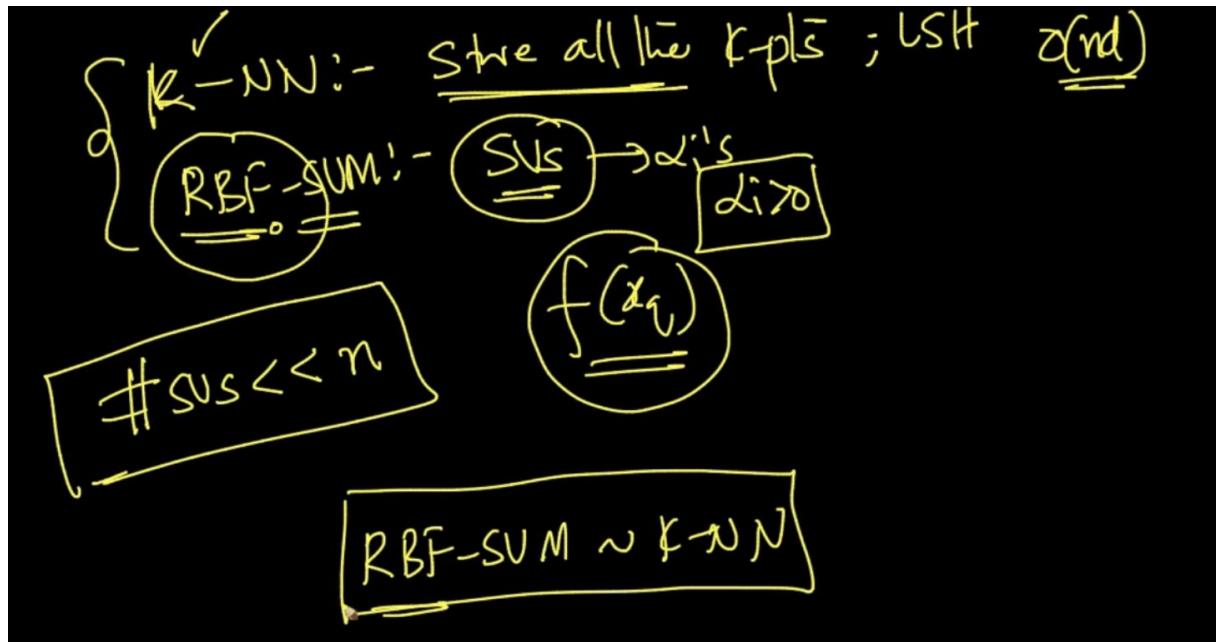
For  $\sigma = 10$



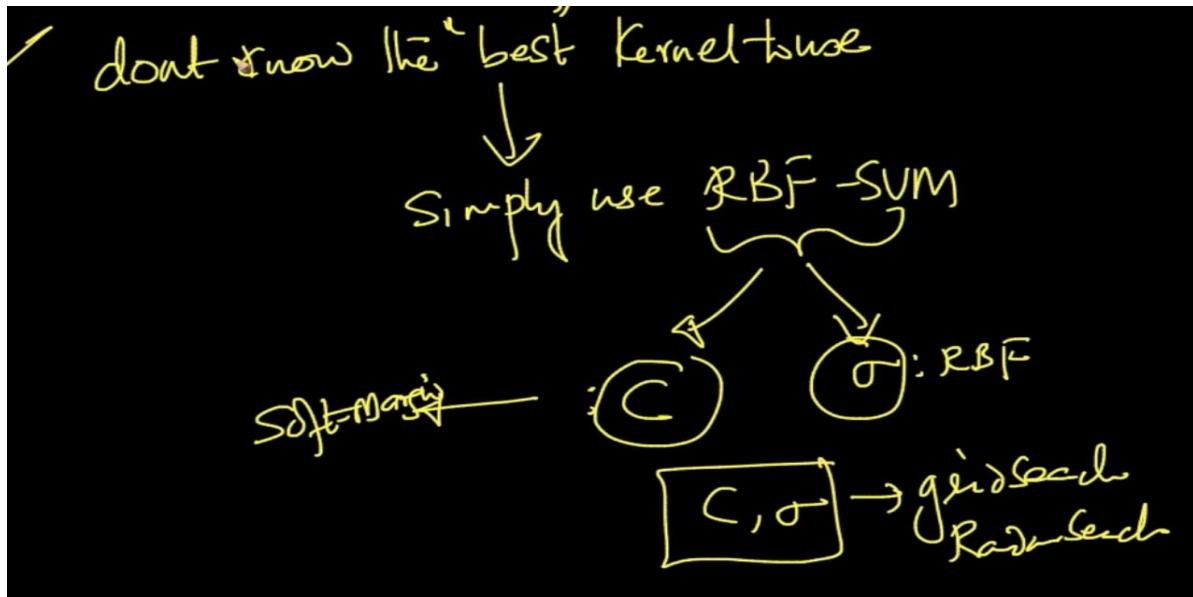
When  $\sigma$  is small we take only the points which are very close and we can take more points when  $\sigma$  is large . So  $\sigma$  can be compared with k in K-NN



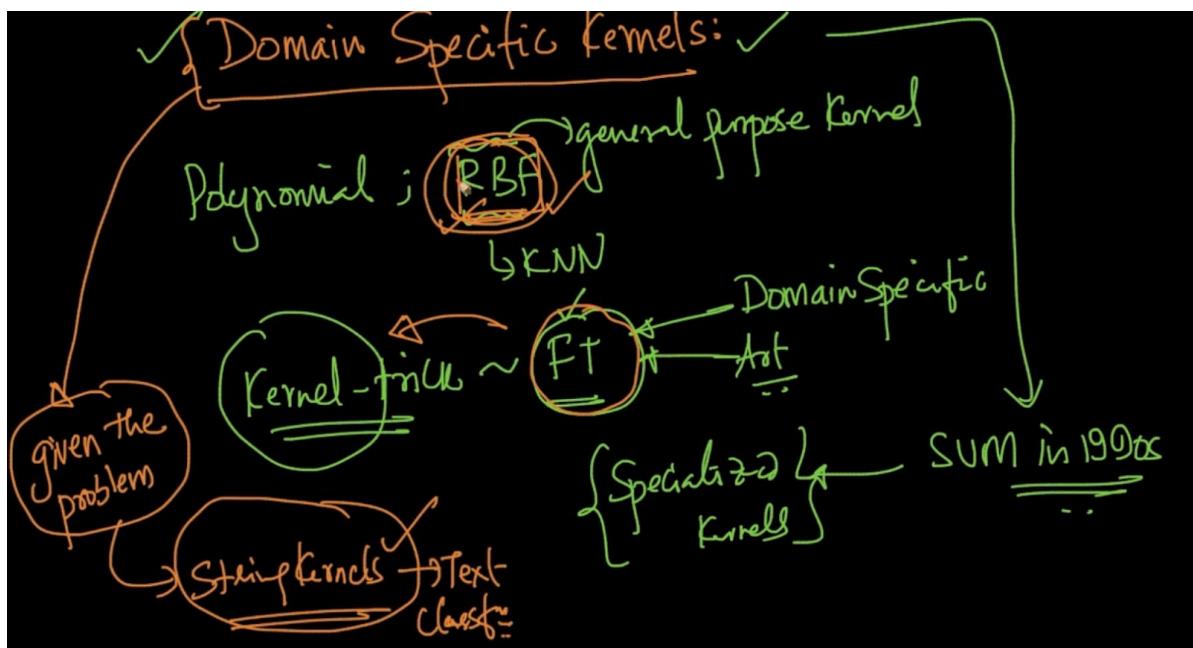
So we can say that RBF-SVM  $\sim$  k-NN



In K-NN we store all the k-Points but in RBF-SVM we store only the support vectors because for  $\alpha > 0$  we only need Support Vectors.

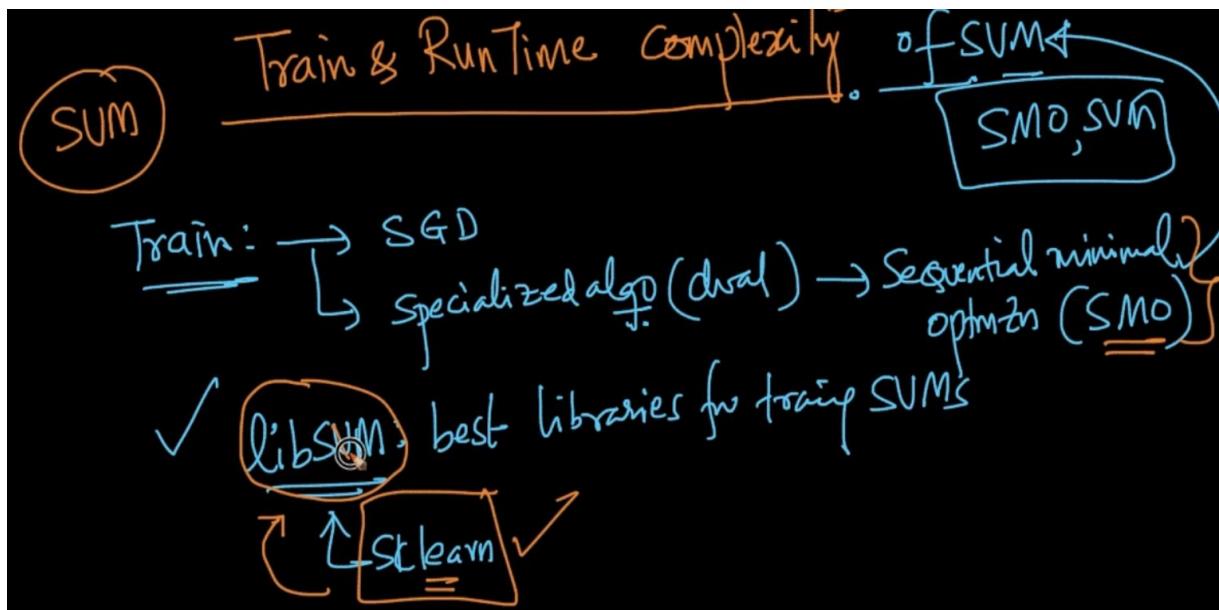


So it can work pretty well in large datasets .Hence , if we don't know which Kernel to use we use RBF-SVM . We do grid search/ random search for both C,  $\sigma$

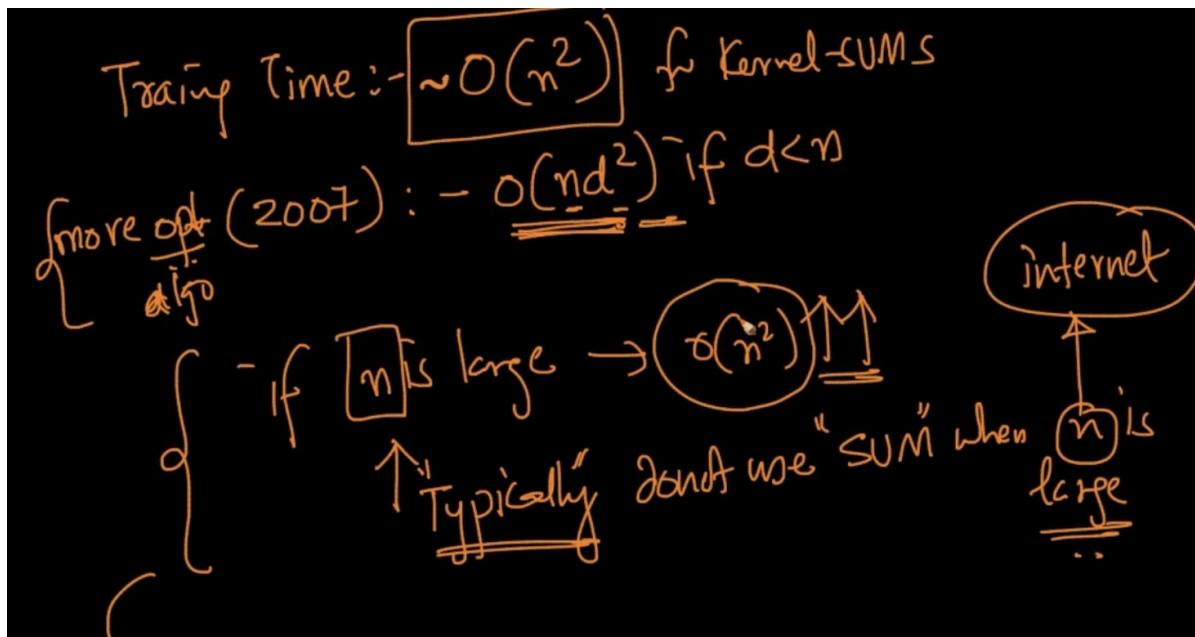


There are domain specific kernels for SVM i.e String Kernel for Text Classification, genome Kernel for Bioinformatics and so on...So, always look the appropriate Kernel i.e if we have Amazon Reviews data it is a text classification task so we can use String Kernel and our RBF-Kernel is obviously there for the backup plan

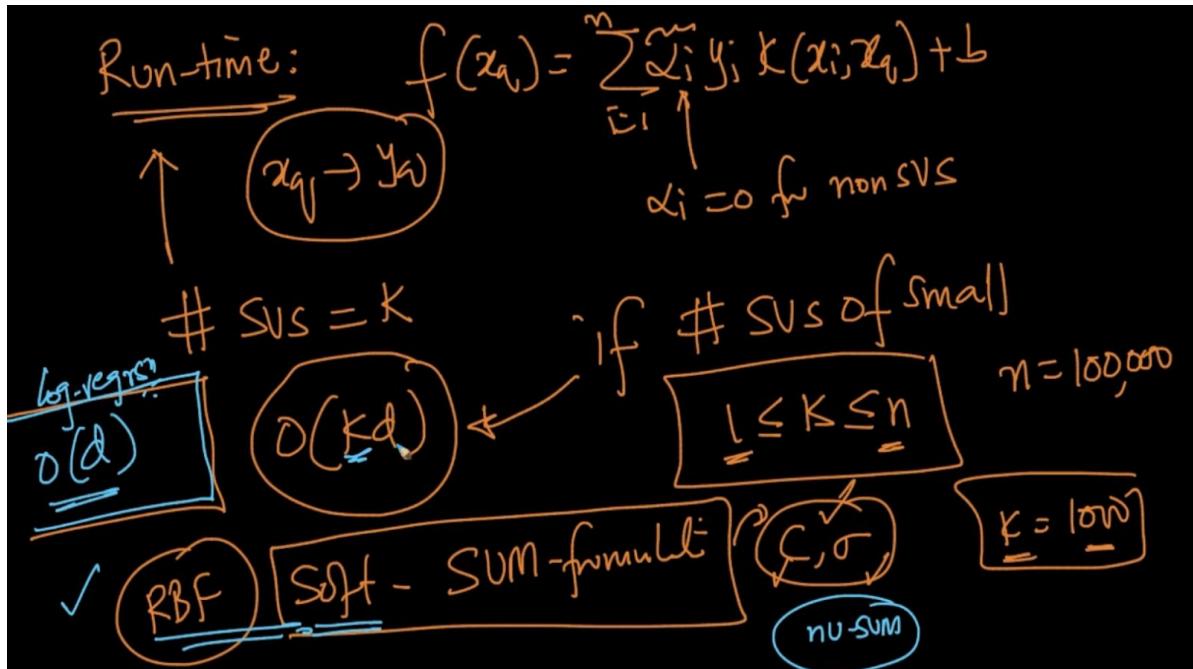
## TRAIN AND RUN-TIME COMPLEXITIES



libSVM is one of the best library for training SVM



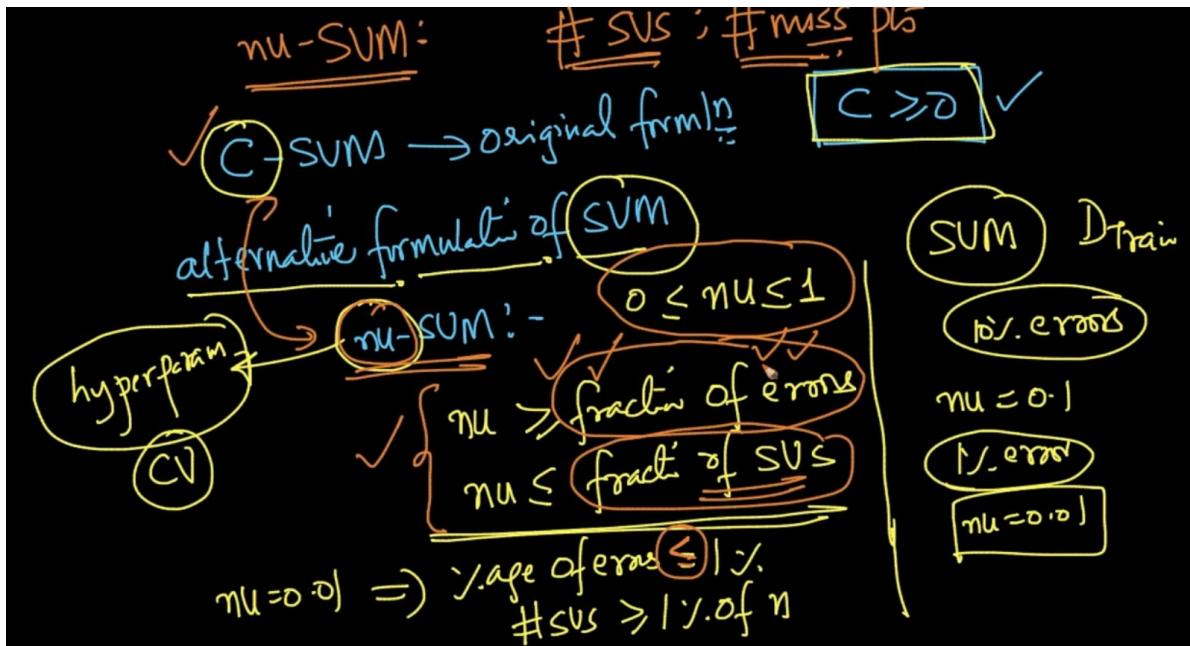
Training Time :  $O(n^2)$  for Kernel- SVMs so when  $n$  is large we typically don't use it since it takes significant time



Runtime Complexity :  $O(kd)$  where  $k$  - No. of Support Vectors and  $d$ - No. of dimensions

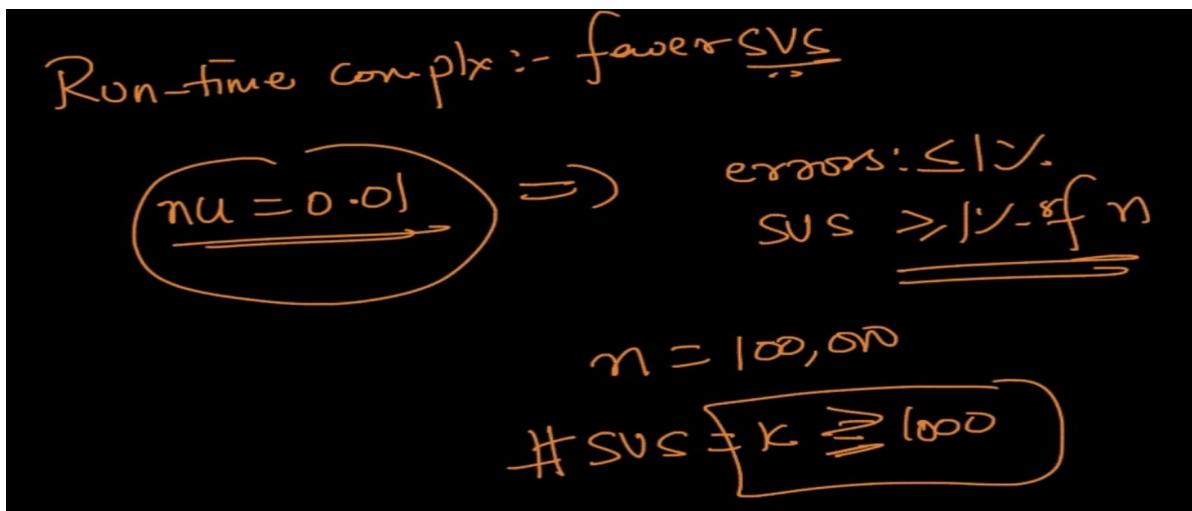
For logistic Regression it's  $O(d)$  which is less than SVM

## nu-SVM CONTROL ERRORS AND SUPPORT VECTORS



We studied C-SVM till now but  $C \geq 0$  so we can't get an accurate C so we use nu-SVM where  $0 \leq \text{nu} \leq 1$  and  $\text{nu} \geq \text{fraction of errors}$ ,  $\text{nu} \leq \text{fraction of Support Vectors}$

So , nu is the threshold of errors i.e suppose we want just 10 % errors just use  $\text{nu} = 0.1$



In run-time complexity we want fewer SVs so if  $\text{nu} = 0.01$  it guarantees errors  $\leq 1\%$  and also Support Vectors  $\geq 1\%$  of n so if  $n = 100,000$   $\text{nu} = 0.01$  , Support Vectors will guaranteed be  $\geq 1000$

## SVM REGRESSION

✓ Support Vector Regression (SVR)  $y_i \in \mathbb{R}$

SVM - Classfn:- SVC  $\rightarrow y_i \in \{+1, -1\}$

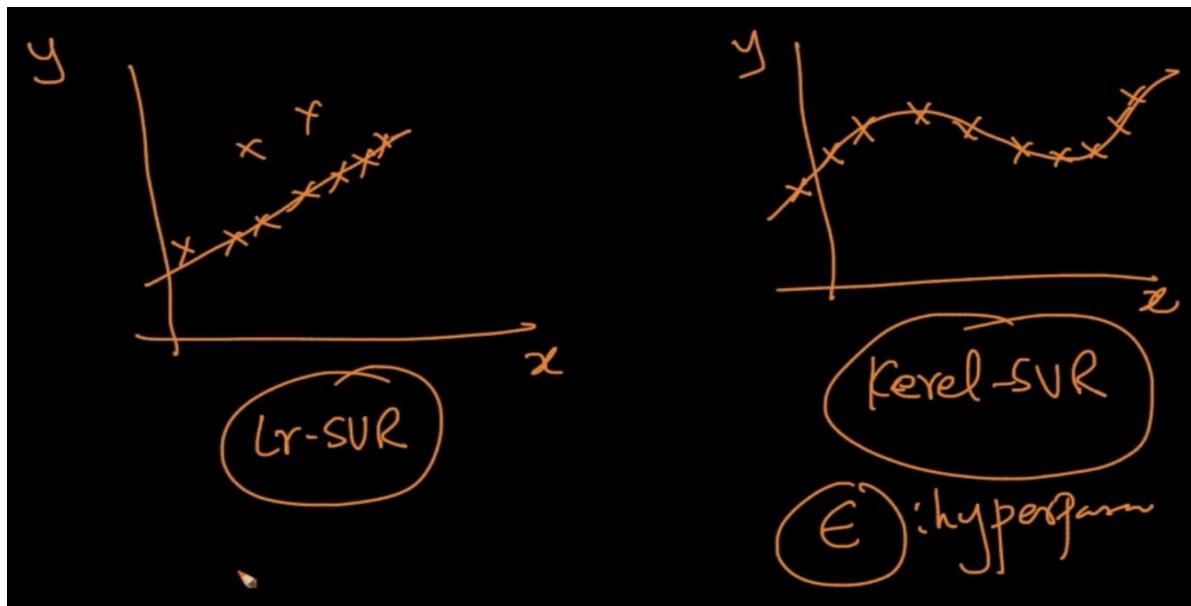
Math:  $\left\{ \begin{array}{l} \text{Lr. form of SVR} \\ \min_{w,b} \frac{1}{2} \|w\|^2 \xrightarrow{\text{reg}} \\ \text{s.t. } y_i - (w^T x_i + b) \leq \epsilon \\ \quad (w^T x_i + b) - y_i \leq \epsilon \\ \quad (\epsilon > 0) \end{array} \right.$

$f(x_i) = w^T x_i + b = y_i$

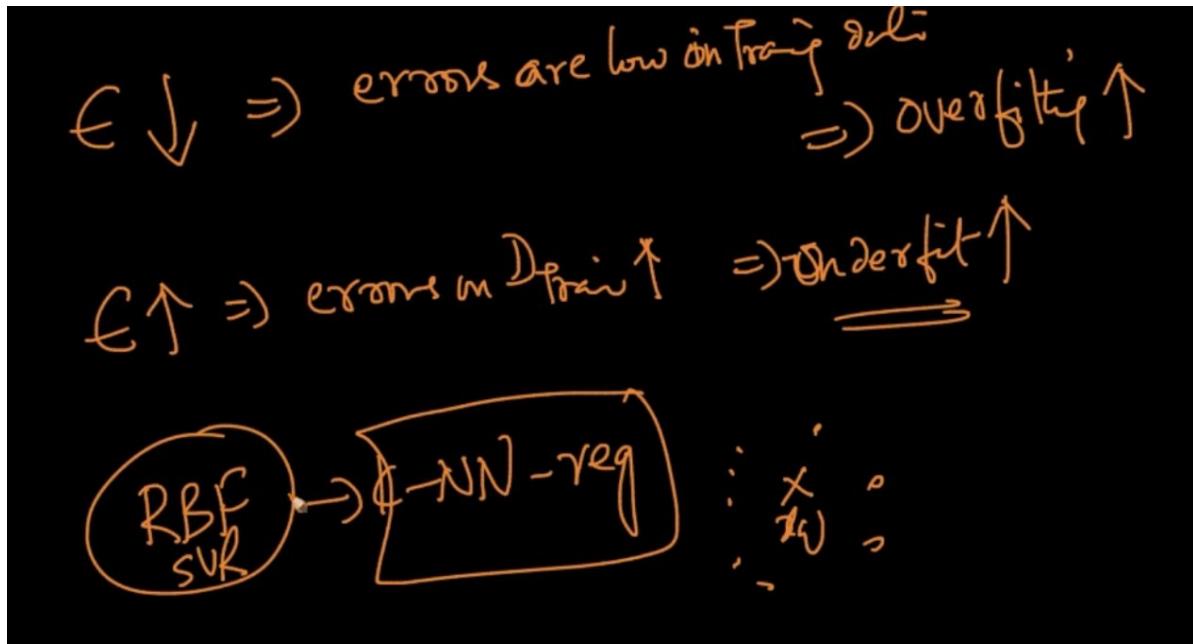
$y_i - \hat{y}_i \leq \epsilon$

$\hat{y}_i - y_i \leq \epsilon$

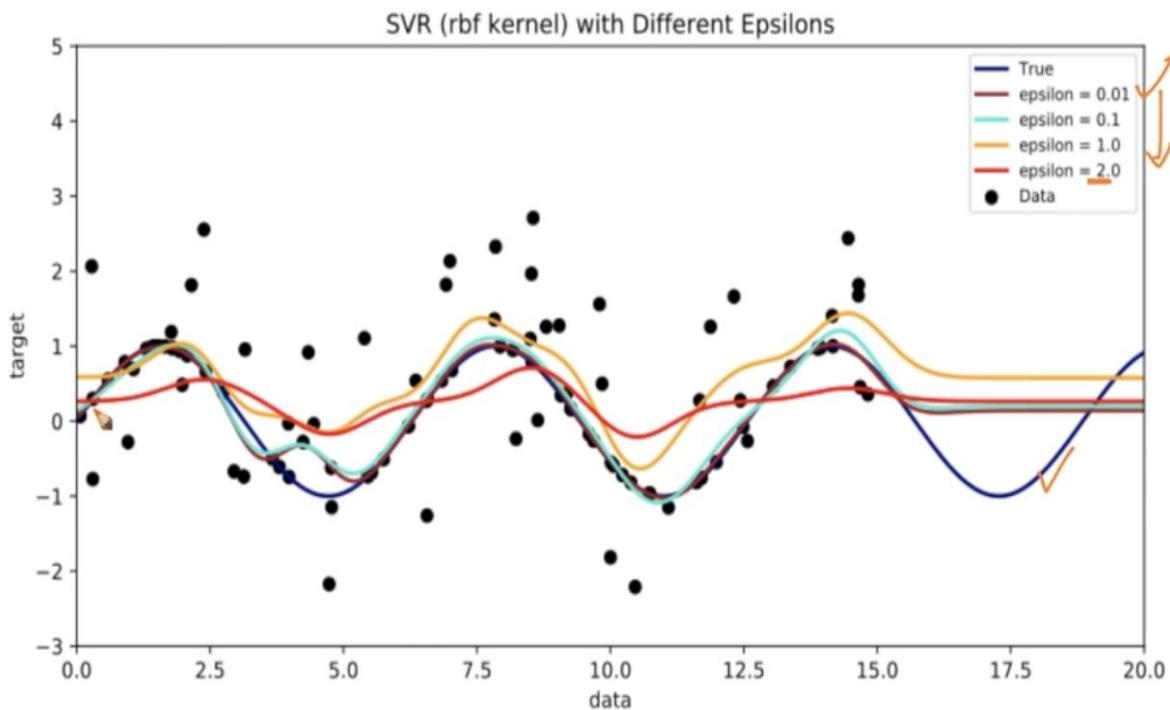
In SVM-Regression,  $\min \frac{1}{2} \|w\|^2 \Rightarrow$  regularizer such that  $y_i - (w^T x_i + b) \leq \epsilon$  and  $(w^T x_i + b) - y_i \leq \epsilon$ .  $\epsilon \geq 0$  and  $\epsilon$  is a hyper-parameter



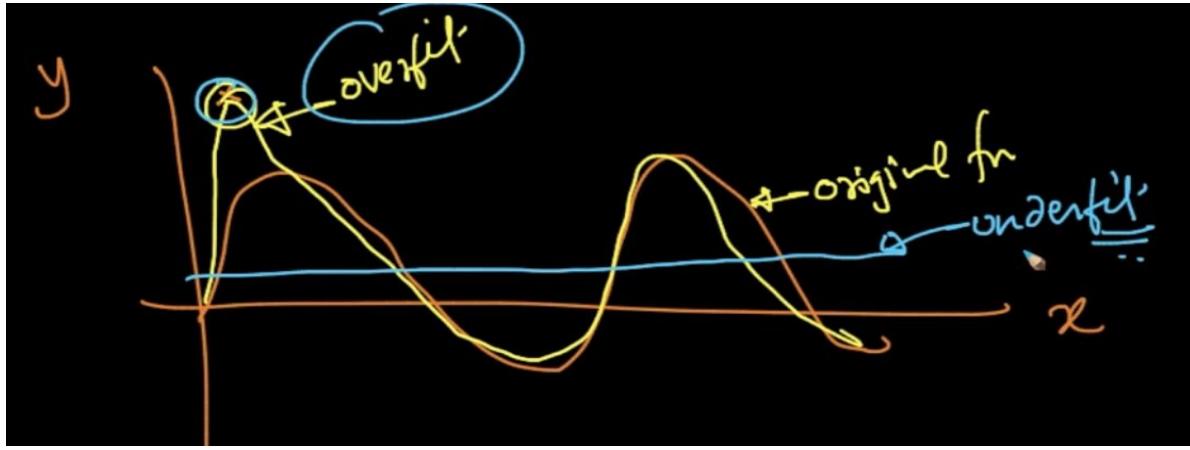
Linear SVR can only fit linear data and Kernel-SVR can fit any non linear data



Effect of  $\epsilon$  on errors

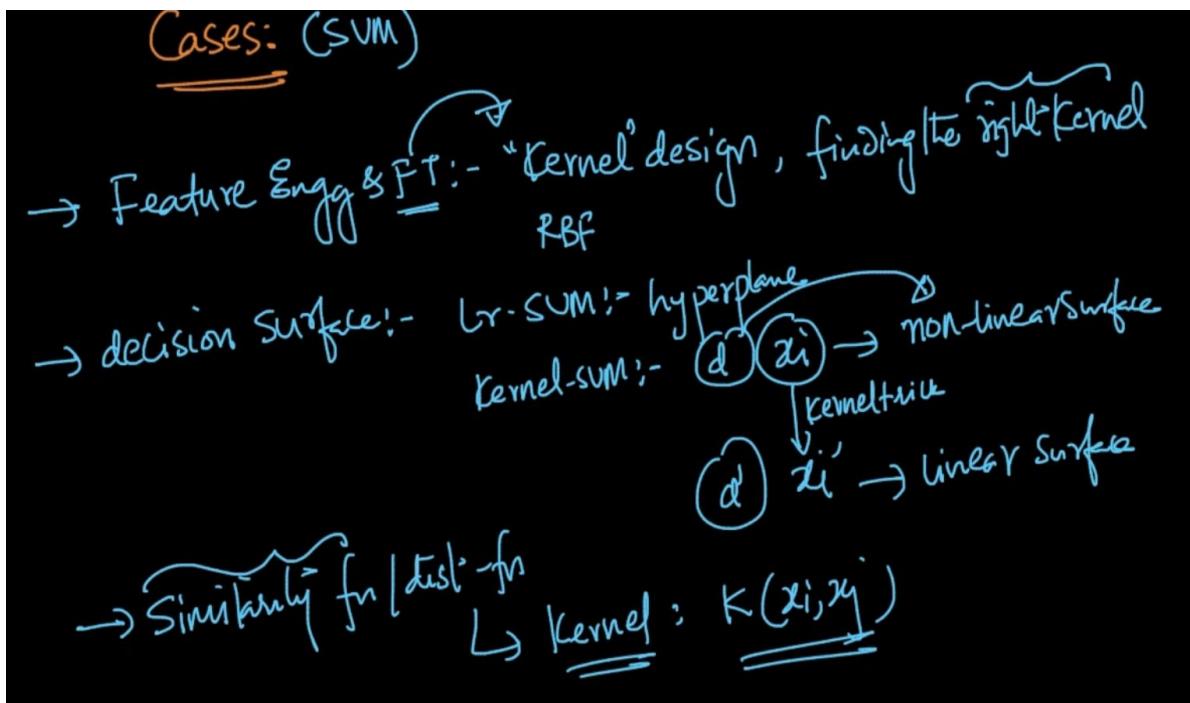


Brown line is overfitting , it has  $\epsilon = 0.01$  and Red line is underfitting  $\epsilon = 2.0$

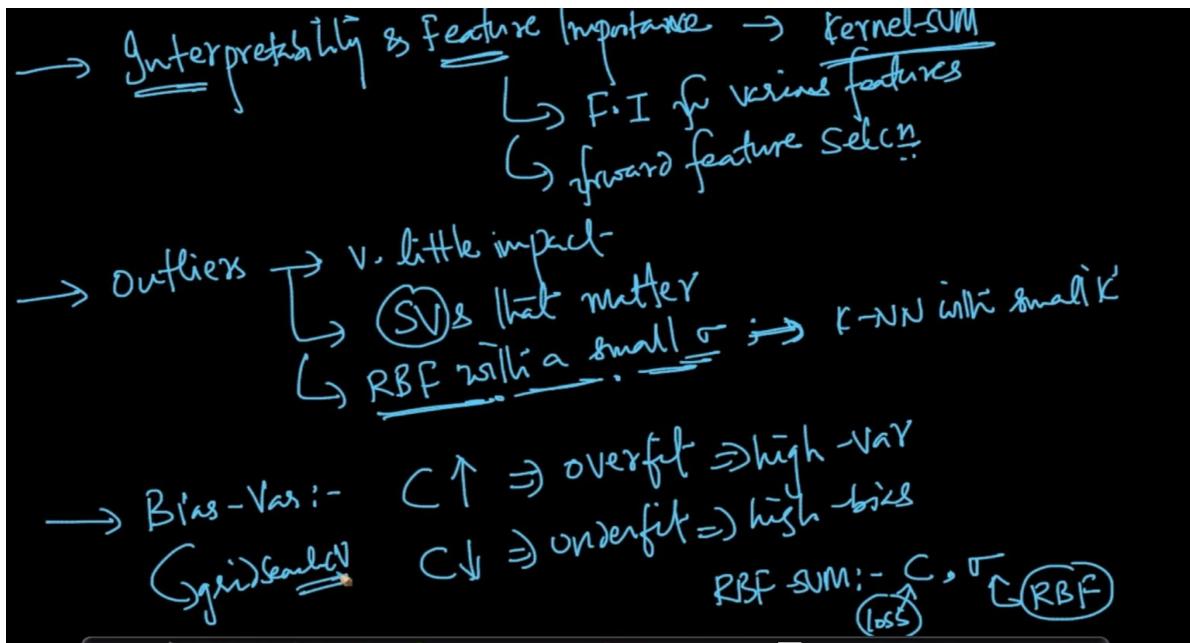


When overfitting happens the function distorts for a single point and in underfitting it just doesn't care.

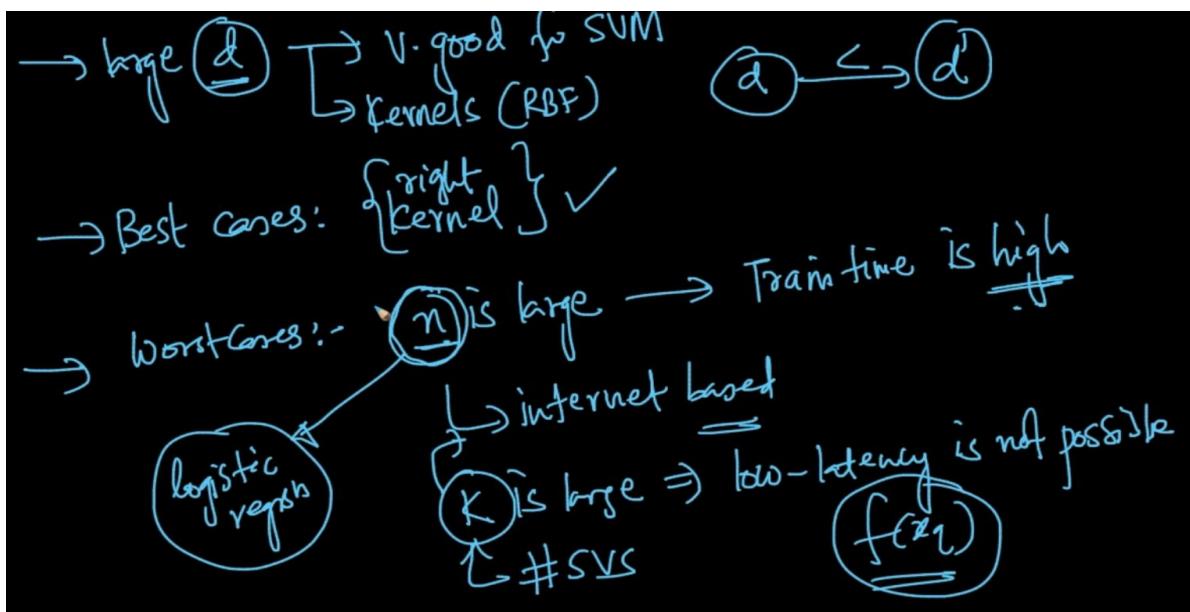
## CASES



- 1) Feature Transformation - In SVM, feature transformation is changed to Kernel Design/finding the right kernel
- 2) Decision Surface - In kernel-SVM we don't get a linear hyperplane with  $x_i$  but with kernel trick it is transformed to  $x'_i$  to a different dimensional which is linearly separable
- 3) Similarity Function - If it's given then it's very easy to convert into Kernel



- 4) Interpretability and Feature importance is very difficult in Kernel-SVM. For Linear SVM we can always apply the same methods as logistic regression
- 5) Outliers has very little impact on SVM
- 6) Bias Variance - We can find the ideal parameters : C,  $\sigma$  using gridSearchCv



When we have a large dimensional data SVM works like a charm. Best case for SVM is Finding the right Kernel and worst case is when 'n' or K is large