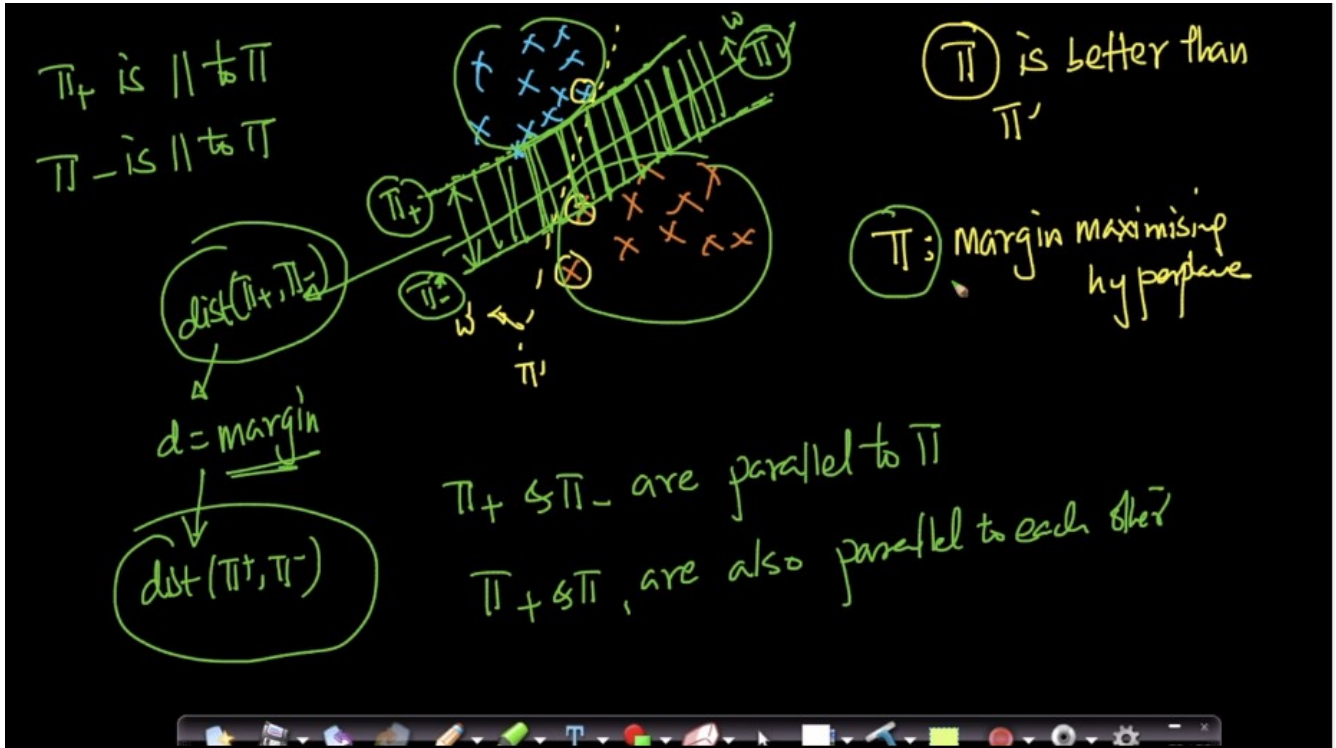


NOTES

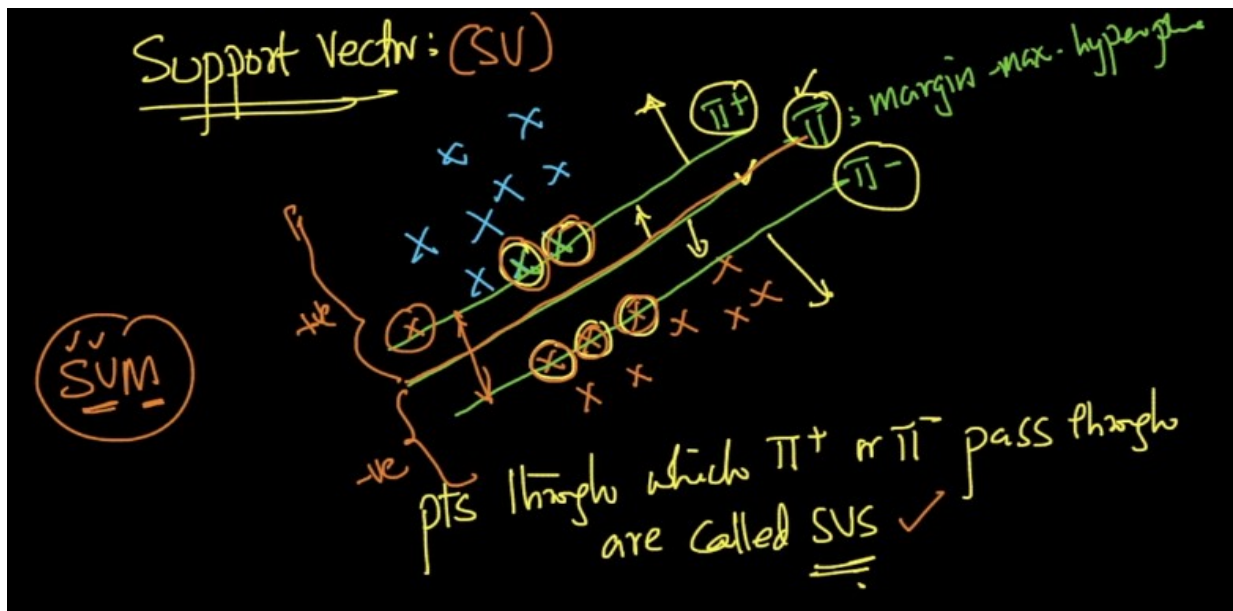
Support vector machines:

Let's assume the data is linearly separable. There will be many hyper planes that can separate the data.

Fitting the hyper plane that can separate the two classes of data as widely as possible.



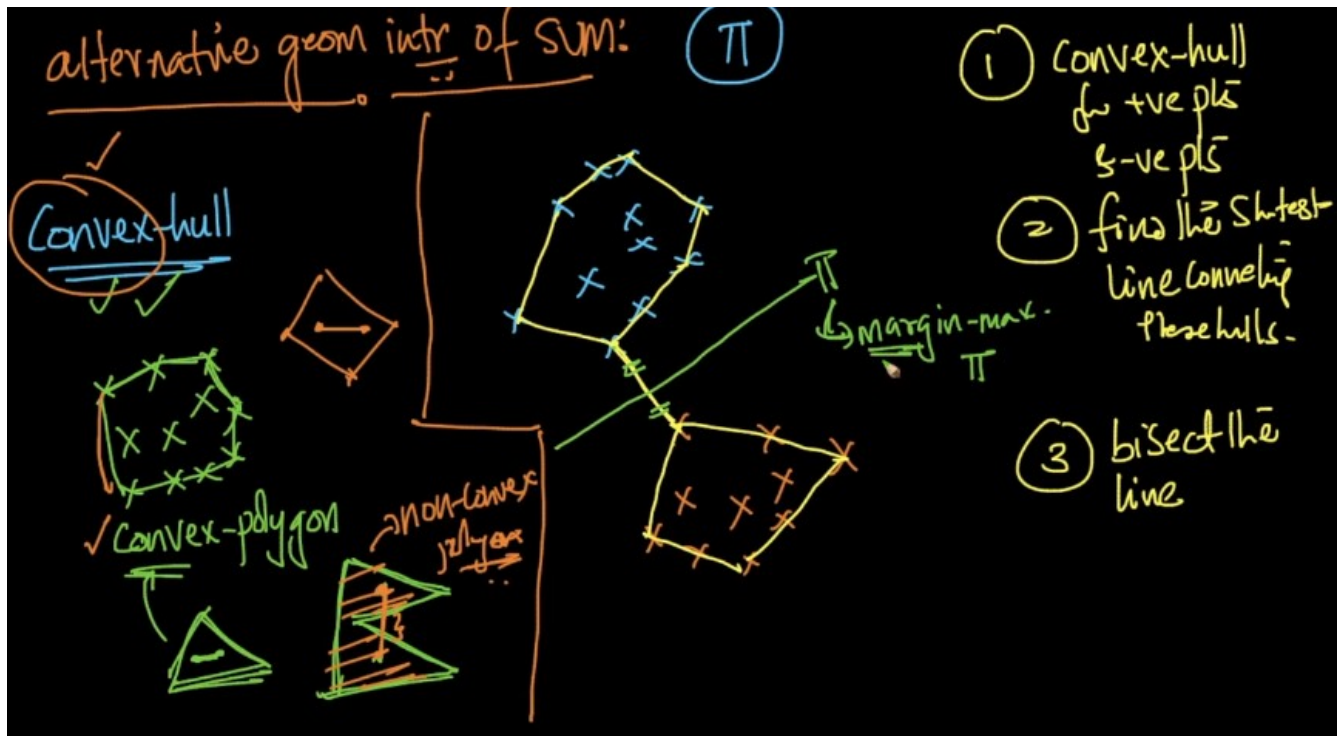
SVM's try to find the hyper plane that maximizes the margin = $dist(p_i^+, p_i^-)$, As margin increases the generalization error decrease.



Alternative geometric intuition of SVM:

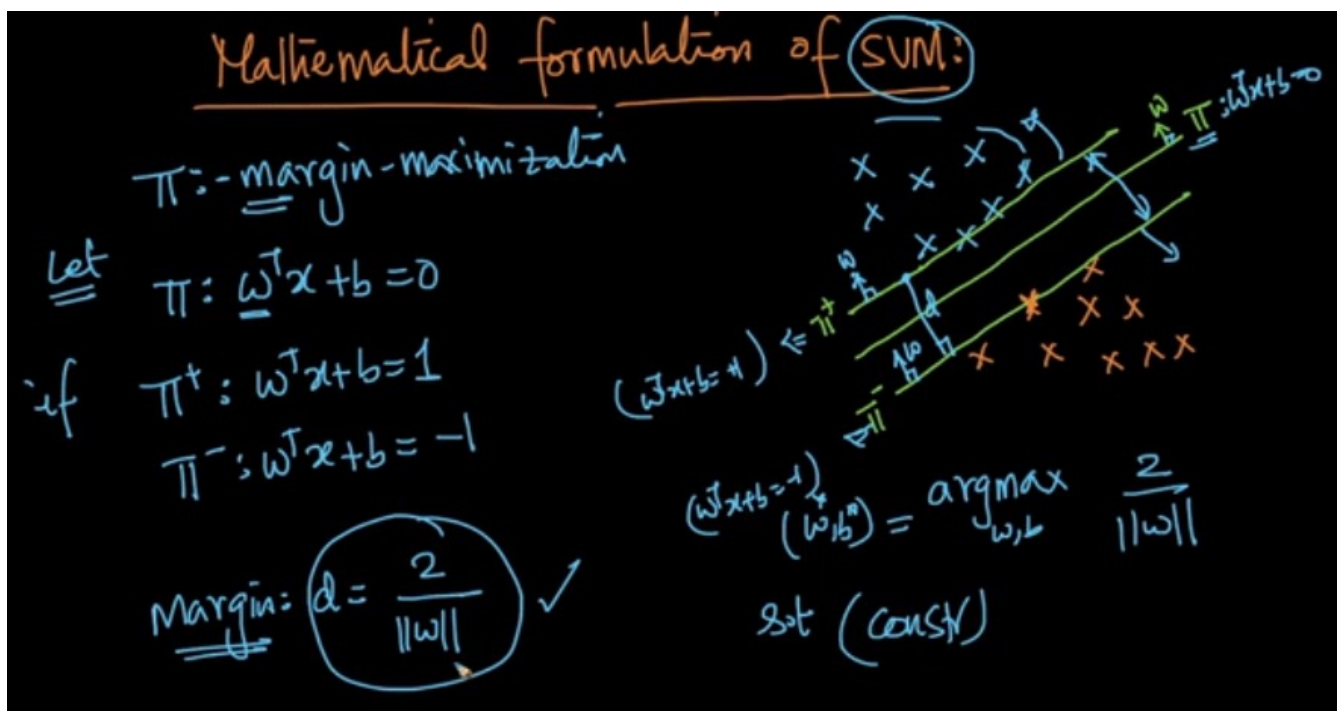
First we make a convex polygon:

The points in the convex hull must be inside or on the convex hull.



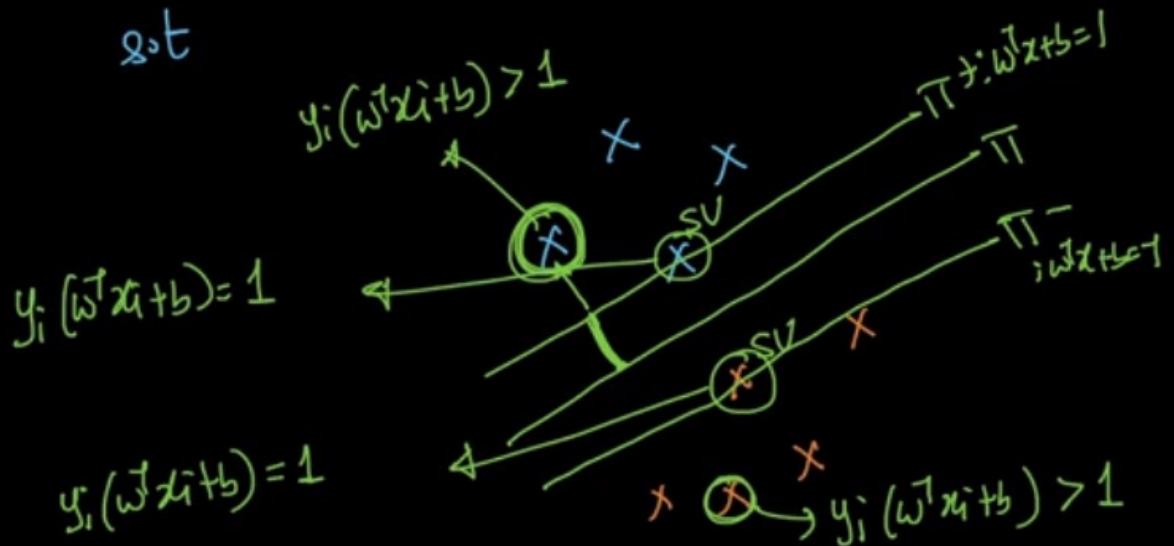
Mathematical derivation:

Objective: We want to find a hyper plane that does margin maximization.



$$(\mathbf{w}^*, b^*) = \arg \max_{\mathbf{w}, b} \frac{2}{\|\mathbf{w}\|} = \text{margin}$$

s.t



The optimization function is find \mathbf{w}^* , b^* .

$$(\mathbf{w}^*, b^*) = \arg \max_{\mathbf{w}, b} \frac{2}{\|\mathbf{w}\|} \leftarrow \text{margin}$$

$$\text{s.t } \underline{y_i(\mathbf{w}^T \mathbf{x}_i + b) \geq 1 \text{ for all } \mathbf{x}_i}$$

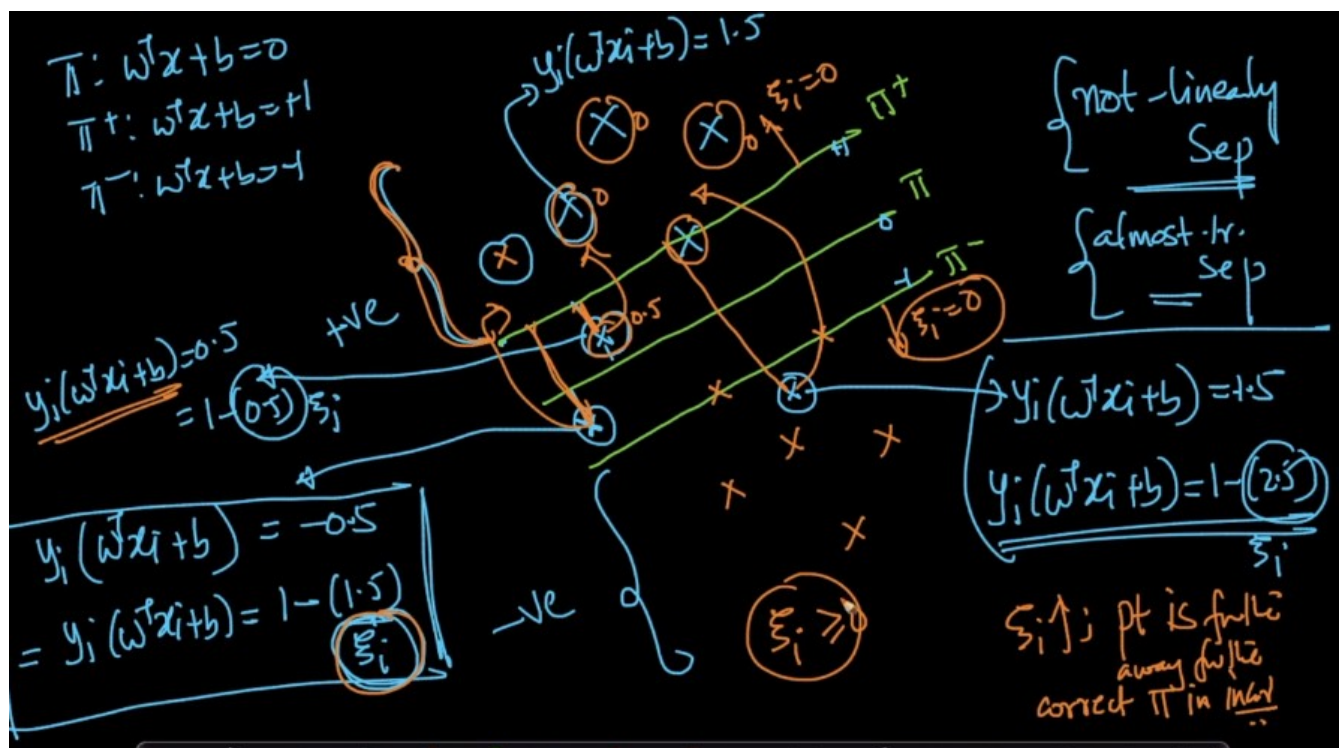
$\underbrace{\quad}_{n\text{-const.}}$

$$\text{Constr. optmzn. prob. of SVM} \left\{ \begin{aligned} & \omega, b^* = \underset{\omega, b}{\operatorname{argmax}} \frac{2}{\|\omega\|} \\ & \text{s.t } \forall i, y_i(\omega^T x_i + b) \geq 1 \end{aligned} \right.$$

This works only when the data is linearly separable, this is also called as hard margin SVM.

Edge cases:

Not – linearly separable data, alternate formulation of the edge points. As ξ_i increases the points will be away from the SVM.



$x_i \rightarrow \xi_i$
 $\checkmark \quad \xi_i = 0$ if $y_i(w^T x_i + b) \geq 1$ ↑ correctly classified
 π^+ & π^-
 $\xi_i > 0$ & it is equal to the same units of
dist away from the correct hyperplane
 in the incor dir

Changes to the formulation.

$(w, b) = \arg \max_{w, b} \left(\frac{2}{\|w\|} \right) = \arg \min_{w, b} \left(\frac{\|w\|}{2} \right)$
 $\max_x f(x) = \min_x \frac{1}{f(x)}$
 $\xi_i = 1$

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

(Annotations: $\frac{\|\vec{w}\|^2}{2}$ is labeled "margin"; $C \cdot \frac{1}{n} \sum \xi_i$ is labeled "avg. dist miss. pts for correct TTS"; $\xi_i > 0$ is labeled "hyperplane"; $y_i(\vec{w} \cdot \vec{x}_i + b) > 1 - \xi_i \quad \forall i$ is labeled "corr. class. pts $\xi_i = 0$ " and "incorr. class. pts $\xi_i > 0$ ");

Minimize error = min misclassification

$\min \sum \xi_i$

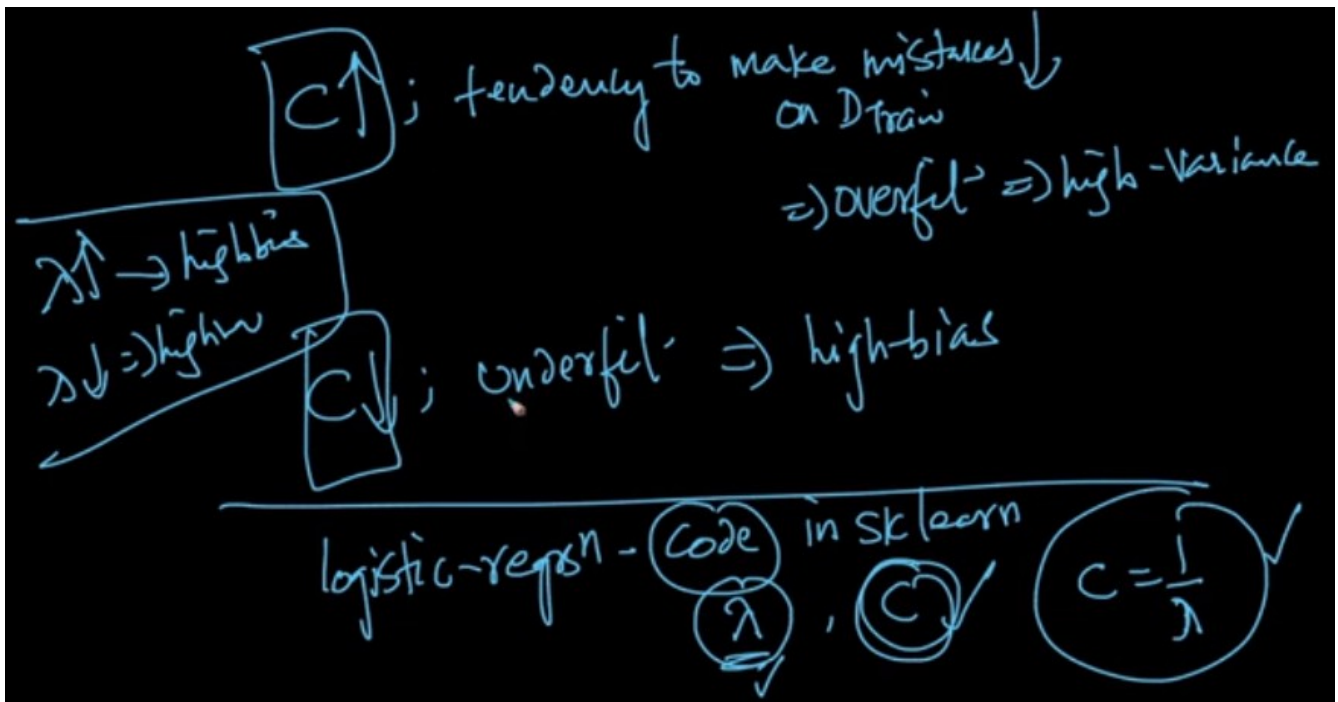
25:42

Modified loss function:

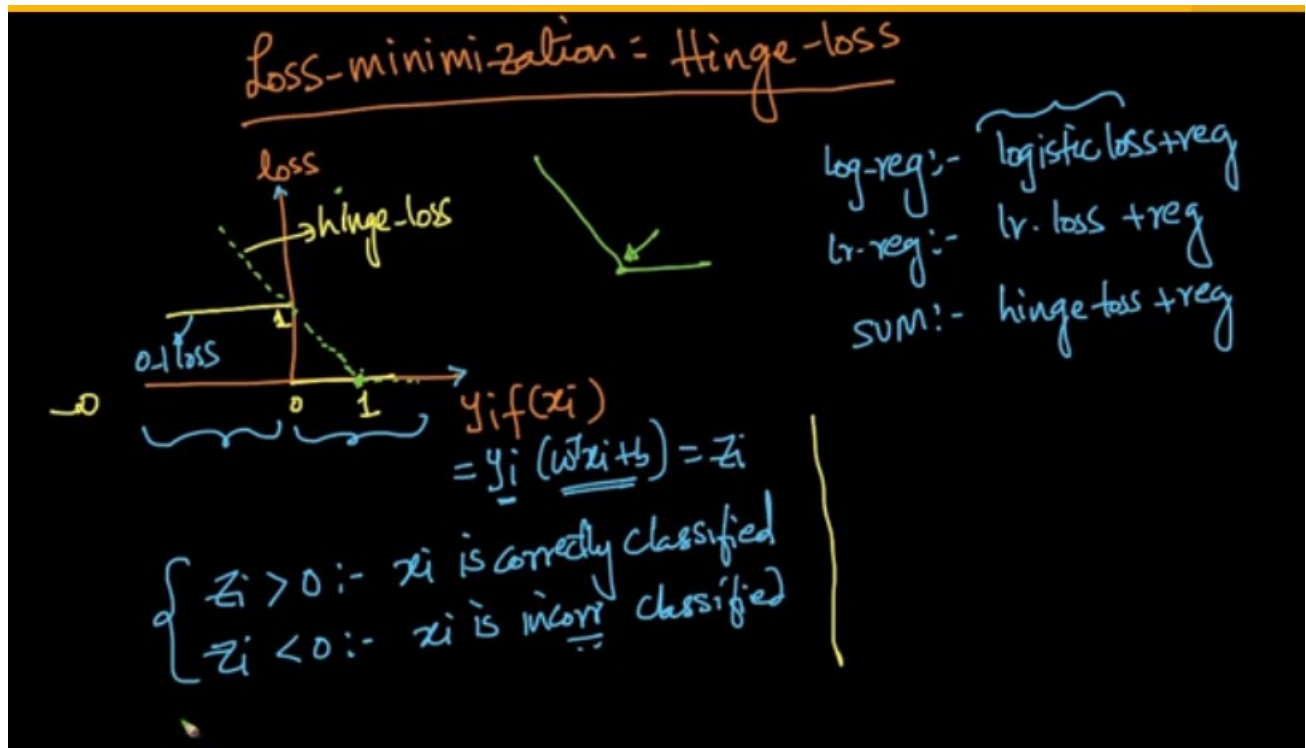
$(\vec{w}, b) = \underset{w, b}{\operatorname{argmin}} \left(\frac{\|\vec{w}\|^2}{2} + C \cdot \frac{1}{n} \sum_{i=1}^n \xi_i \right)$

(Annotations: $\frac{\|\vec{w}\|^2}{2}$ is labeled "margin"; $C \cdot \frac{1}{n} \sum \xi_i$ is labeled "loss"; $\frac{\|\vec{w}\|^2}{2}$ is also labeled "reg" (regularization); $C \cdot \frac{1}{n} \sum \xi_i$ is labeled "hyperplane"; $y_i(\vec{w} \cdot \vec{x}_i + b) > 1 - \xi_i \quad \forall i$ is labeled "constraint"; $\xi_i > 0$ is labeled "hyperplane"; $\min_{w, b} C \cdot \frac{1}{n} \sum \xi_i$ is labeled "minimize loss + reg");

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



Loss function(Hings Loss) based interpretation:



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

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

$\begin{cases} \text{Case 1: } z_i \geq 1 \Rightarrow 1 - z_i \text{ is -ve value} \Rightarrow \max(0, 1 - z_i) = 0 \\ \text{Case 2: } z_i < 1; 1 - z_i > 0 \Rightarrow \max(0, 1 - z_i) = 1 - z_i \end{cases}$

This is alternate formulation of Hinge Loss.

geom-formulation of loss-min.

Corr class's pt $\rightarrow \sum_i z_i = 0 \leftarrow z_i$

$d_j = 1 - \underbrace{y_i (\omega^T x_i + b)}_{z_j} = 1 - z_j$

$\xi_j = \text{dist fr } x_j \text{ to the } \Pi^+ = d_j = 1 - z_j$

$(\omega^T x_j + b) : -ve$

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

Soft sum:

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

s.t. $(1 - y_i(w^T x_i + b)) \geq \xi_i \quad \forall i$
 $\xi_i \geq 0$

$\begin{cases} C \uparrow \Rightarrow \text{overfit} \\ C \downarrow \Rightarrow \text{underfit} \end{cases}$

loss-min:

$$\min_{w, b} \left[\sum_{i=1}^n \max(0, 1 - y_i(w^T x_i + b)) \right] + \lambda \|w\|^2$$

$\begin{cases} \lambda \uparrow \Rightarrow \text{underfit} \\ \lambda \downarrow \Rightarrow \text{overfit} \end{cases}$

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

Dual for of SVM formulation:

Dual form of SVM:

soft-margin SVM

$$\min_{w, b} \frac{1}{2} \|w\| + C \sum_{i=1}^n \xi_i$$

s.t. $y_i(w^T x_i + b) \geq 1 - \xi_i \quad \forall i$
 $\xi_i \geq 0$

Equivalent

Dual

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

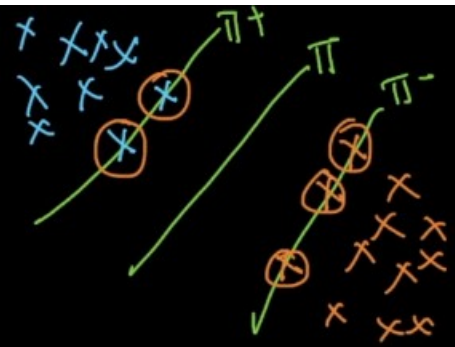
s.t. $\alpha_i \geq 0$ (4) $\alpha_i \geq 0$ only for SVs
 $\sum_{i=1}^n \alpha_i y_i = 0$

Primal of SVM

$\checkmark x_{qj} = \frac{w^T x_j + b}{\|w\|} = f(x_q)$

(1) $x_i \rightarrow \alpha_i$
 (2) x_i 's only occur in the form of $x_i^T x_j$
 (3) $f(x_q) = \sum_{i=1}^n \alpha_i y_i x_i^T x_q + b$

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



SVs :- $\alpha_i > 0$

non SVs :- $\alpha_i = 0$

$f(x_q)$:- only pts that matter are SVs

SVM

Kernel Trick:

$$\begin{cases} \max_{\alpha} \sum_{i=1}^n \alpha_i - \frac{1}{2} \sum_{i=1}^n \sum_{j=1}^n \alpha_i \alpha_j y_i y_j x_i^T x_j \\ \text{s.t. } \sum_{i=1}^n \alpha_i y_i = 0, \alpha_i \geq 0 \end{cases}$$

$x_i^T x_j \rightarrow \text{Sim}(x_i, x_j)$
 $\rightarrow K(x_i, x_j)$
 \hookrightarrow kernel fn

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

\rightarrow The most imp. idea in SVM is kernel-trick

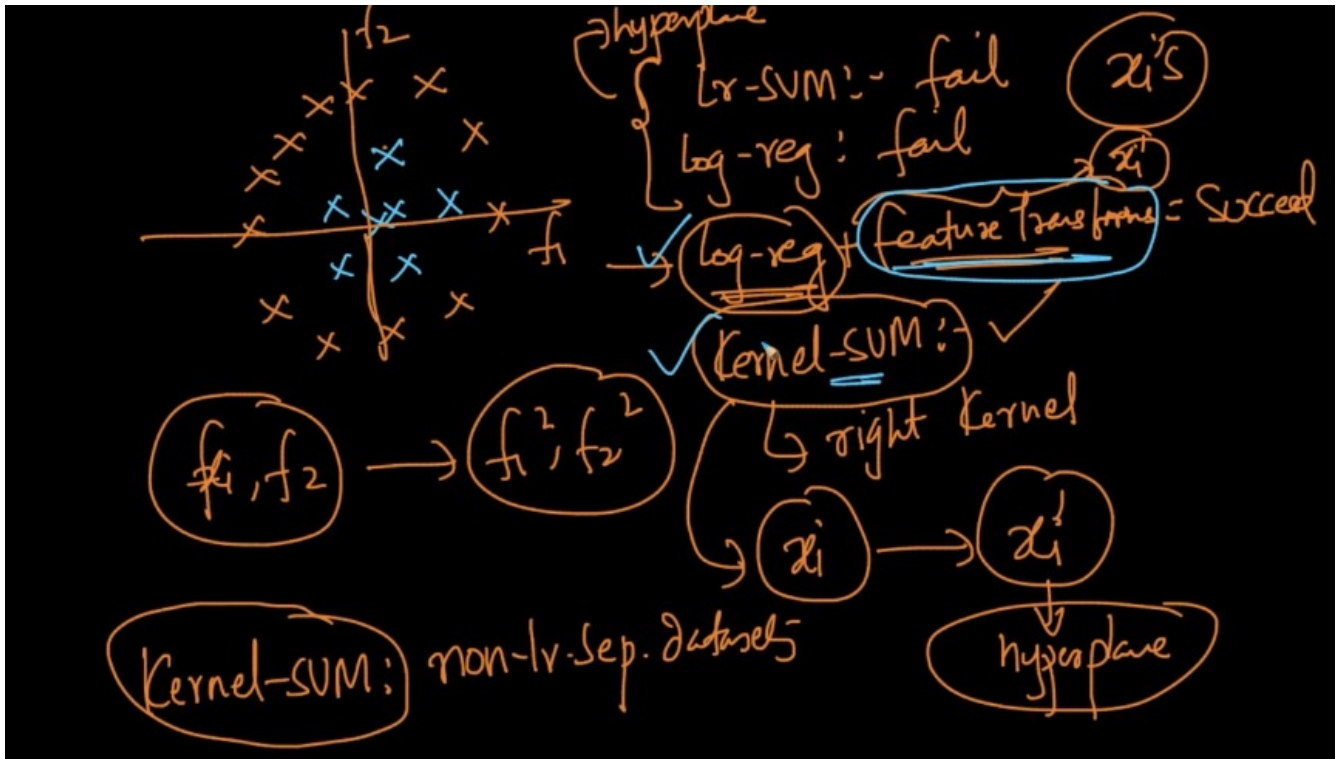
Soft-SVM-hyperplanes \approx log-reg
 \hookrightarrow margin-max.

lr-SVM:- $x_i^T x_j$

kernel-SVM:- $K(x_i, x_j)$

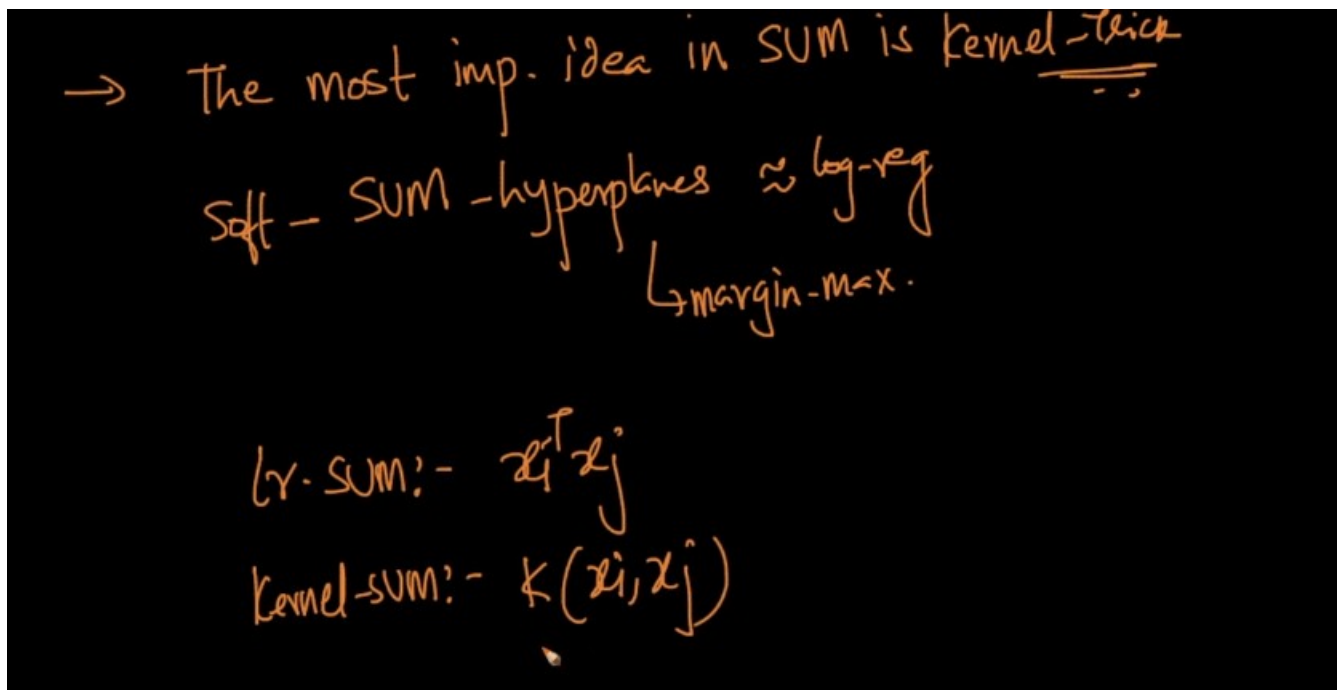
$$\checkmark K(x_i, x_j) = x_i^T x_j$$

Failure cases of SVM:

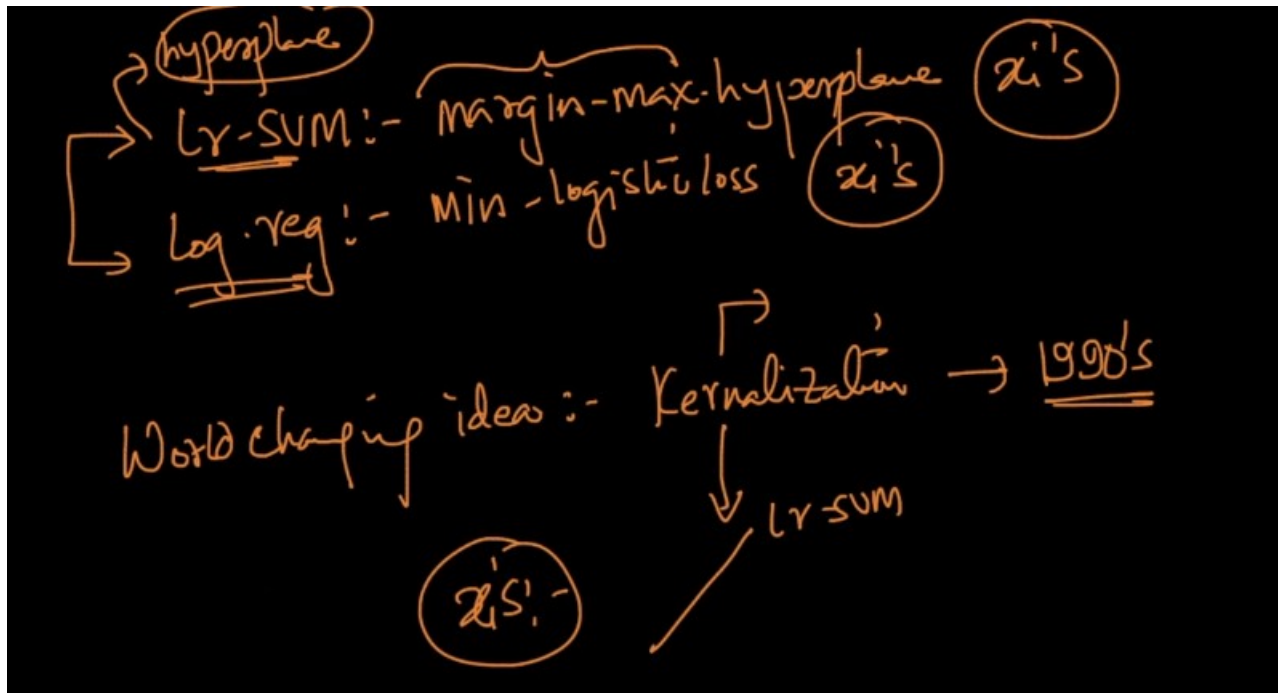


Kernalization makes SVM handle non – linearly separable data sets.

We use kernels in the dual form of SVM to make different relationships other than the default linear assumption.

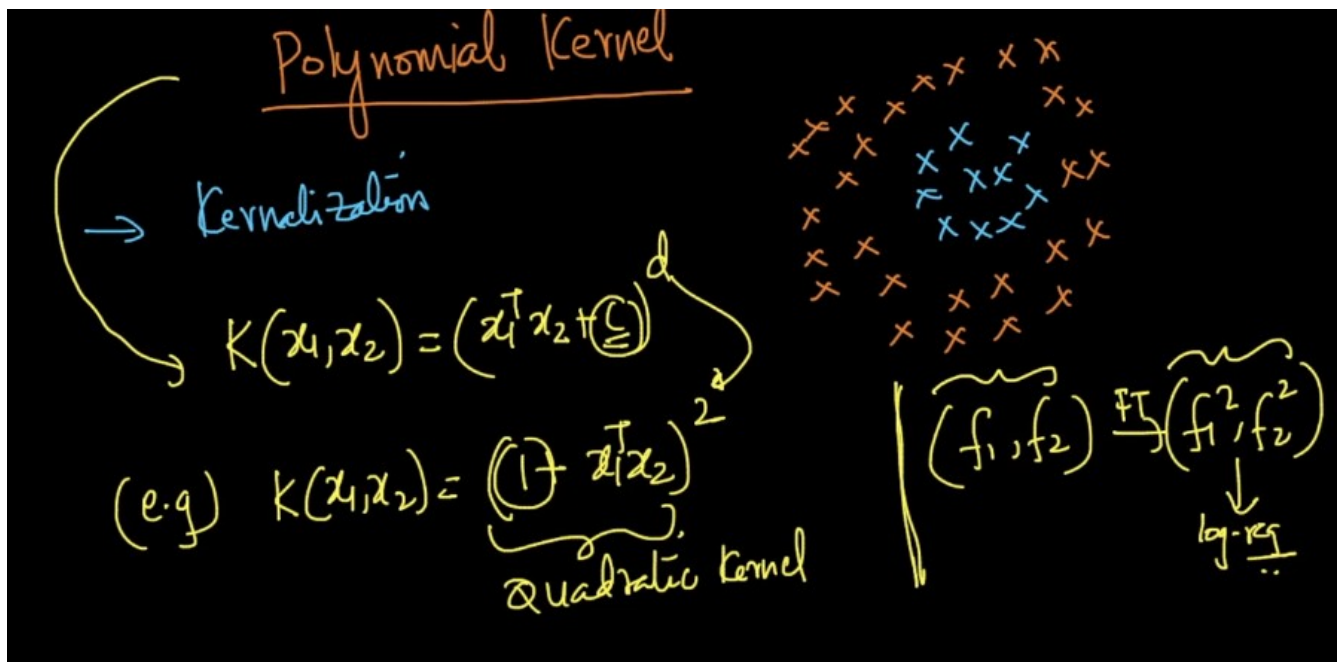


The key idea of SVM's is Kernel trick.



Non – linearly seperable data sets can be classified using SVM and kernels.

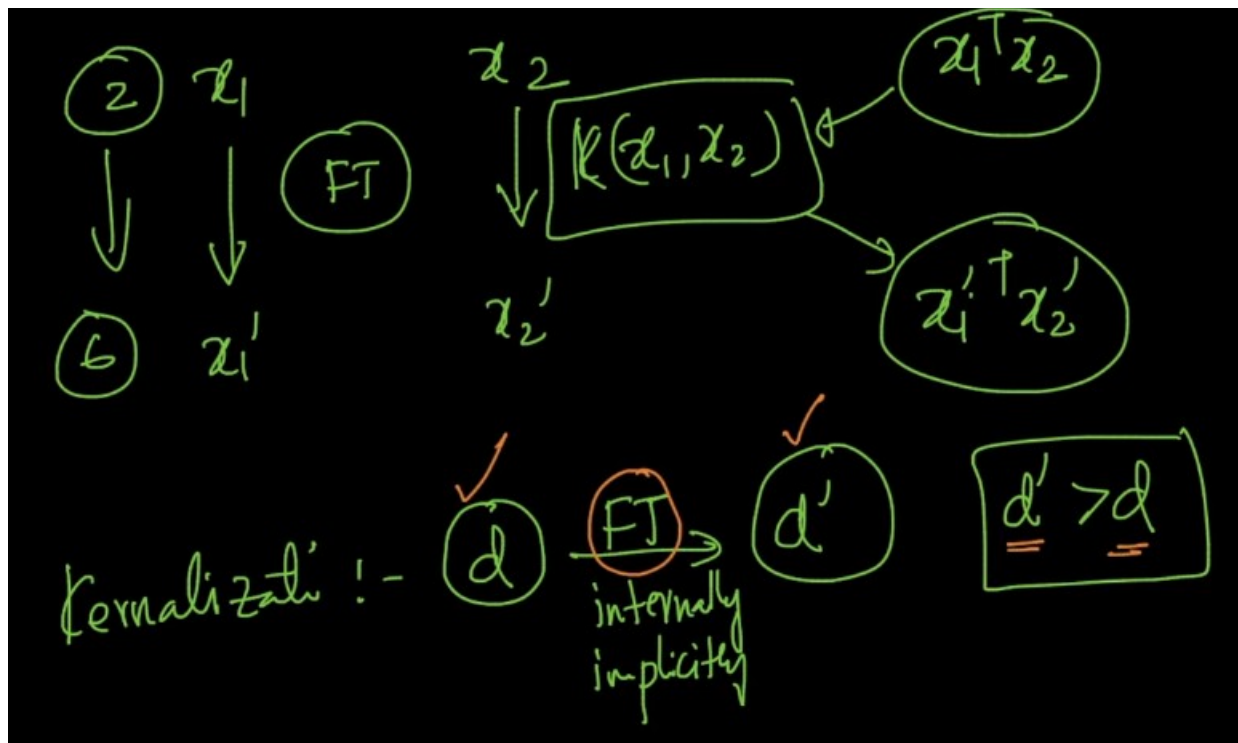
Polynomial kernel:



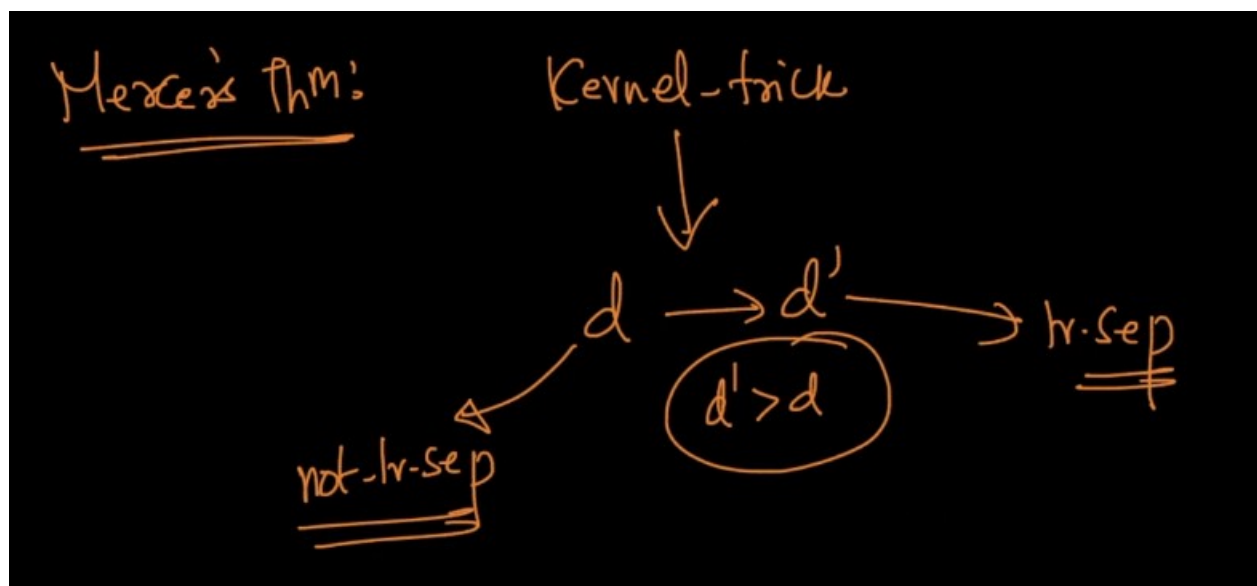
Example:

$$\begin{aligned}
 k(x_1, x_2) &= (1 + x_1^T x_2)^2 & x_1 &= \langle x_{11}, x_{12} \rangle \\
 &= (1 + x_{11}x_{21} + x_{12}x_{22})^2 & x_2 &= \langle x_{21}, x_{22} \rangle \\
 &= 1 + \underline{x_{11}^2} \underline{x_{21}^2} + \underline{x_{12}^2} \underline{x_{22}^2} + 2x_{11}x_{21} + 2x_{12}x_{22} + 2x_{11}x_{21}x_{12}x_{22} \\
 \text{Let } [1, x_{11}^2, x_{12}^2, \sqrt{2}x_{11}, \sqrt{2}x_{12}, \sqrt{2}x_{11}x_{12}] &: x'_1 \\
 [1, x_{21}^2, x_{22}^2, \sqrt{2}x_{21}, \sqrt{2}x_{22}, \sqrt{2}x_{21}x_{22}] &: x'_2 \\
 &= (x'_1)^T (x'_2)
 \end{aligned}$$

Kernelization to the feature transform internally using the kernel trick.



The data will become separable after kernel trick. After applying mercers theorem the data will be more linearly separable transforming to d' dimensional space.

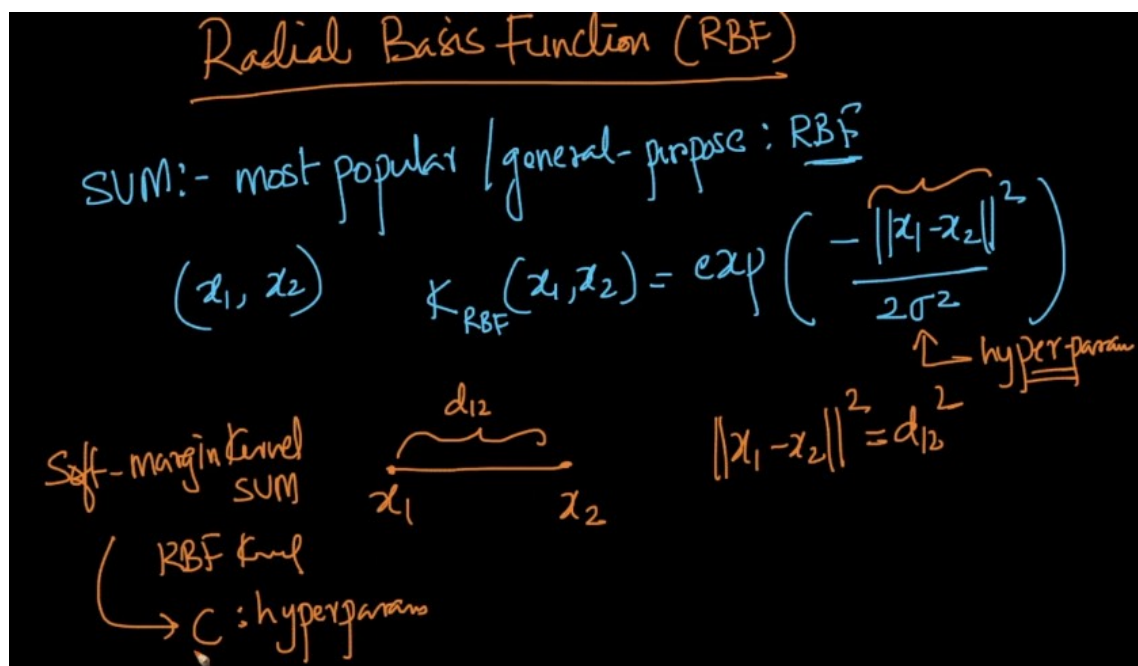


Kernel trick does the transformation internally than that of the feature transform that is done externally in logistic regression.

What is a right Kernel to apply?

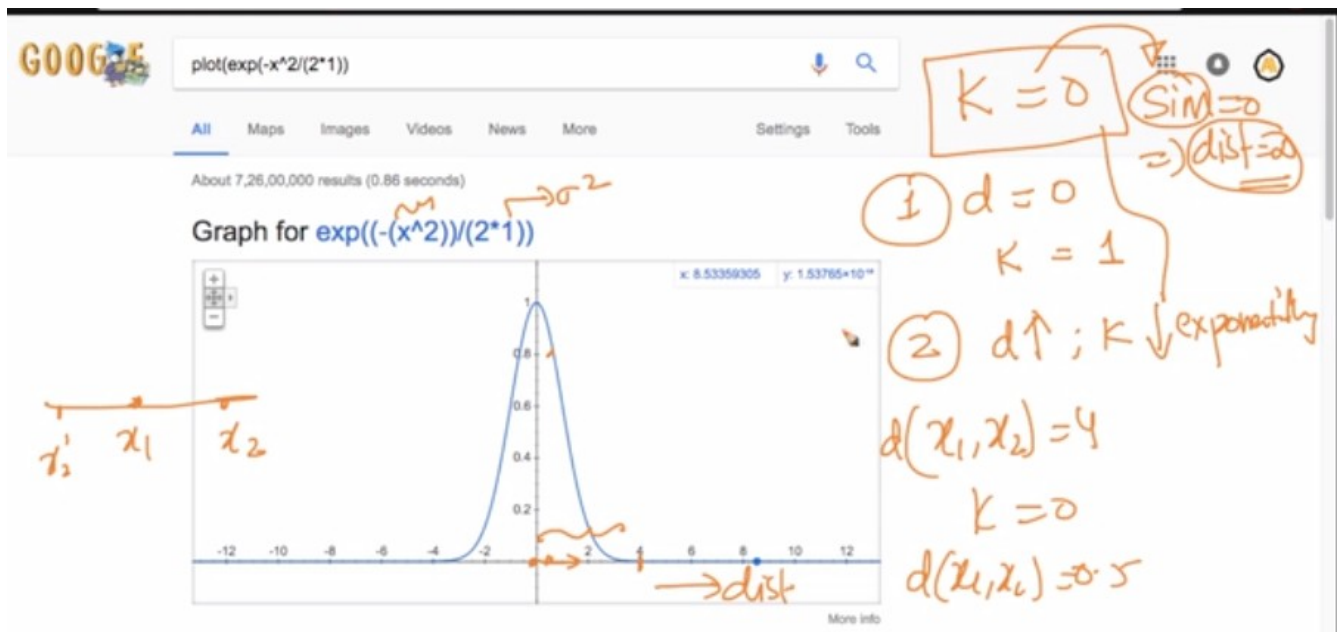
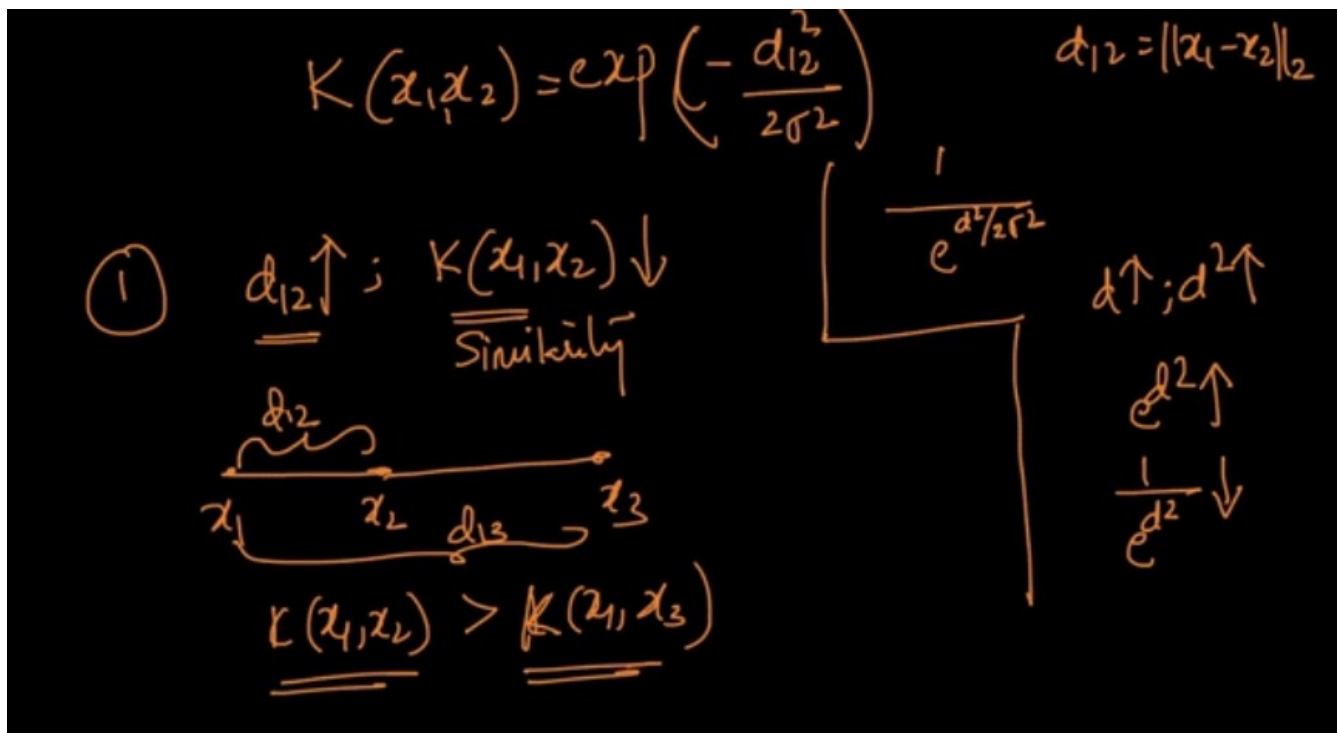
The problem of explicit feature transform is changed to find the right kernel in SVM than that of Logistic regression.

RBF – kernel:



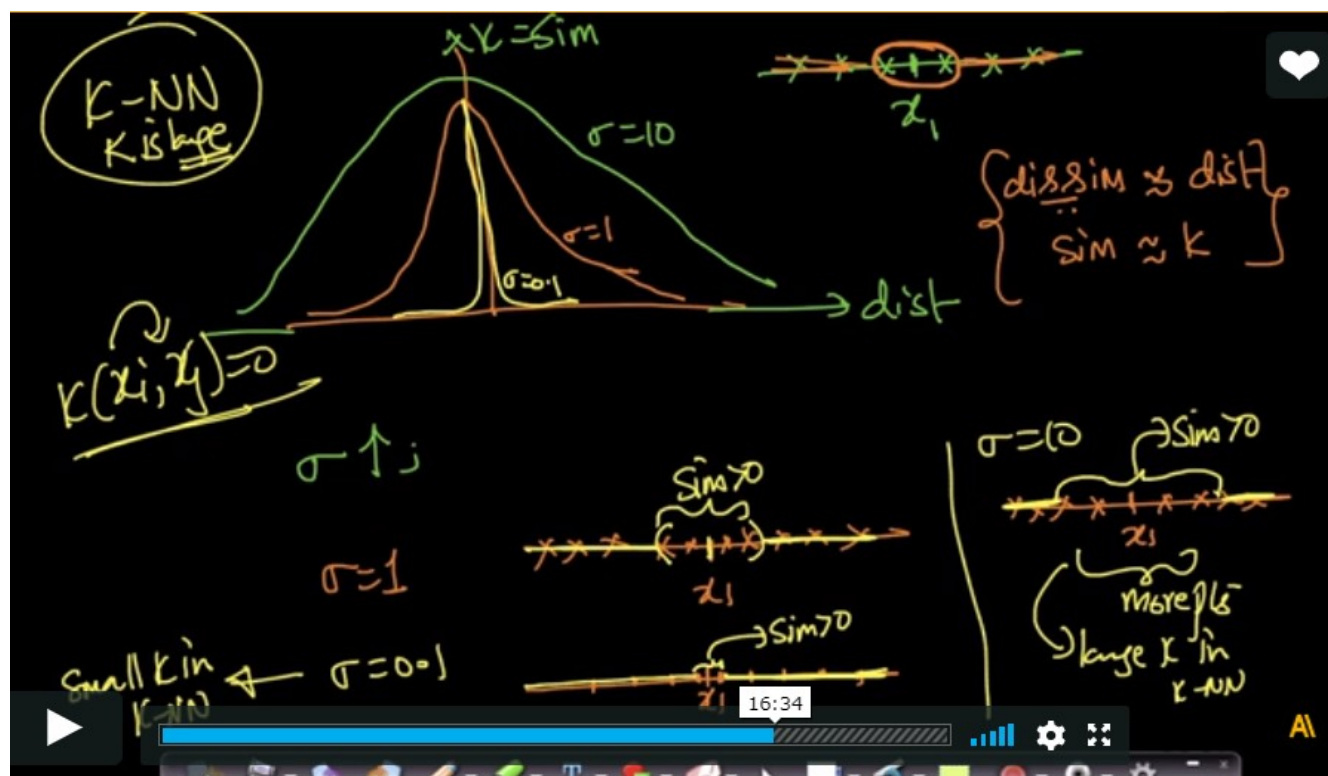
Here sigma is the hyper parameter in RBF. As distance between the points increases then the kernel decreases.

Kernels work like similarity distance measures.

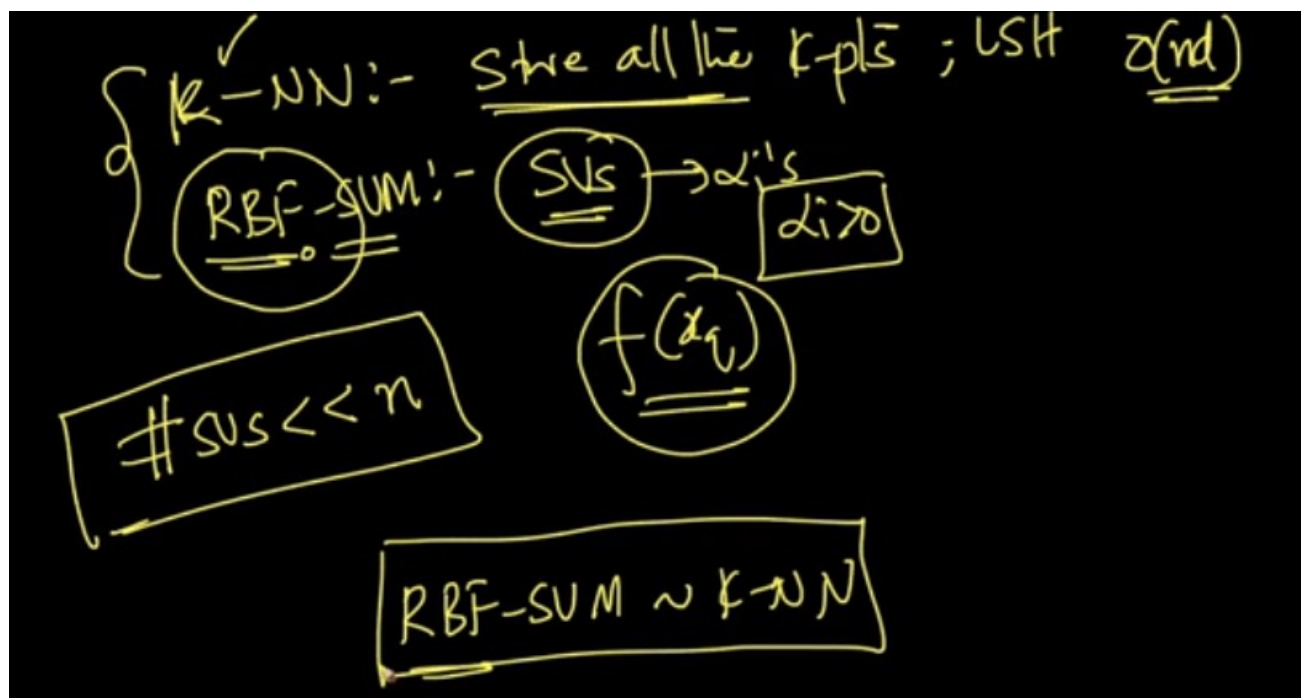


Kernels measure the similarity.

Distance measure the distance between the data points.

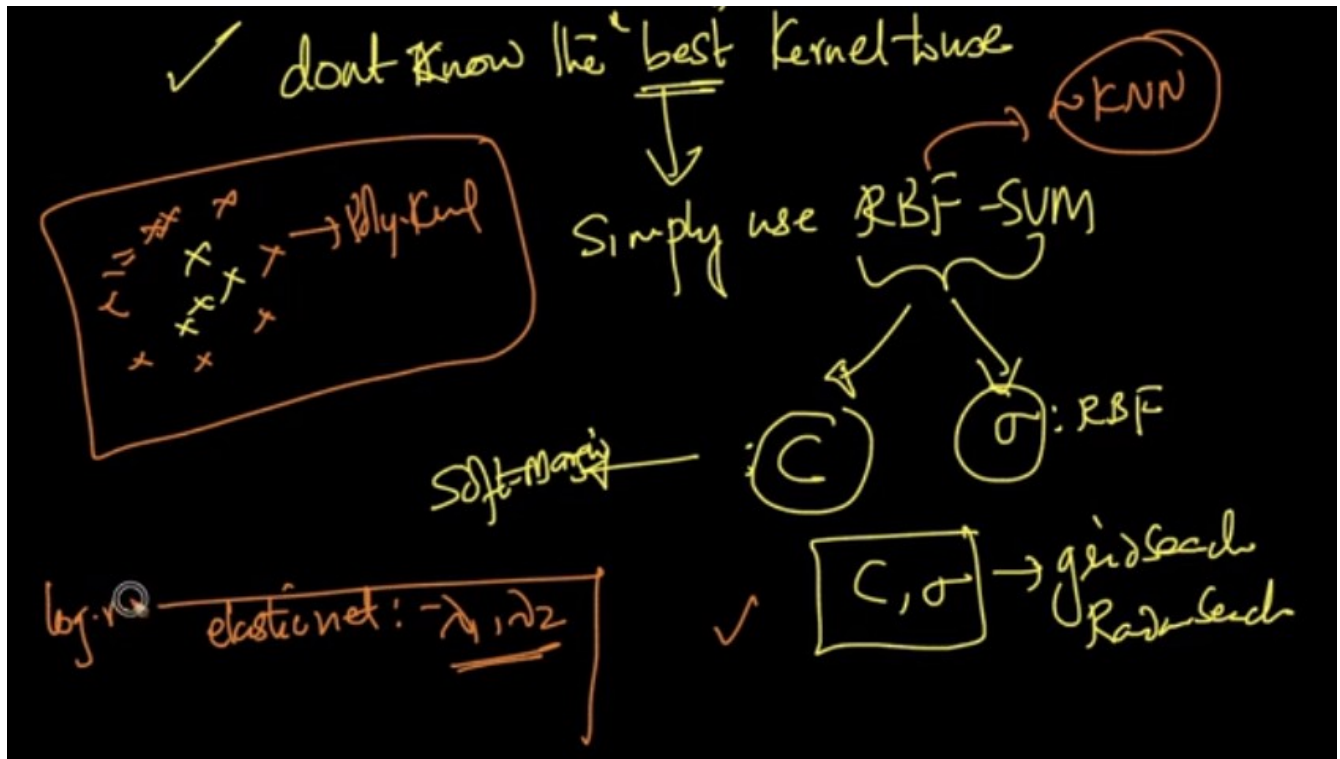


KNN and RBF- Kernel



If there is no clue of choosing the kernel, RBF will be the good choice.

Elastic nets have two hyper parameters



Domain specific Kernels:

There are a lot of specify kernels,

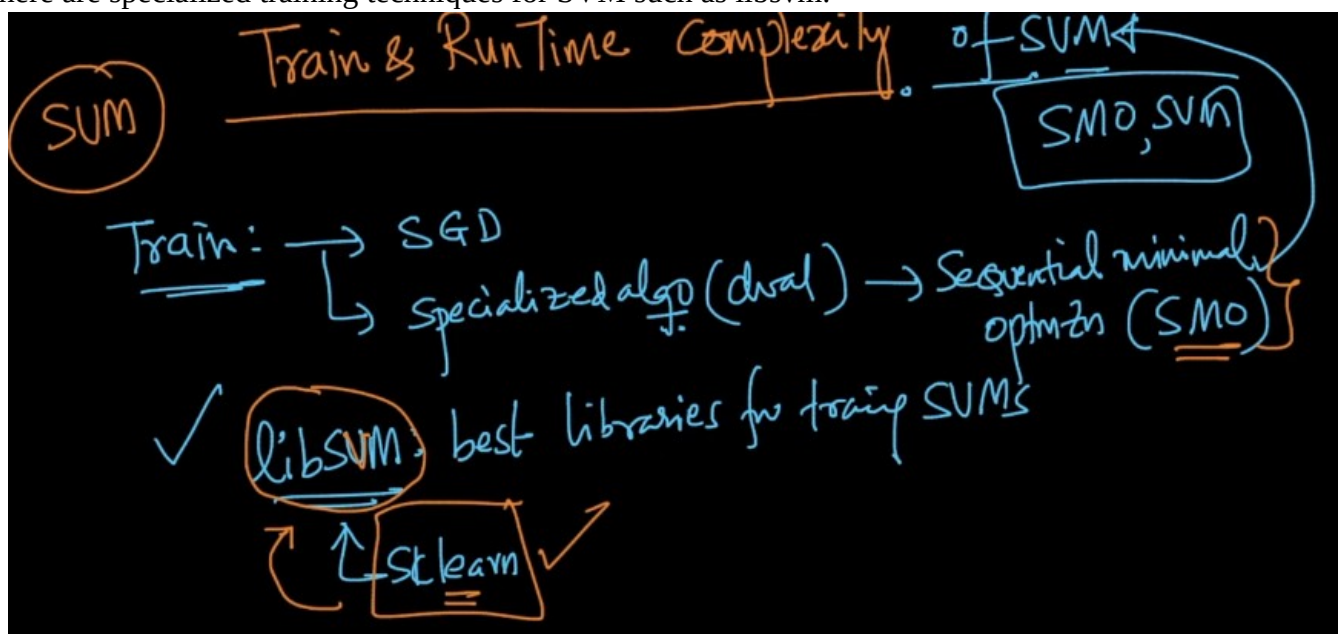
String kernels – If the features contains text.

Geonmic kernels

Graph kernels

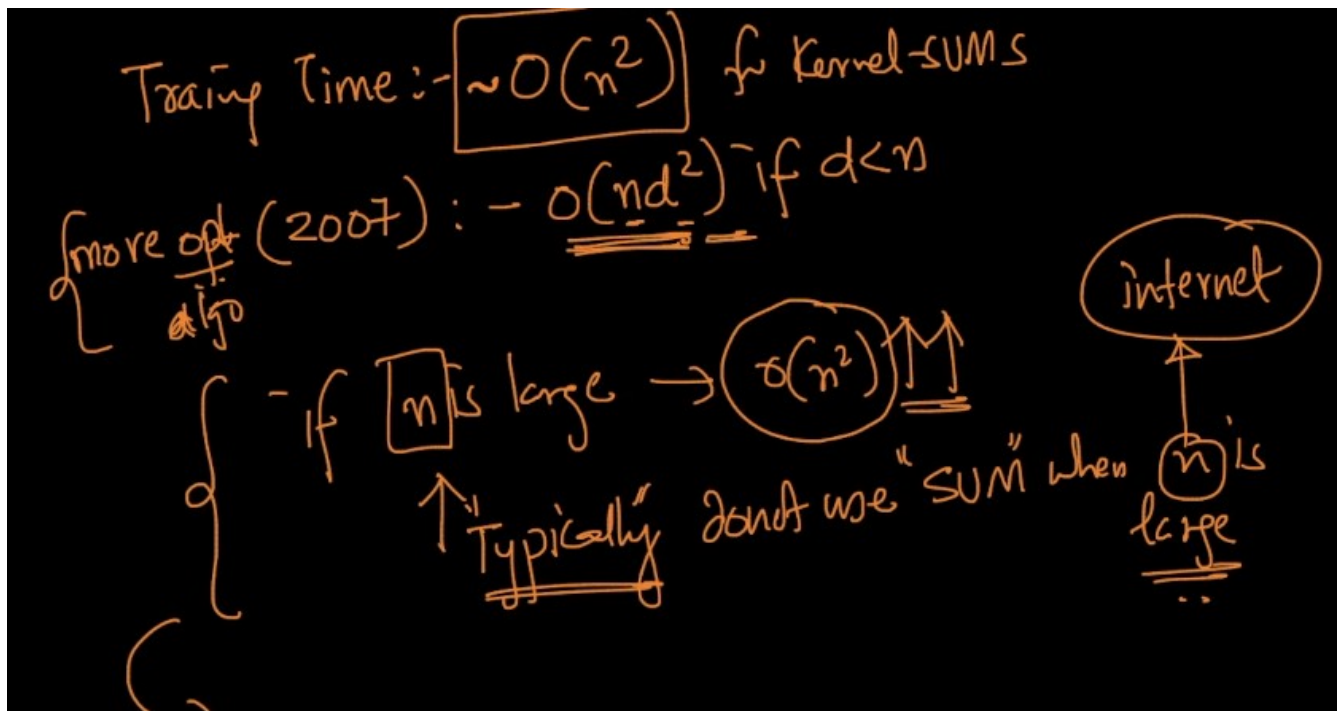
Train and run time complexities of SVM.

There are specialized training techniques for SVM such as libsvm.

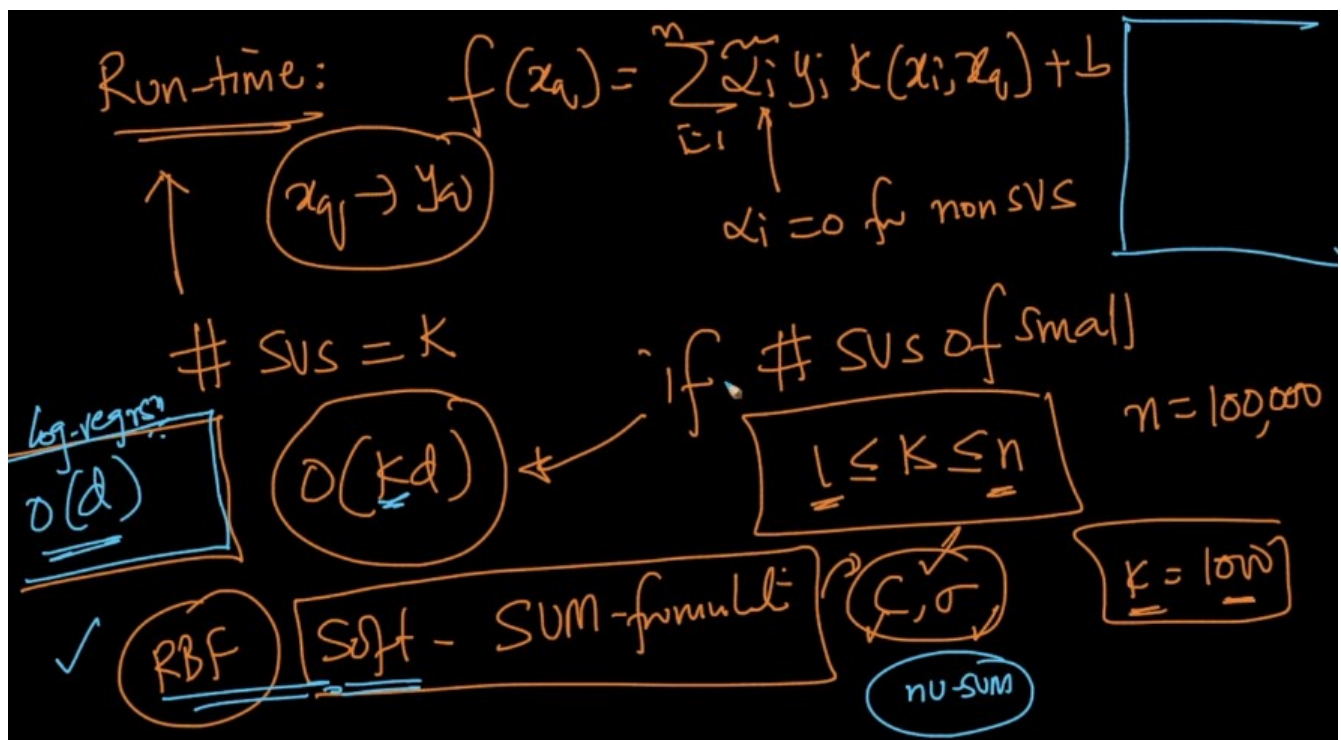


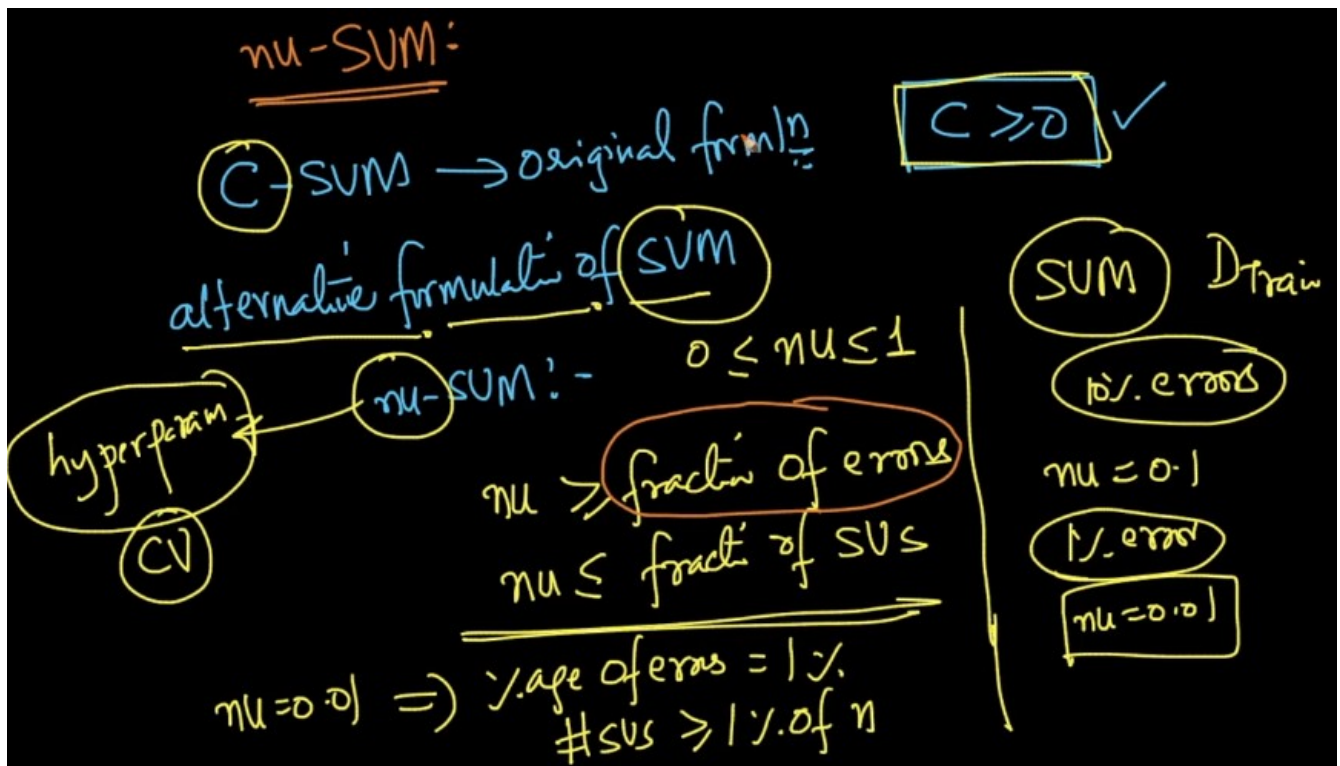
The training time of the SVM is roughly around $O(n^2)$.

Typically, we don't use SVM when number of data points are large.

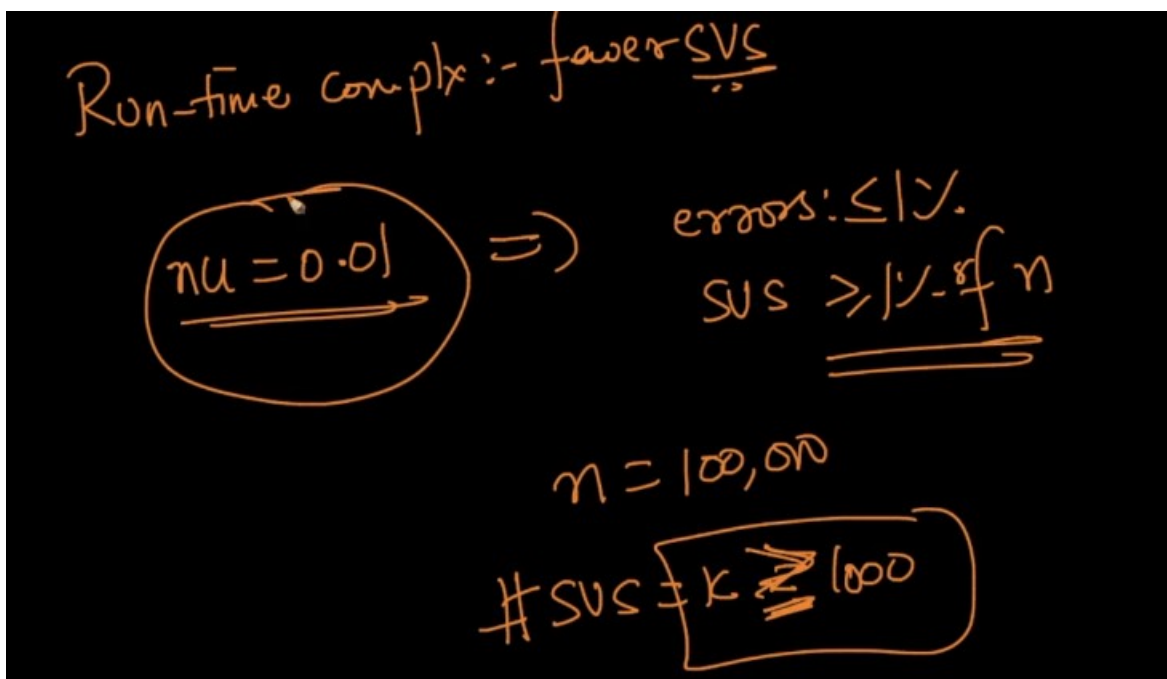


Run time complexity:





nu is the value of the error made by the SVM.



SVM regression:

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

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

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

hyperparam

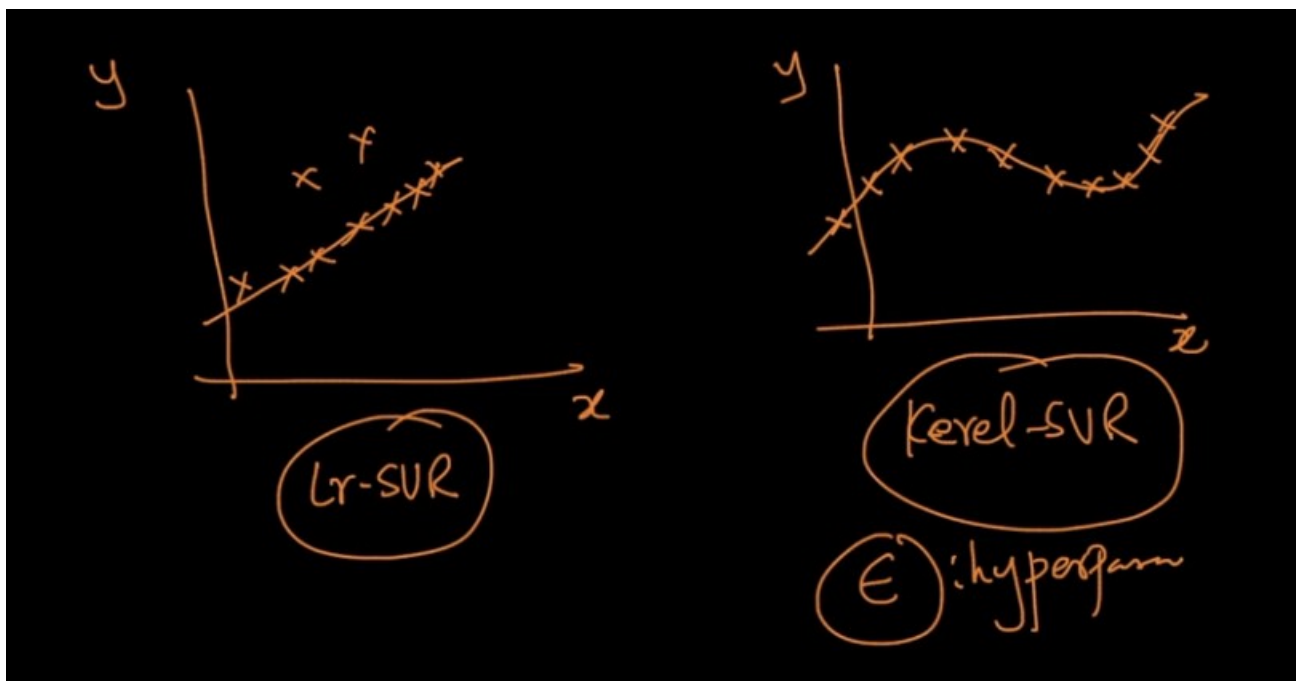
$f(x_i) = \hat{w}^T x_i + b = \hat{y}_i$

$\left\{ \begin{array}{l} y_i - \hat{y}_i \leq \epsilon \\ \hat{y}_i - y_i \leq \epsilon \end{array} \right.$

Cr. form of SVR

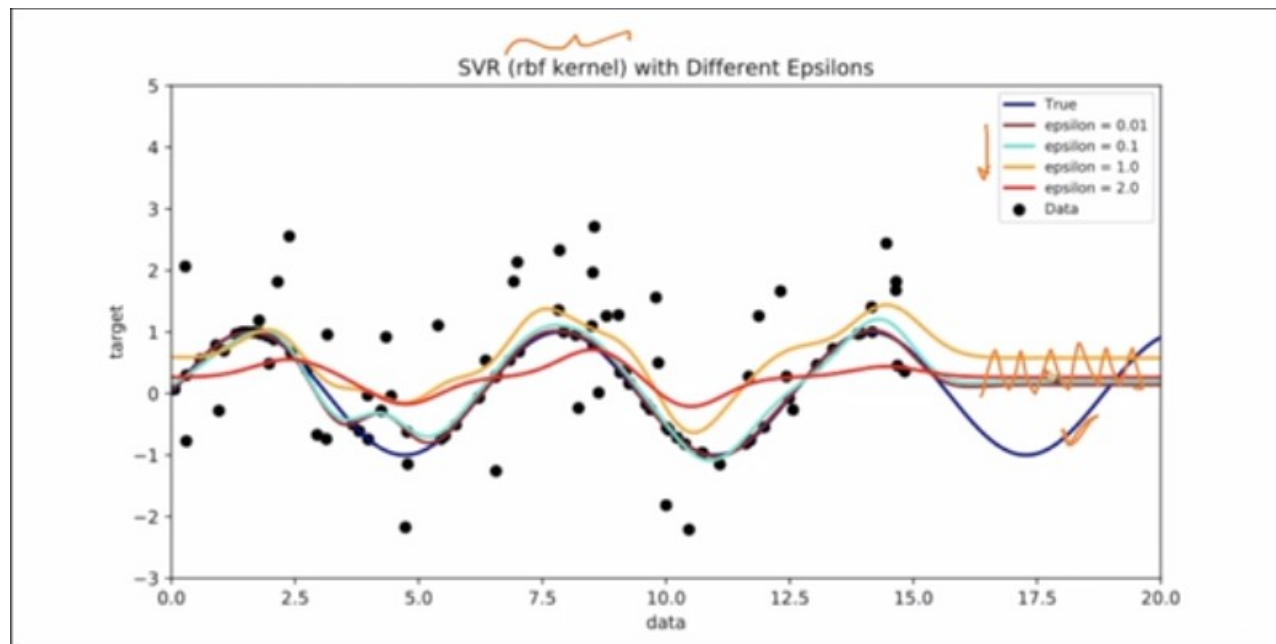
Here, we obtain the difference between the y_i and \hat{y}_i and made to less than or equal to epsilon, where epsilon is a hyper parameter.

Where the $\|w\|^2$ can be kernalized.



If epsilon is low then the error are low on the training data, over fitting will increase.
If epsilon is high then errors on training data will increase, that means under fit the data.

RBS represents the KNN-reg roughly.



As epsilon is very less then the likelihood of making mistakes decreases, then the model will over fit the data, choosing the right value of epsilon helps in making the model to generalize well.

Cases:

feature engineering: Kernel design, finding the right kernel.

Decision surface: Lr sum – Hyperplane.

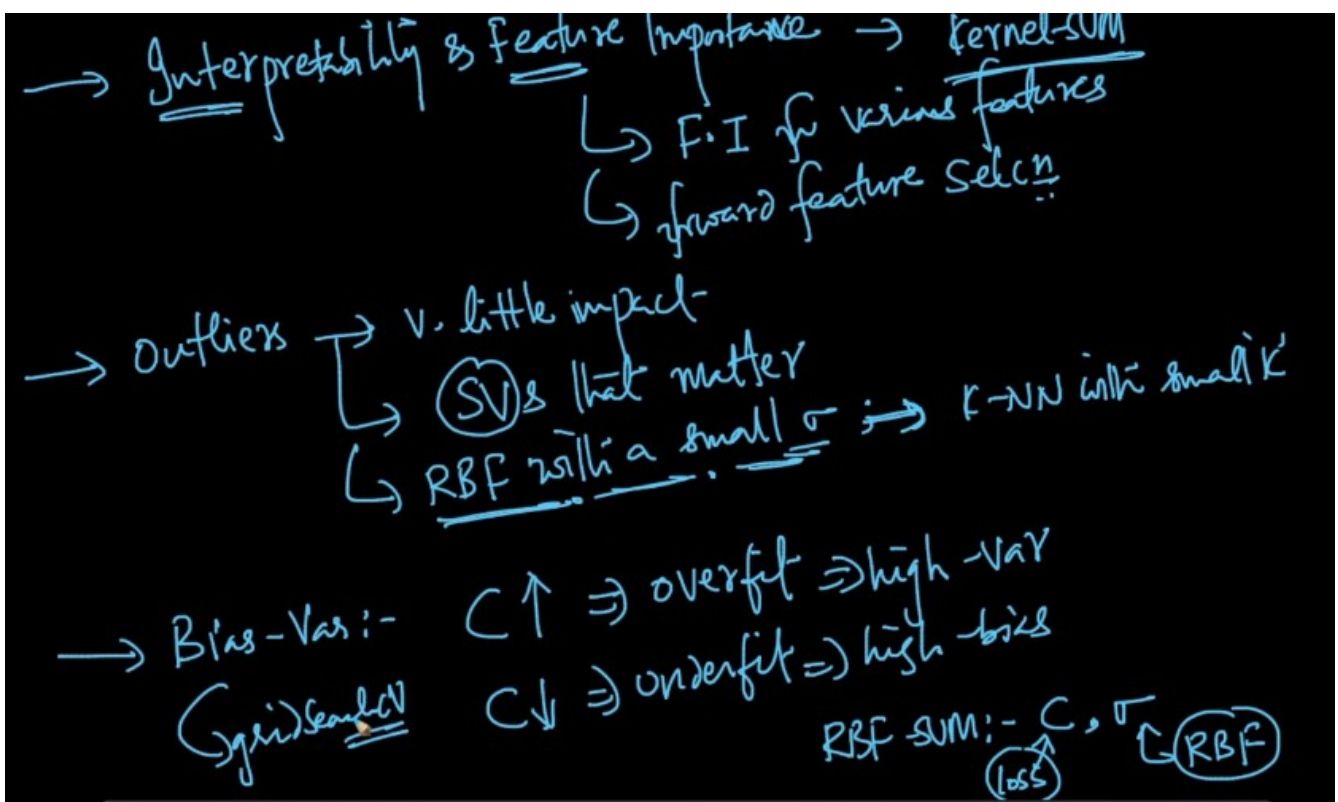
Kernel - sum

Cases in SVM:

SVM's can work with similarity matrix.

Interpretability: For kernel SVM's we cannot get the feature importance.

Outliers: Outliers have little impact because we only consider the support vectors, RBF kernels get influenced by outliers, because it behaves as KNN.



The best case for SVM is when we have right kernel.

