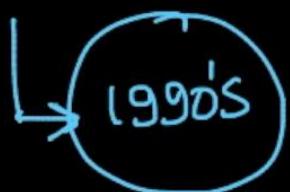


Support Vector Machines (SVM)

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

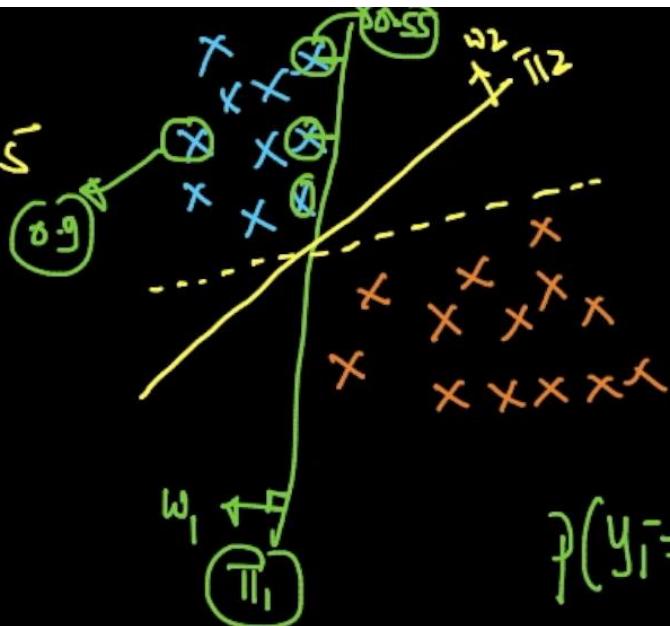


geom-intuition

→ Many TIs that separate +ve from -ve pls

→ Key idea of SVM:-

TI that separates +ve from -ve pls as 'widely' as possible

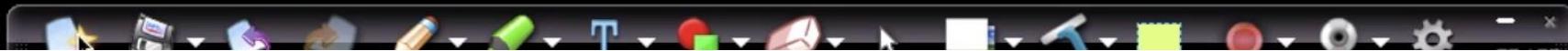
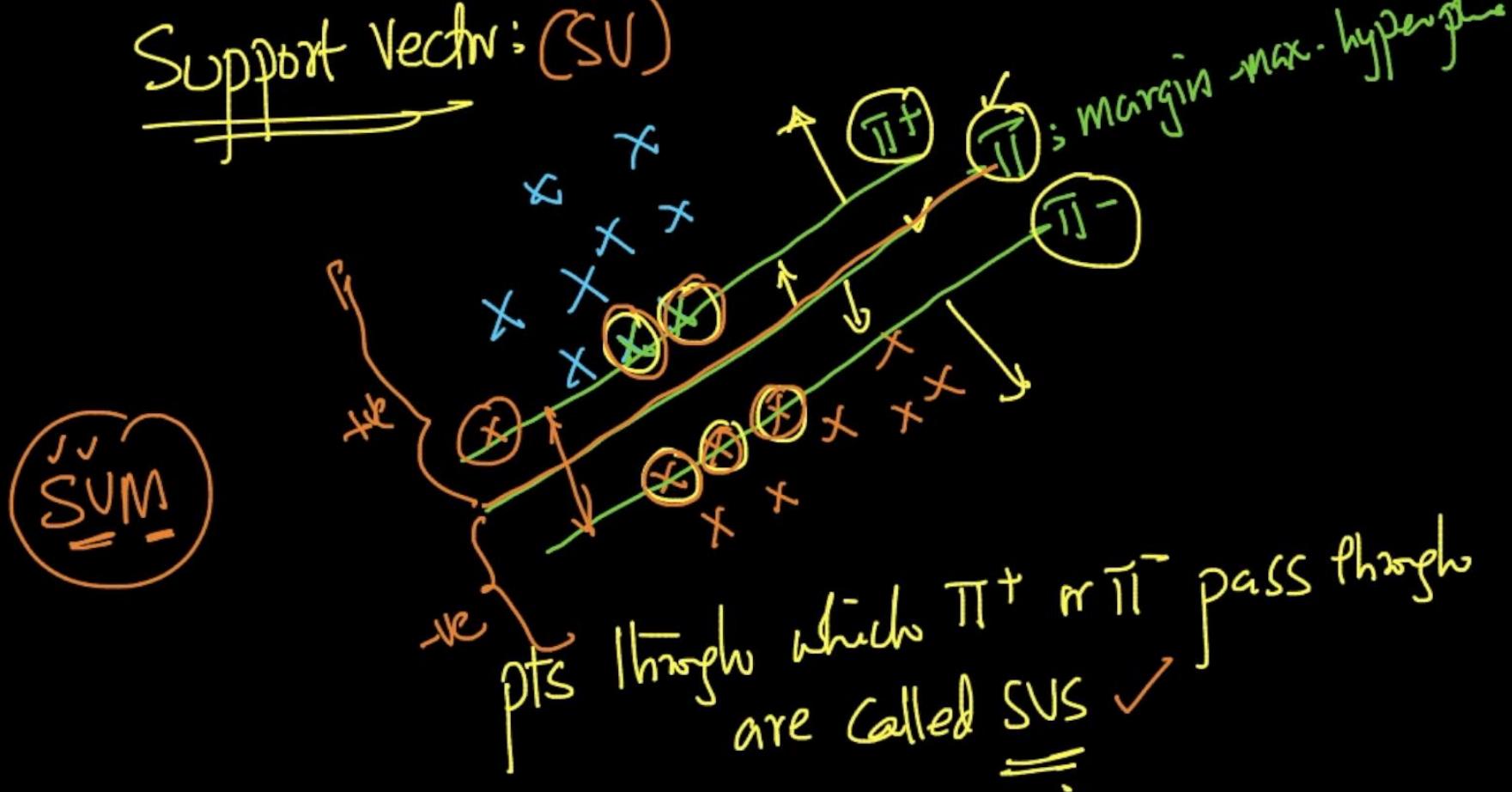


$$\hat{y}(y_i=1) = \sigma(\omega^T x_i)$$

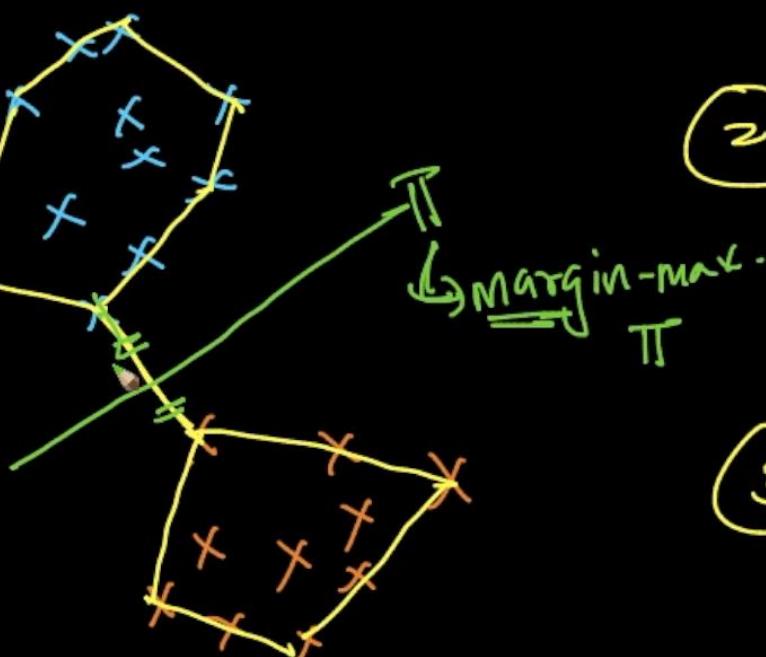
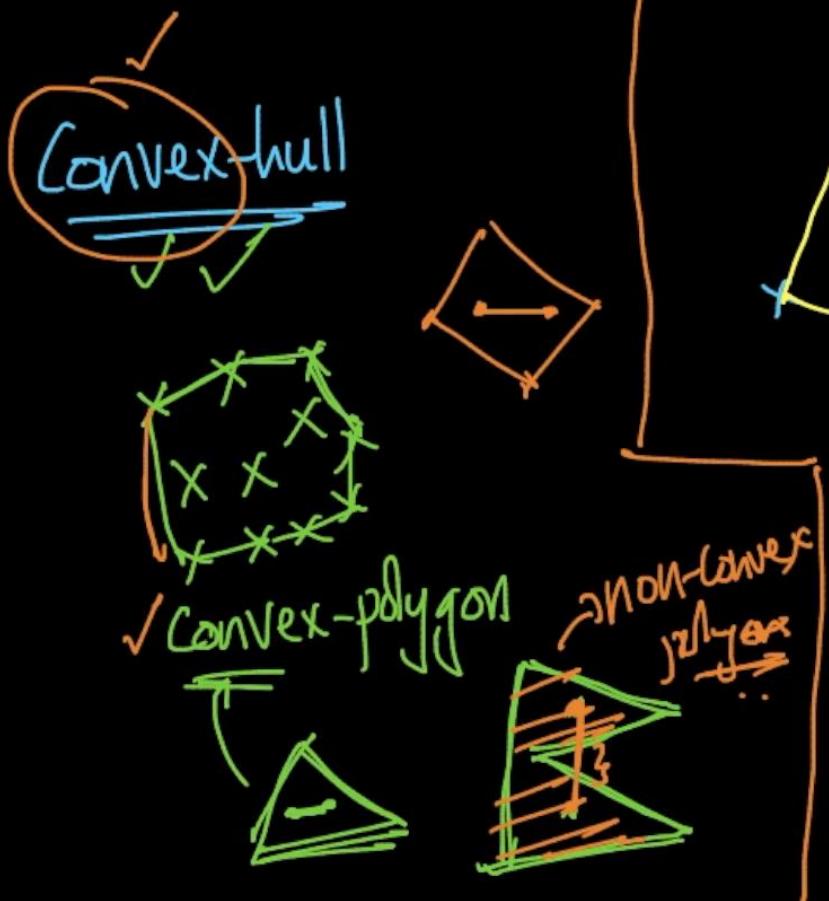
SUM :- Try to find a Π that maximizes the
margin = $\text{dist}(\Pi^+, \Pi^-)$

✓ Margin ↑ \Rightarrow generalization acc ↑

Support Vectn: (SV)



alternative geom intr of sum: π



- ① convex-hull
for +ve pls
&-ve pls
- ② find the shortest
line connecting
these hulls.

- ③ bisect the
line

Mathematical formulation of SVM:

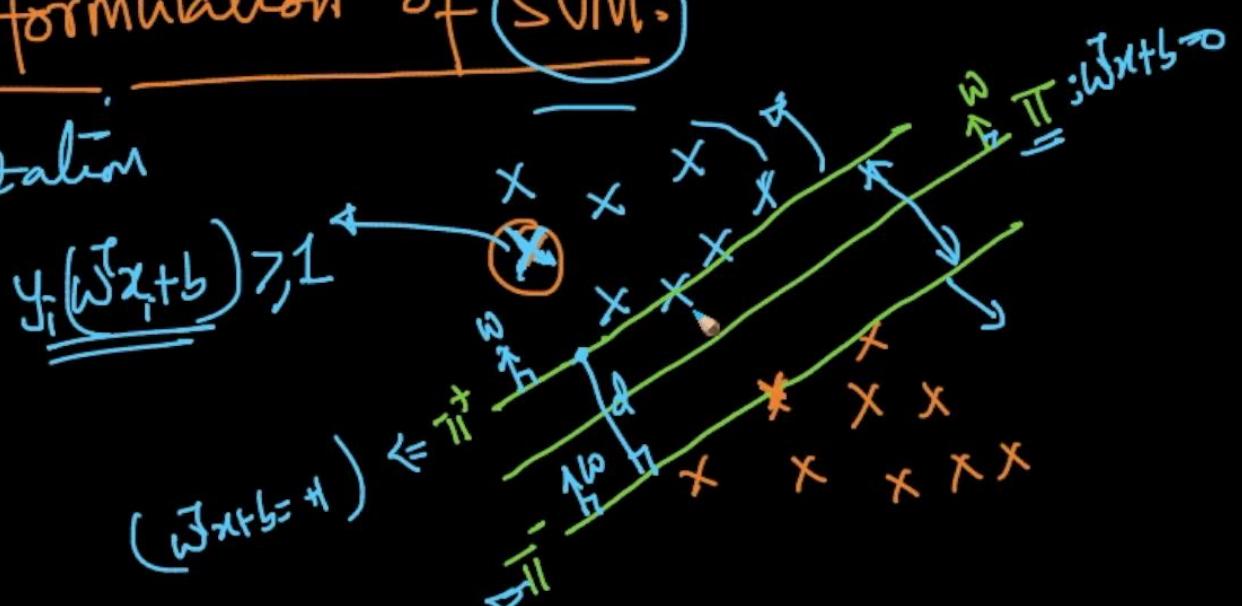
Π^- : margin maximization

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

if $\Pi^+: \underline{\omega}^T x + b = 1$

$$\Pi^-: \underline{\omega}^T x + b = -1$$

Margin: $\underline{d} = \frac{2}{\|\omega\|}$

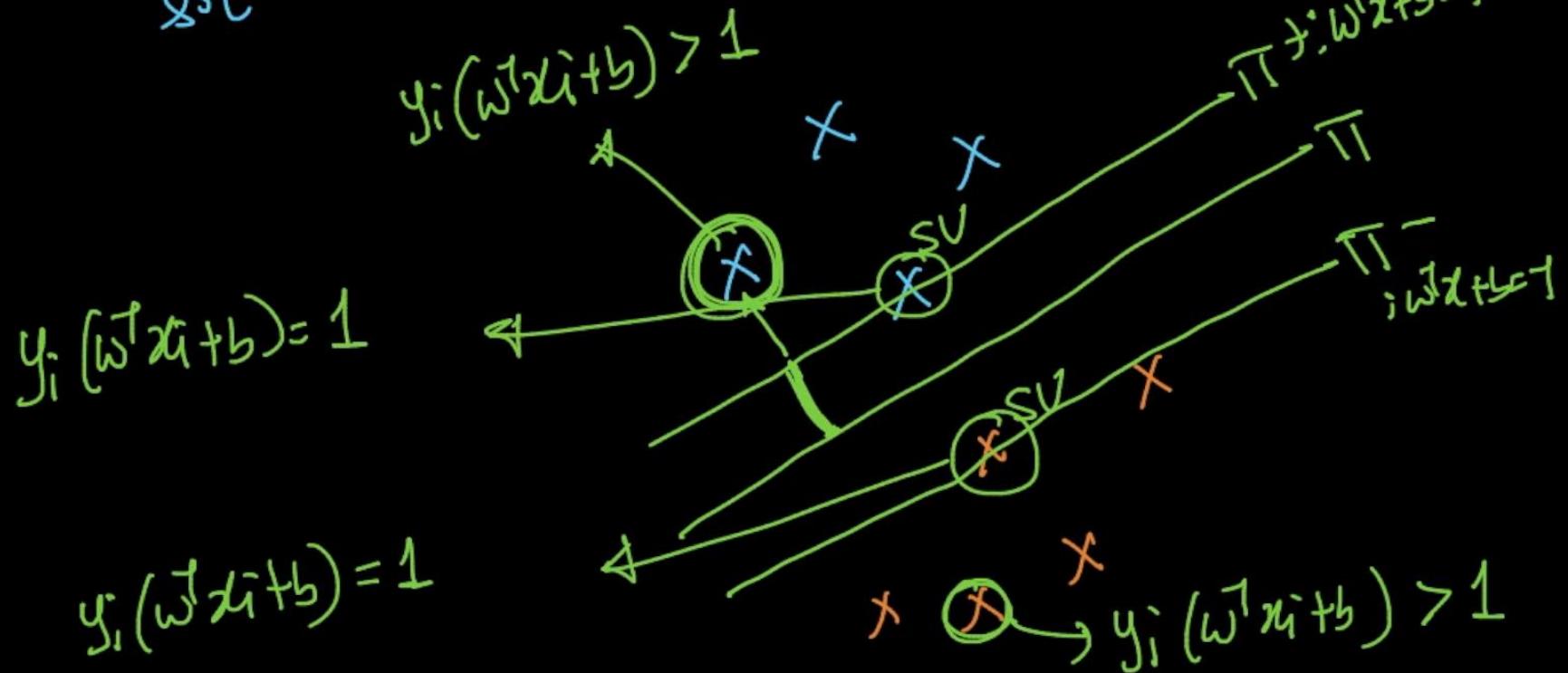


$$(\underline{\omega}^T x + b = 1) \quad (\underline{\omega}, b) = \underset{\omega, b}{\operatorname{argmax}} \quad \frac{2}{\|\omega\|}$$

Set (const)

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

so that



$$(\vec{w}^*, b^*) = \underset{\vec{w}, b}{\operatorname{argmax}} \quad \frac{2}{\|\vec{w}\|} \leftarrow \text{Margin}$$

Set $y_i (\vec{w}^T \vec{x}_i + b) \geq 1$ for all \vec{x}_i

n-const



Constn,
Optimzn.
prob,
of SVM

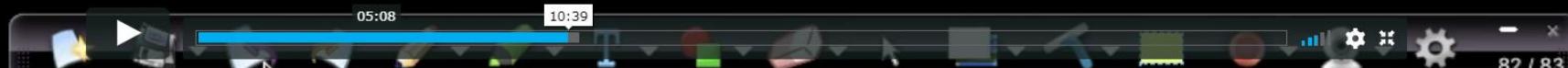
$$\left\{ \begin{array}{l} \vec{w}, b^* = \underset{\vec{w}, b}{\operatorname{argmax}} \frac{2}{\|\vec{w}\|} \\ \text{s.t } \forall i, y_i(\vec{w}^T \vec{x}_i + b) \geq 1 \end{array} \right.$$

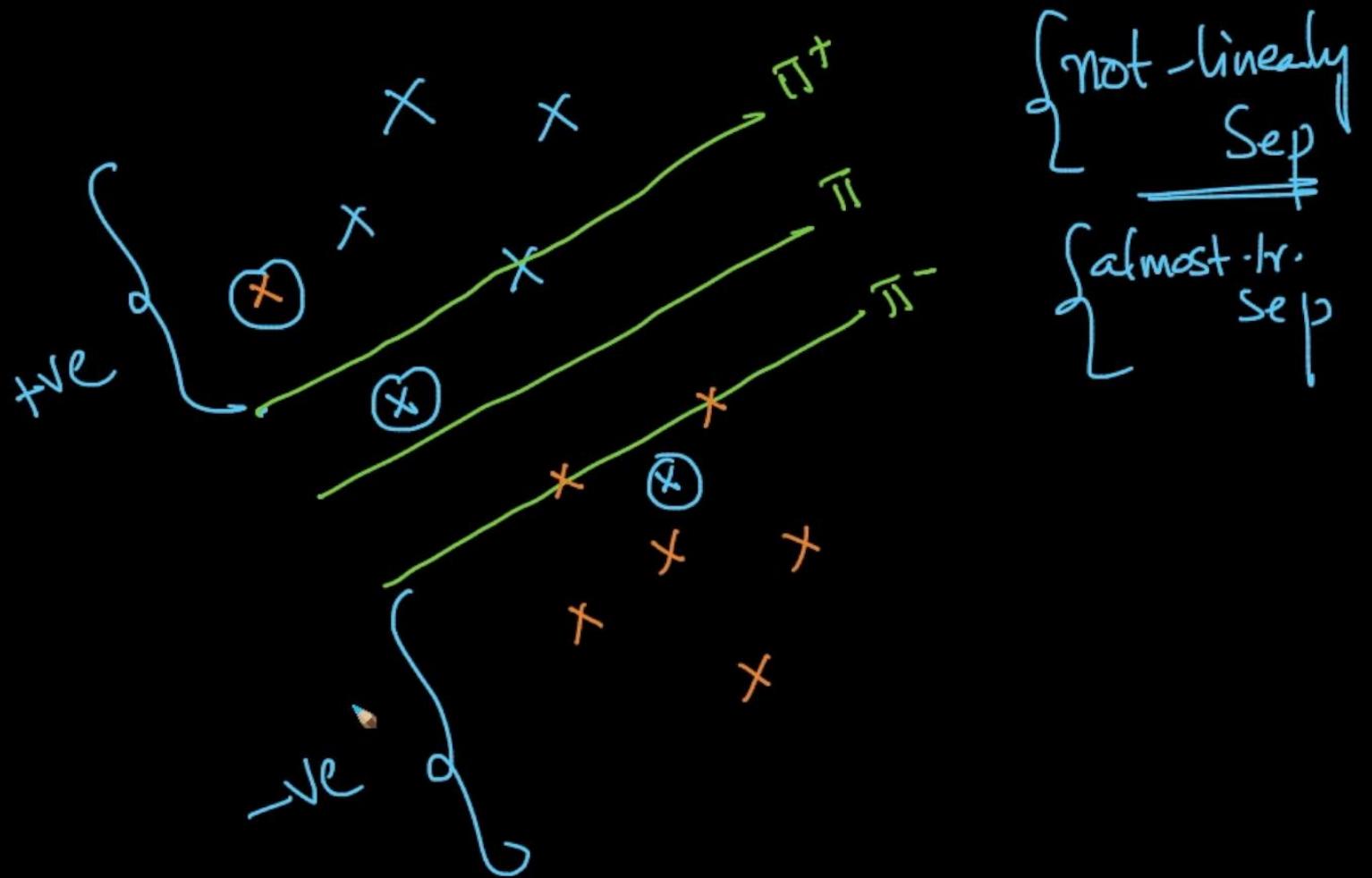


Constn,
Optimzn.
prob,
of SVM

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

data is linearly sep





$$\Pi: \vec{w}^T x + b = 0$$

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

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

$\vec{w} \neq \text{Unit-Norm}$

$$y_i(\vec{w}^T x_i + b) = 0.5$$

$$y_i(\vec{w}^T x_i + b) = -0.5$$

$$= y_i(\vec{w}^T x_i + b) = 1 - (1.5)$$



{not-linearly
Sep}

{almost-lin.
Sep}

$$y_i(\vec{w}^T x_i + b) = +0.5$$

$$y_i(\vec{w}^T x_i + b) = 1 - (2.5)$$

$\xi_i \uparrow$; pt is further away from the correct Π in margin



21:04



23:38
09-09-2019

$$(\vec{w}, b) = \arg \max_{\vec{w}, b} \left(\frac{2}{||\vec{w}||} \right)$$

$$\sum_i \xi_i$$

$$= \arg \min_{\vec{w}, b} \left(\frac{||\vec{w}||}{2} \right)$$

$$\max_x f(x) = \min_x -f(x)$$



22:14

23:40
09-09-2019

$$(\hat{\omega}, \hat{b}) = \underset{\omega, b}{\operatorname{argmin}} \quad \text{margin} \cdot \frac{\|\omega\|}{2} + C \cdot \frac{1}{n} \sum_{i=1}^n \xi_i$$

avg. dist of miss. pts from hyperspace

s.t.

$$\left\{ \begin{array}{l} y_i (\hat{\omega}^T x_i + \hat{b}) \geq 1 - \xi_i \quad \forall i \\ \xi_i \geq 0 \end{array} \right\}$$

corr. classif. pts $\xi_i = 0$
incorrect. class. pts $\xi_i > 0$

minimize errors = min misclassif.

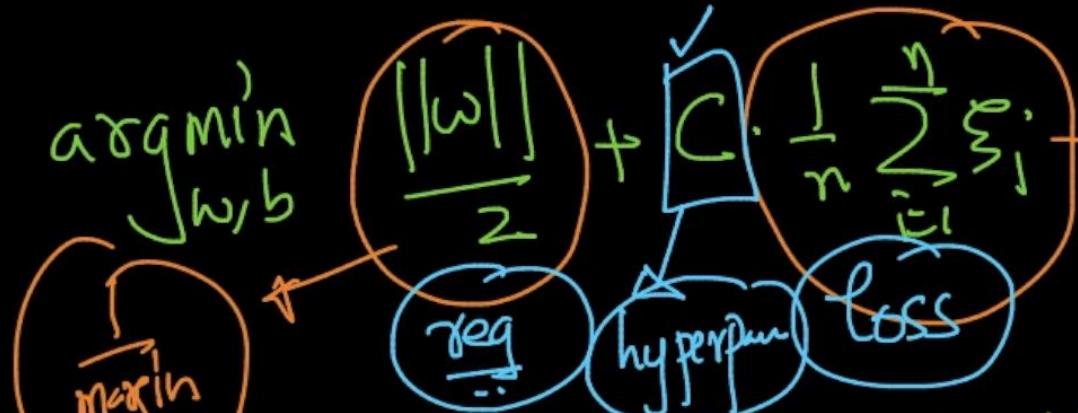
$$\min \|\sum \xi_i\|$$



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

$$\min_{w, b} C \cdot \text{margin}$$

margin



$$\text{s.t. } y_i (\hat{w}^T \hat{x}_i + \hat{b}) \geq 1 - \xi_i \quad \begin{cases} \xi_i > 0 \\ C: \text{+ve} \end{cases}$$

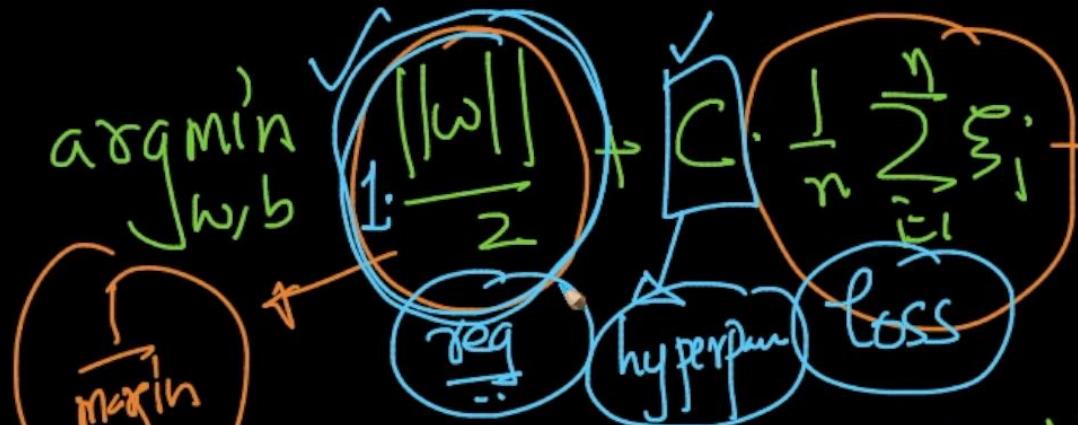
$$\min_w (\text{logistic-loss}) + \lambda (\text{reg})$$



$$(\hat{w}, \hat{b}) = \underset{w, b}{\operatorname{arg\,min}} \left[C \cdot \frac{1}{n} \sum_{i=1}^n \xi_i + \frac{1}{2} \|w\|^2 \right]$$

avg. dist for misscls

$\min_{w, b}$ margin loss yes



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

margin

$$\xi_i > 0$$

C: +ve

Soft-Margin SVM

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



$C \uparrow$

; tendency to make mistakes ↓
on D_{train}

\Rightarrow overfit \Rightarrow high-variance

$\lambda \uparrow \Rightarrow$ high bias

$\lambda \downarrow \Rightarrow$ high var

; underfit \Rightarrow high-bias

logistic-regsn -  in sklearn

λ

C

$C = \frac{1}{\lambda}$



23:49
09-09-2019 ENG

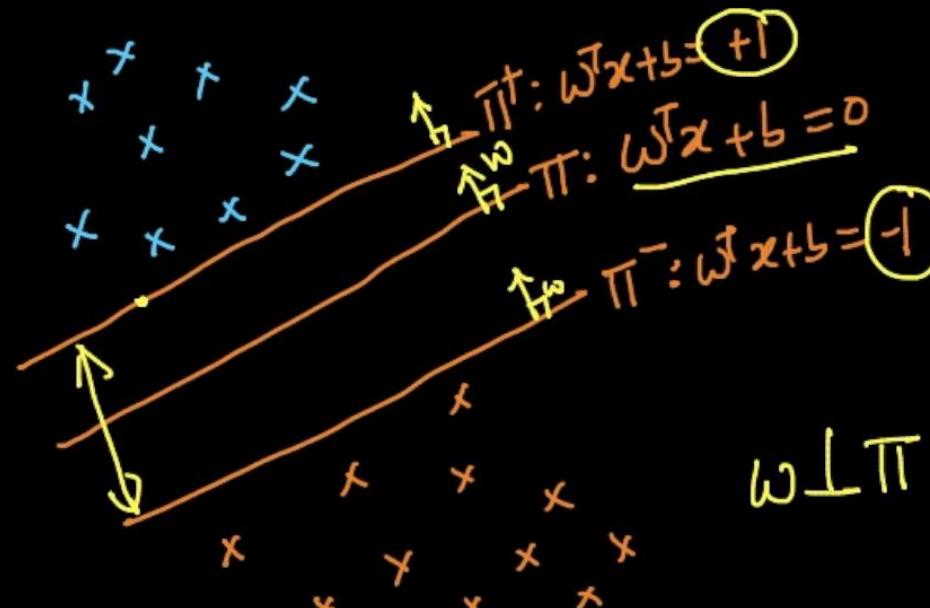
SVM:

(Q) Why $+1 \leq -1$ on the RHS
of $\pi^+ \leq \pi^-$

$$\text{Margin: } -\frac{2}{\|\omega\|}$$

$$\omega^T, b^* = \underset{\omega, b}{\arg \max} \frac{2}{\|\omega\|}$$

constraints



$\left\{ \begin{array}{l} \|\omega\| \neq 1 \\ \end{array} \right.$ (any vecr)
need not be
Unit-vecr



① $\pi^+ :- w^T x + b = K \checkmark$

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

$K > 0$

$K : \text{constant}$

$b = -4$

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

$$\underset{w,b}{\operatorname{argmax}} \frac{2}{\|w\|} = \underset{w,b}{\operatorname{argmax}} \frac{2K}{\|w\|} = \frac{8}{\|w\|}$$

② $\pi^{\text{f}, -} \omega^T x + b = k$

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

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

$$\omega \perp \pi$$

$\|\omega\|$ need not be 1



SVM:

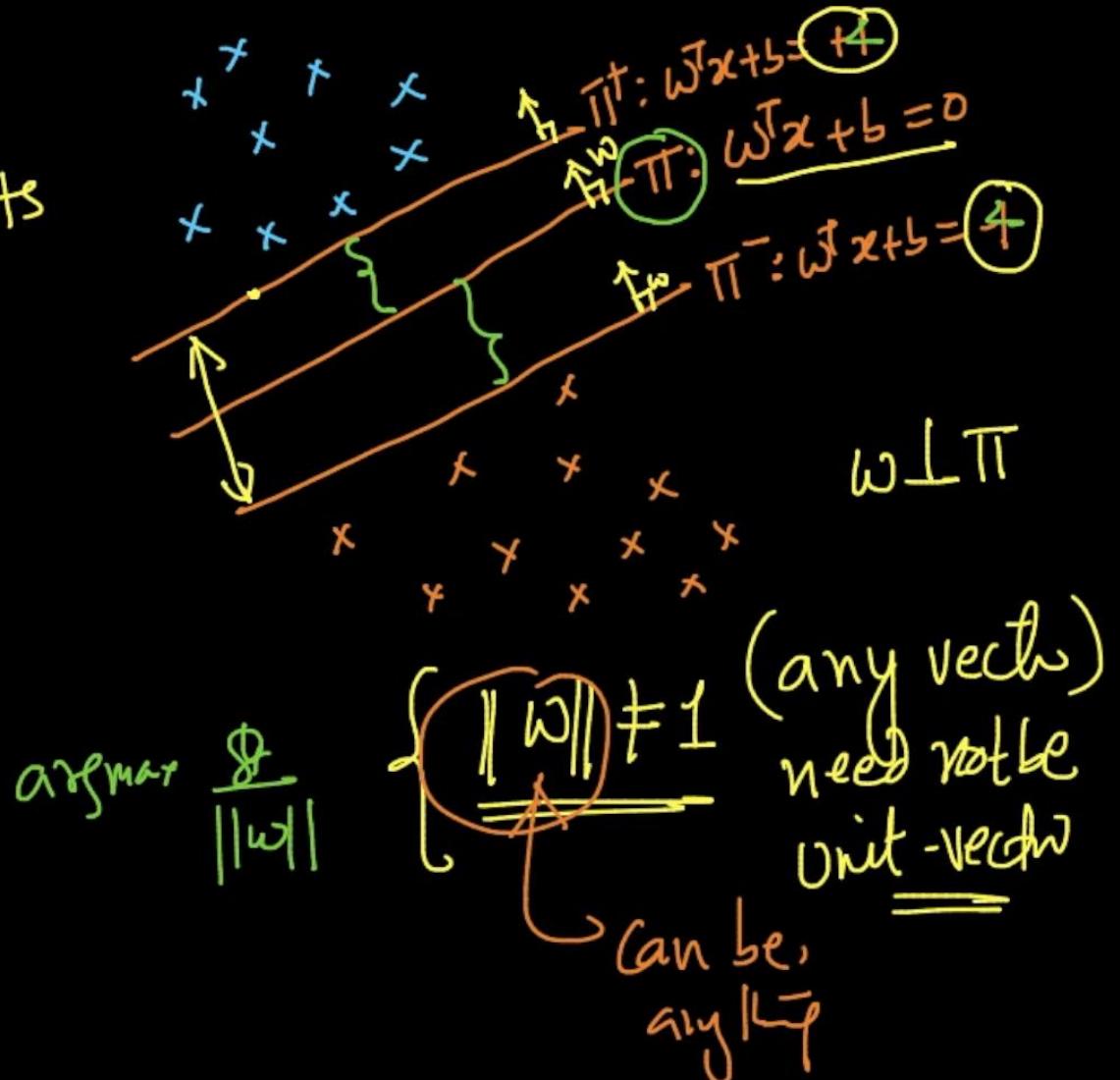
Convenience

(Q) Why $+1 \leq -1$ on the RHS
of $\Pi^+ \leq \Pi^-$

Margin: - $\frac{2}{\|\mathbf{w}\|}$

$$\mathbf{w}^\top, b^\star = \arg \max_{\mathbf{w}, b} \frac{2}{\|\mathbf{w}\|}$$

constants



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

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

$\left\{ \begin{array}{l} \underline{\text{Case1: } z_i > 1 \Rightarrow | -z_i | \text{ is -ve value} \Rightarrow \max(0, | -z_i |) = 0} \\ \underline{\text{Case2: } z_i < 1; | -z_i | > 0 \Rightarrow \max(0, | -z_i |) = | -z_i |} \end{array} \right.$



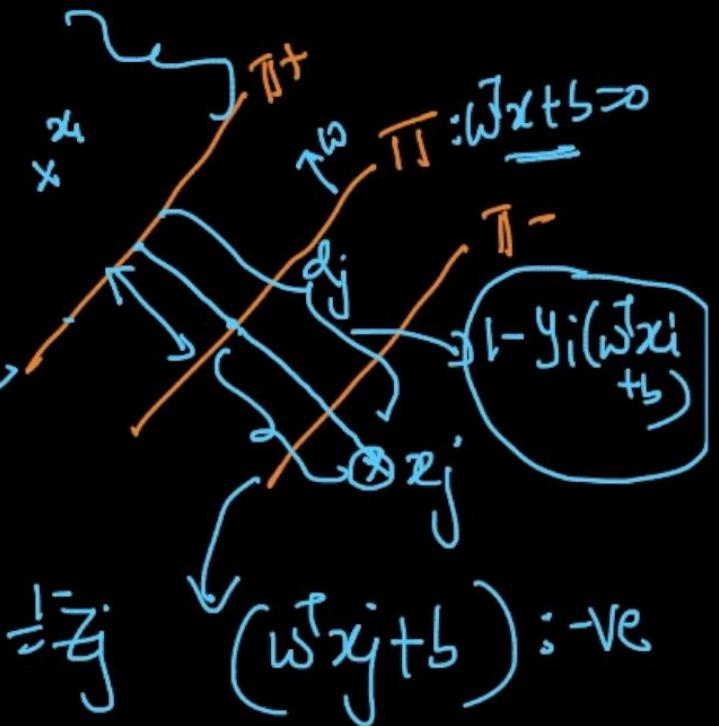
geom-formulation & loss-min-

corr classif \rightarrow $\sum_i \xi_i = 0$ $\leftarrow z_i$

$$\xi_j = 1 - \underbrace{y_i(\vec{w}^T \vec{x}_i + b)}_{z_i} = 1 - z_j$$

$$\xi_j = \text{dist from } \vec{x}_j \text{ to the } \Pi^+ = d_j = \vec{z}_j$$

$$\xi_j = 1 - z_j \rightarrow \text{when } \vec{x}_j \text{ is misclassified}$$



Soft-SVM:

$\left\{ \begin{array}{l} C \uparrow \Rightarrow \text{Overfit} \\ C \downarrow \Rightarrow \text{Underfit} \end{array} \right.$

Loss-Min:

$\left\{ \begin{array}{l} \lambda \uparrow \Rightarrow \text{Underfit} \\ \lambda \downarrow \Rightarrow \text{Overfit} \end{array} \right.$

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

s.t. $(1 - y_i(\omega^\top x_i + b)) > \xi_i \quad \forall i$

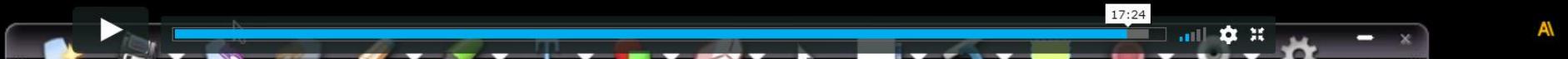
$\xi_i > 0$

$\sum_i \cdot = 0$; corr.
 $\xi_i > 0$; inner

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

$$\|\omega\| > 0 \Rightarrow \min \frac{\|\omega\|}{2}$$

is same as $\min \|\omega\|^2$



Dual form of SVM :: Equivalent :: Dual

Soft-margin SVM

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

s.t. $y_i(w^\top z_i + b) \geq 1 - \xi_i$, $\xi_i \geq 0$

Primal of SVM

$$\checkmark x_{qj}: \underset{\|y_q\|}{\text{Sign}}(w^\top x_{qj} + b) = f(x_{qj})$$

Dual

$$\max_{\alpha} \sum_{i=1}^n \alpha_i - \frac{1}{2} \sum_{i=1}^n \sum_{j=1}^n \alpha_i \alpha_j y_i y_j z_i^\top z_j$$

s.t. $\alpha_i \geq 0$

- ④ $\alpha_i > 0$ only for SVs
- $\sum_{i=1}^n \alpha_i y_i = 0$
- $\alpha_i = 0$ for non-SVs

① $x_i \rightarrow \alpha_i$

② x_i 's only occur in the form of $x_i^\top x_j$

③ $f(x_q) = \sum_{i=1}^n \alpha_i y_i x_i^\top x_q + b$

→ The most imp. idea in SUM is kernel-trick

Soft-SUM-hyperplanes \approx log-reg
↳ margin-max.

by-SUM:- $x_i^T x_j$

$$K(x_i, x_j) = \underline{x_i^T x_j}$$

Kernel-SUM:- $\underline{\underline{K(x_i, x_j)}}$

hyperplane

Lr-SUM :- margin-max-hyperplane

log-reg :- min-logistic loss

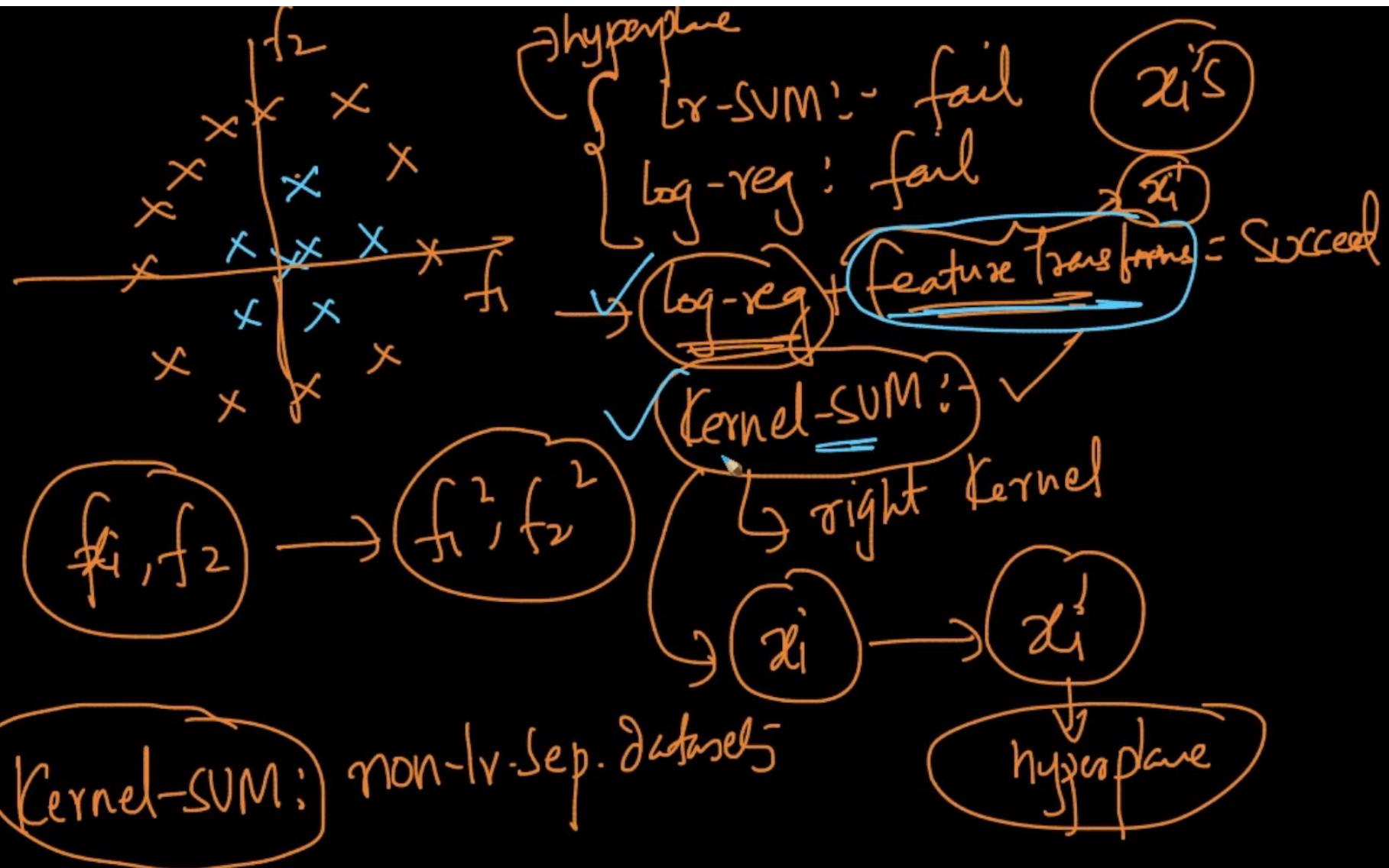
x_i 's

x_i 's

World changing idea :- Kernelization \rightarrow 1990's

x_i 's :-

Lr-SVM



Kernalization:- SUM handle non-lv. Sep. datasets

Kernel logistic-regn

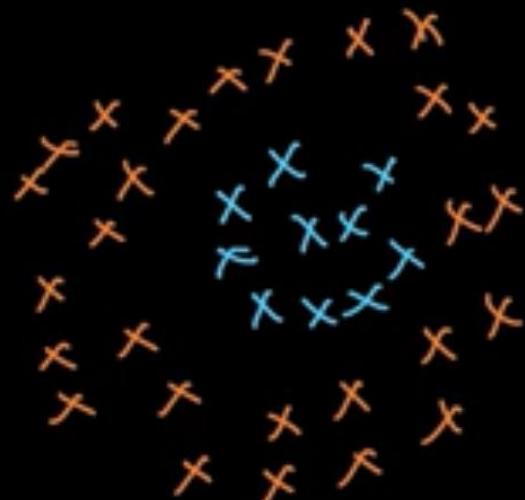
$$x_i^T x_j \rightarrow k(x_i, x_j)$$

Polynomial Kernel

→ Kernelization

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

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



$$\left| \begin{array}{l} (\tilde{f}_1, \tilde{f}_2) \xrightarrow{\text{map}} (\tilde{f}_1^2, \tilde{f}_2^2) \\ \downarrow \text{log-reg} \end{array} \right.$$

$$k(x_1, x_2) = \left(1 + \underline{x}_1^\top \underline{x}_2\right)^2$$

$\underline{x}_1 = \langle x_{11}, x_{12} \rangle$ (2D)

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

$\underline{x}_2 = \langle x_{21}, x_{22} \rangle$ (2D)

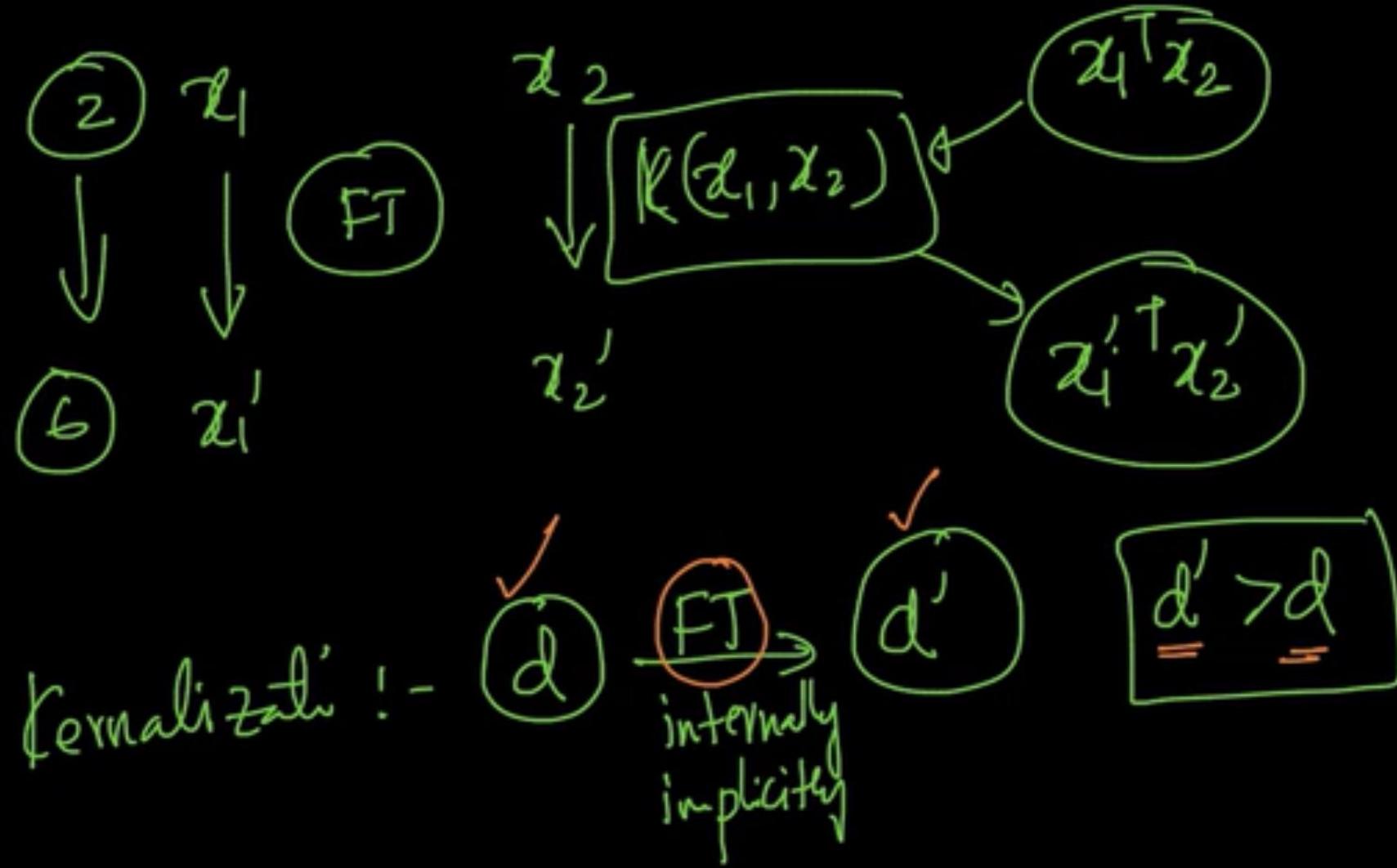
$$= \overbrace{1 + \underline{x}_{11}^2 x_{21}^2 + \underline{x}_{12}^2 x_{22}^2 + 2\underline{x}_{11}x_{21} + 2\underline{x}_{12}x_{22}}^{+ 2x_{11}x_{21}x_{12}x_{22}}$$

Let $\begin{bmatrix} 1, \underline{x}_{11}^2, \underline{x}_{12}^2, \sqrt{2}x_{11}, \sqrt{2}x_{12}, \sqrt{2}x_{11}x_{12} \end{bmatrix} : \underline{x}_1'$

$\begin{bmatrix} 1, x_{21}^2, x_{22}^2, \sqrt{2}x_{21}, \sqrt{2}x_{22}, \sqrt{2}x_{12}x_{22} \end{bmatrix} : \underline{x}_2'$

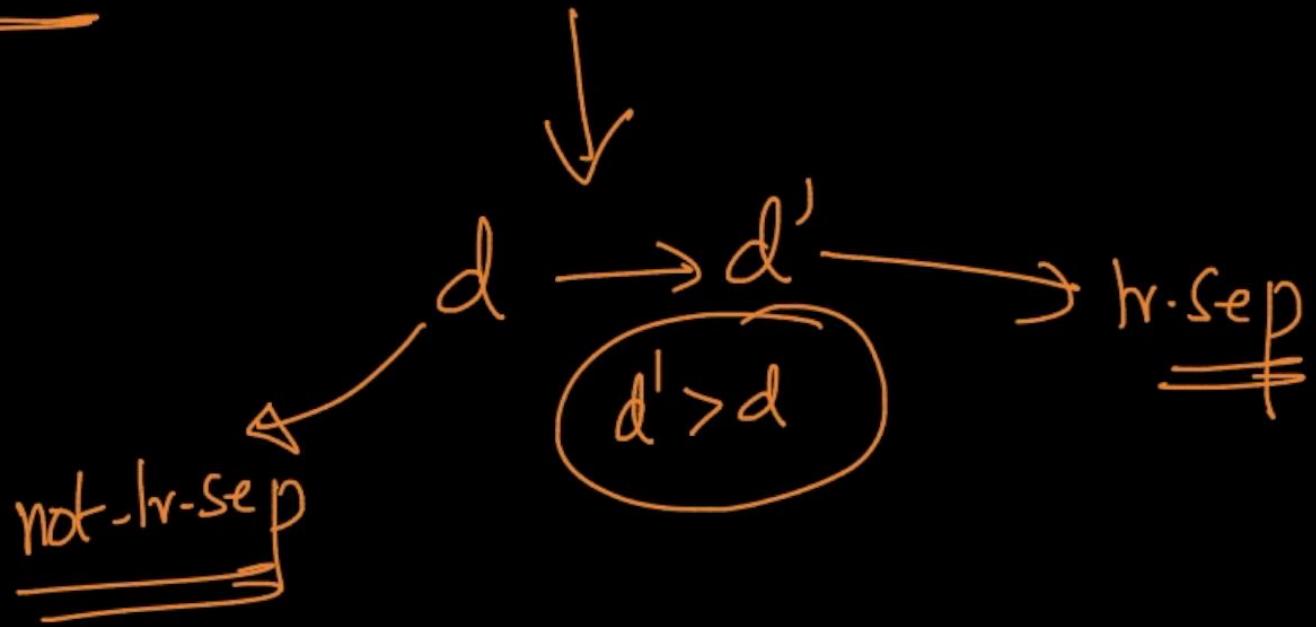
$$= (\underline{x}_1')^\top (\underline{x}_2')$$





Mercer's Thm:

Kernel-trick



Polynomial Kernel

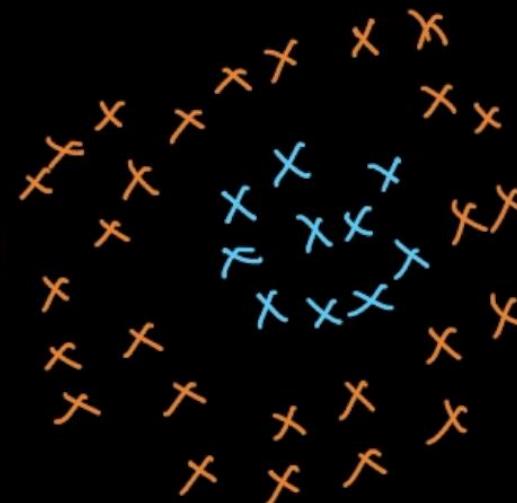
→ Kernelization

→ implicit FT

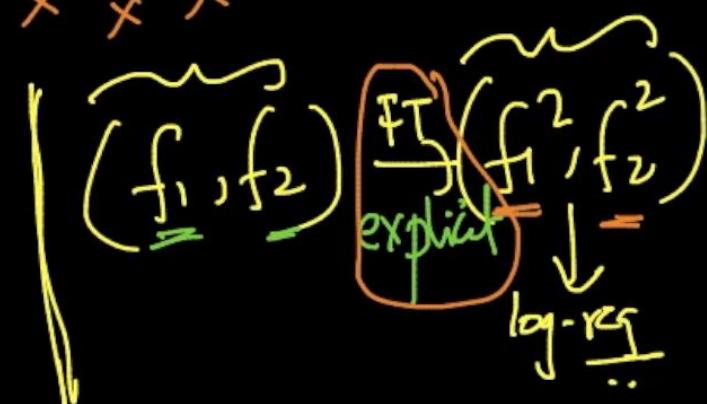
$$K(x_1, x_2) = (x_1^T x_2 + C)^d$$

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

bias-var



Poly K of degree 2



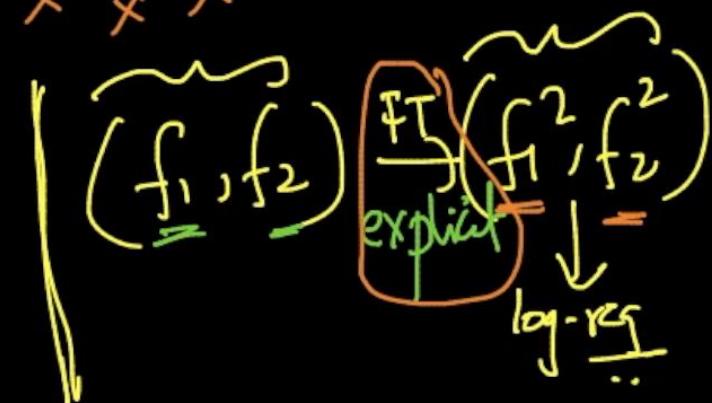
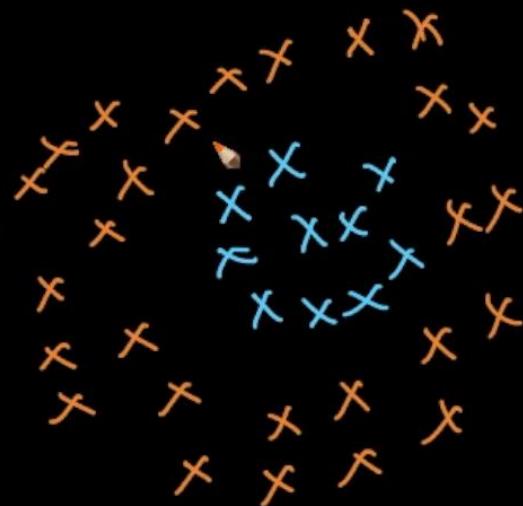
Polynomial Kernel

→ Kernelization

→ implicit FT

$$K(x_1, x_2) = \underline{(x_1^T x_2 + \underline{\underline{C}})}^d$$

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



Radial Basis Function (RBF)

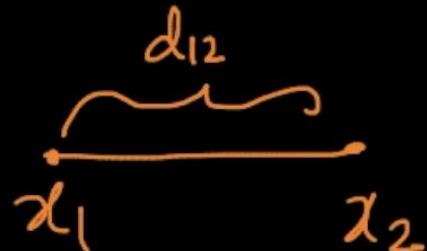
SVM:- most popular / general-purpose : RBF

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

\uparrow hyper-param

Soft-Margin Kernel
SVM

(RBF Kernel
 $\rightarrow C$: hyperparam



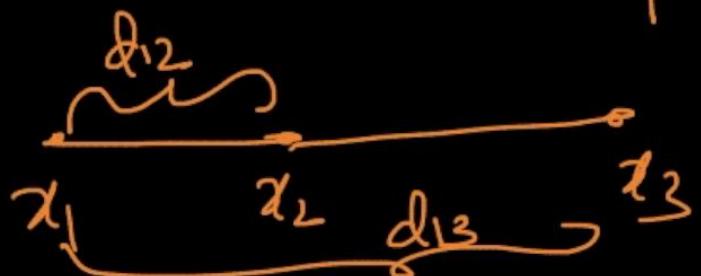
$$\|x_1 - x_2\|^2 = d_{12}^2$$

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

①

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

Similitud



$$\underline{K(x_1, x_2)} > \underline{K(x_1, x_3)}$$

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

$$d \uparrow; d^2 \uparrow$$

$$d^2 \uparrow$$

$$\frac{1}{e^{d^2}} \downarrow$$



② Γ :-

$$\left\{ \begin{array}{l} \Gamma = 1 \\ \Gamma = 0.1 \\ \Gamma = 10 \end{array} \right.$$



Contents - Google Docs plot(exp(-x^2/(2*1)) - Google... plot(exp(-x^2/(2*0.01)) - Google... plot(exp(-x^2/(2*100)) - Google... Chekuri Srikan...

Secure | https://www.google.co.in/search?ei=Jw03WuWqGcmLvQTi1oSwBg&q=plot%28exp%28-x%5E2%2F%282*1%29%29&oq=plot%28exp...

GOOG

plot(exp(-x^2/(2*1))

All Maps Images Videos News More Settings Tools

About 7,26,00,000 results (0.86 seconds)

Graph for $\exp(-(x^2)/(2*1))$

x: 8.53359305 y: 1.53765e-16

-12 -10 -8 -6 -4 -2 2 4 6 8 10 12

More info

plot x e^-x, x^2 e^-x, x=0 to 8 - Wolfram|Alpha

APPLIED COURSE

Contents - Google Docs plot(exp(-x^2/(2*1)) - Google S... plot(exp(-x^2/(2*0.01)) - Google S... plot(exp(-x^2/(2*100)) - Google S... Chekuri Srikan...

Secure | https://www.google.co.in/search?ei=Jw03WuWqGcmLvQTi1oSwBg&q=plot%28exp%28-x%5E2%2F%282*1%29%29&oq=plot%28exp...

GOOG

plot(exp(-x^2/(2*1))

All Maps Images Videos News More Settings Tools

About 7,26,00,000 results (0.86 seconds)

Graph for $\exp(-(x^2)/(2*1))$

$\text{R} = | ;$

$\text{RBF} \sim \text{gaussian PDF}$

More info

plot x e^-x, x^2 e^-x, x=0 to 8 - Wolfram|Alpha

Contents - Google Docs plot(exp(-x^2/(2*1)) - Google... plot(exp(-x^2/(2*0.01)) - Google... plot(exp(-x^2/(2*100)) - Google... Chekuri Srikan...

Secure | https://www.google.co.in/search?ei=MQ03Ws2jBoH6vgSBoI3YCw&q=plot%28exp%28-x%5E2%2F%282*0.01%29%29&oq=plot%28e...

GOOG

plot(exp(-x^2/(2*0.01))

All Maps Videos Images News More Settings Tools

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

Graph for $\exp(-(x^2)/(2*0.01))$

$\sigma = 0.1$

$\sigma^2 = 0.01$

$dist > 1; k = 0$

y = e^{-x^2} - MATLAB Answers - MATLAB Central - MathWorks

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

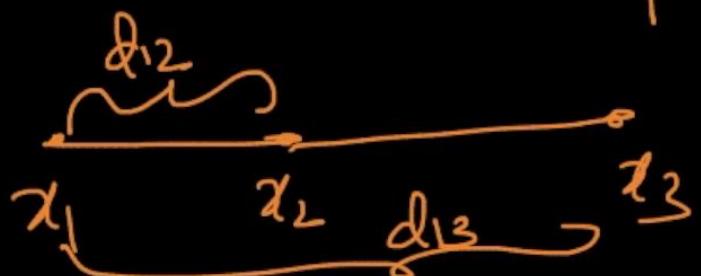
$$\|x_1 - x_2\|_2$$

①

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

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

$$d \uparrow; d^2 \uparrow$$

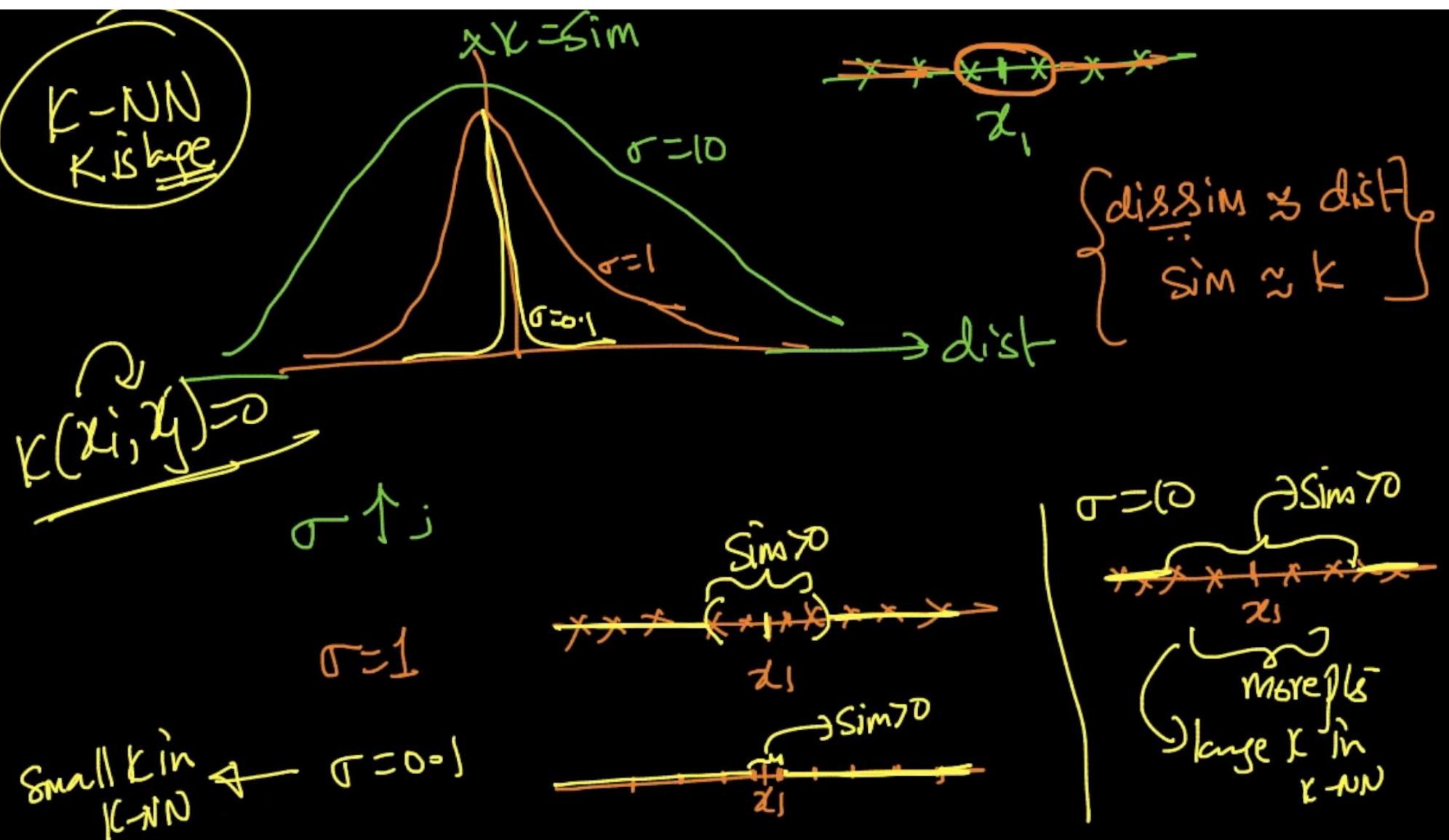


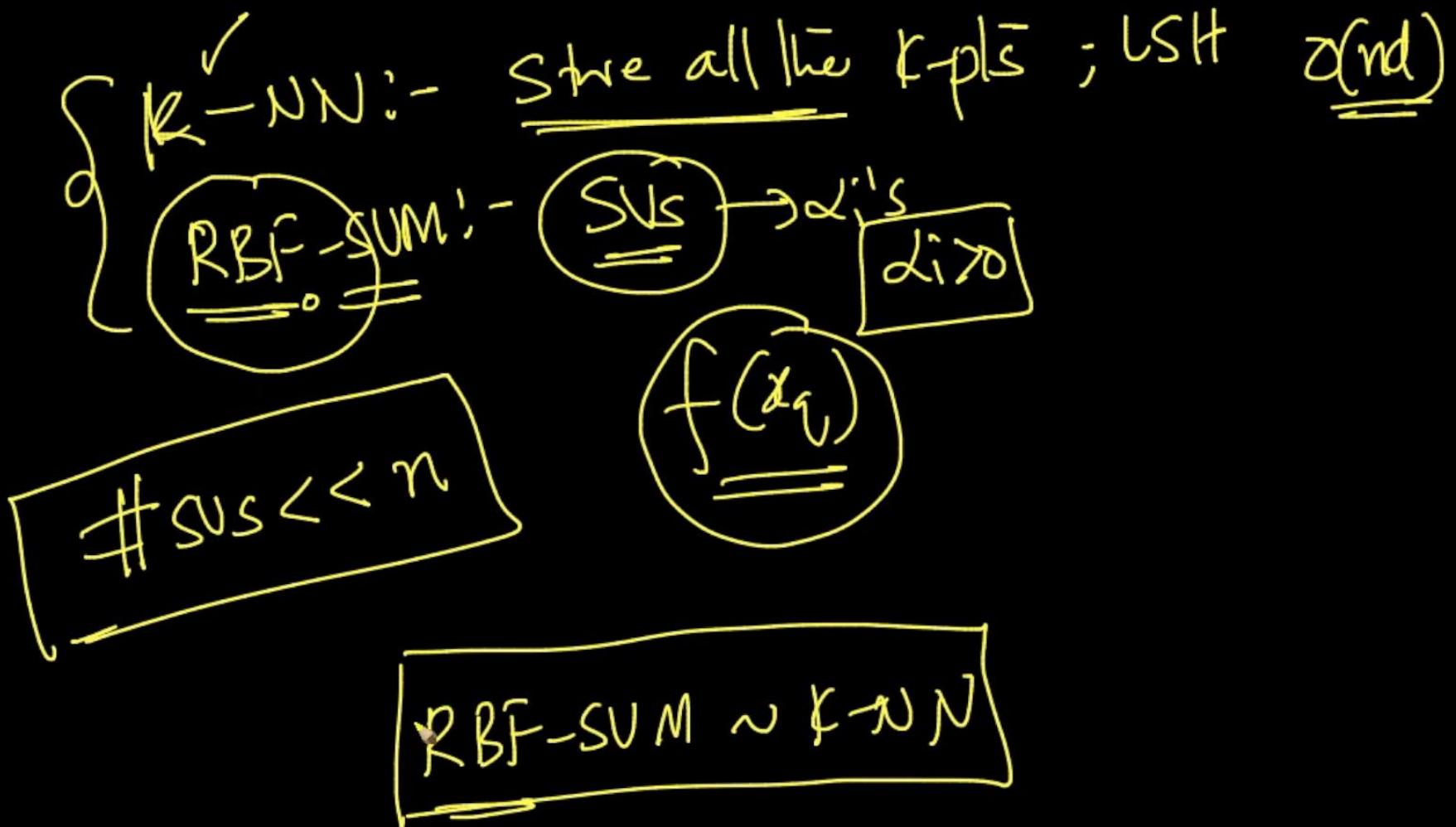
$$\underline{K(x_1, x_2)} > \underline{K(x_1, x_3)}$$

$$d^2 \uparrow$$

$$\frac{1}{e^{d^2}} \downarrow$$





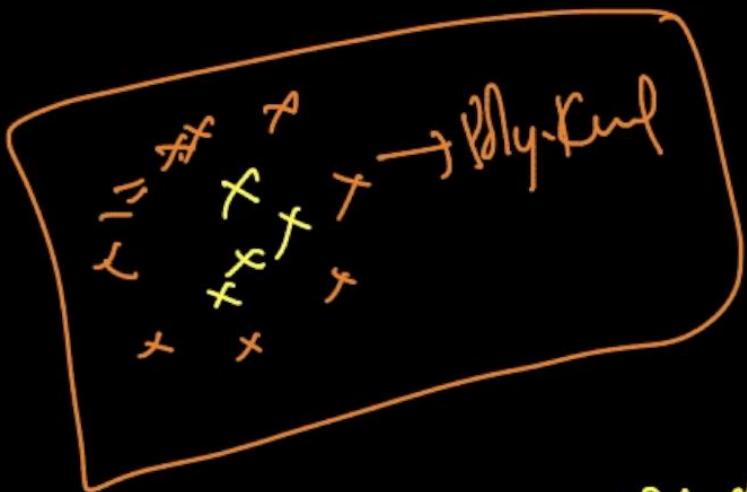


✓ don't know the "best" kernel to use



Simply use RBF-SUM

~kNN



Soft-Margin

(C)

(σ) : RBF

log-reg elastic net: λ_1, λ_2

✓ $C, \sigma \rightarrow$ grid search
Rad search



✓ Domain Specific Kernels: ✓

Polynomial ;

RBF

general purpose kernels

KNN

given the
problem

Kernel - function

FT

Domain Specific

Alt

{ Specialized }
Kernels

SUM in 1900s

Stein Kernels

Text
Classification

✓ Domain Specific Kernels: ✓

Polynomial ;



↳ KNN

given the problem

Kernel - find

FT

Domain Specific

Alt

{ Specialized }
Kernels

SUM in 1900s

Stein Kernels

Text
Classif

SUM

Train & RunTime complexity of SVMs
SMO, SVM

Train: \rightarrow SGD
 \downarrow specialized algo (dual) \rightarrow Sequential minimal optimization (SMO)

✓ libSVM: best libraries for training SVMs

sklearn

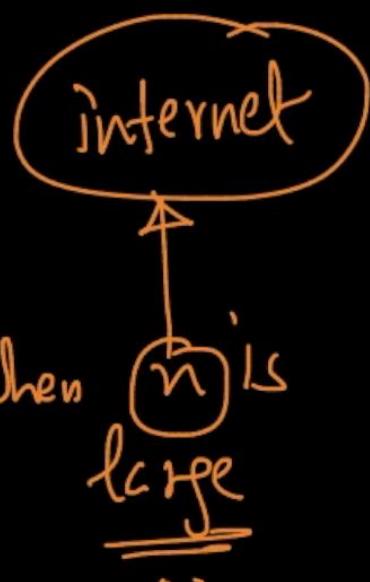


Training Time:- $\sim O(n^2)$ for Kernel-SUMS

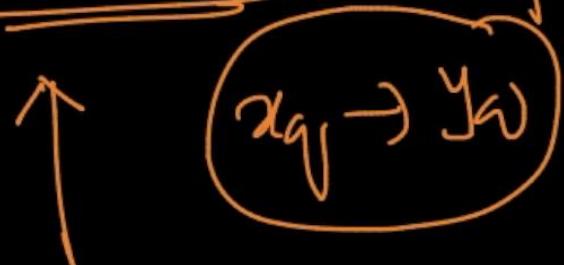
{more opt (2007) :- $O(nd^2)$ if $d < n$

{ - if n is large $\rightarrow O(n^2) \underline{\underline{M}}$

↑ "Typically" don't use "SUM" when n is large ..



Run-time:



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

$x_i = 0$ for non SVs

SVs = k

~~log-regression~~
D(d)

D(Fd)

if # SVs of small

$l \leq k \leq n$

n = 100,000

✓ RBF

Soft-SUM-formula

C, σ

$k = \frac{1}{n}$

nu-SUM



nu-SVM:

(C) SVM \rightarrow original form \vdash

alternative formulation of SVM

hyperparam

CV

nu-SVM:-

$$0 \leq \text{nu} \leq 1$$

$\text{nu} >$ fraction of errors

$\text{nu} \leq$ fraction of SVs

$\text{nu} = 0.0 \Rightarrow$ %age of errors = 1%
 $\# \text{SVs} \geq 1\% \text{ of } n$

$$C \gg 0 \quad \checkmark$$

SUM \rightarrow D Train

% of errors

$$\text{nu} = 0.1$$

% error

$$\text{nu} = 0.01$$

Run-time complx :- fewer SVS

$$nu = 0.01 \Rightarrow$$

errors: $\leq 1\%$.
SVS $> 10^{-8} f n$

$$n = 100, 000$$

$$\# SVS = k \geq 100$$



✓ Support Vector Regression (SVR)

$y_i \in \mathbb{R}$

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

Kernalized

Math:
(x, form of
SVR)

$$\min_{w,b} \frac{1}{2} \|w\|_n^2$$

$$s.t \quad y_i - (\bar{w}^T x_i + b) \leq \epsilon$$

hyperplane

$\epsilon > 0$ ✓

$$f(x_i) = \bar{w}^T x_i + b \\ = y_i$$

$$\begin{cases} y_i - \hat{y}_i \leq \epsilon \\ \hat{y}_i - y_i \leq \epsilon \end{cases}$$



→ Interpretability & feature importance → Kernel-SUM
↳ F.I for various features
↳ forward feature Selcn.

→ Outliers → v. little impact
↳ SVs that matter
↳ RBF with a small σ → k-NN with small K

→ Bias-Vari:-
C ↑ ⇒ overfit ⇒ high var
(grindsmall)
C ↓ ⇒ underfit ⇒ high bias

RBF-SUM:- C, Γ , ζ ↳ RBF
(loss)

→ large \underline{d} \rightarrow v. good for SUM
 \rightarrow kernels (RBF) $d \hookleftarrow d'$

→ Best cases: {right kernel} ✓

→ Worst cases:- n is large \rightarrow Train time is ~~high~~

logistic regn

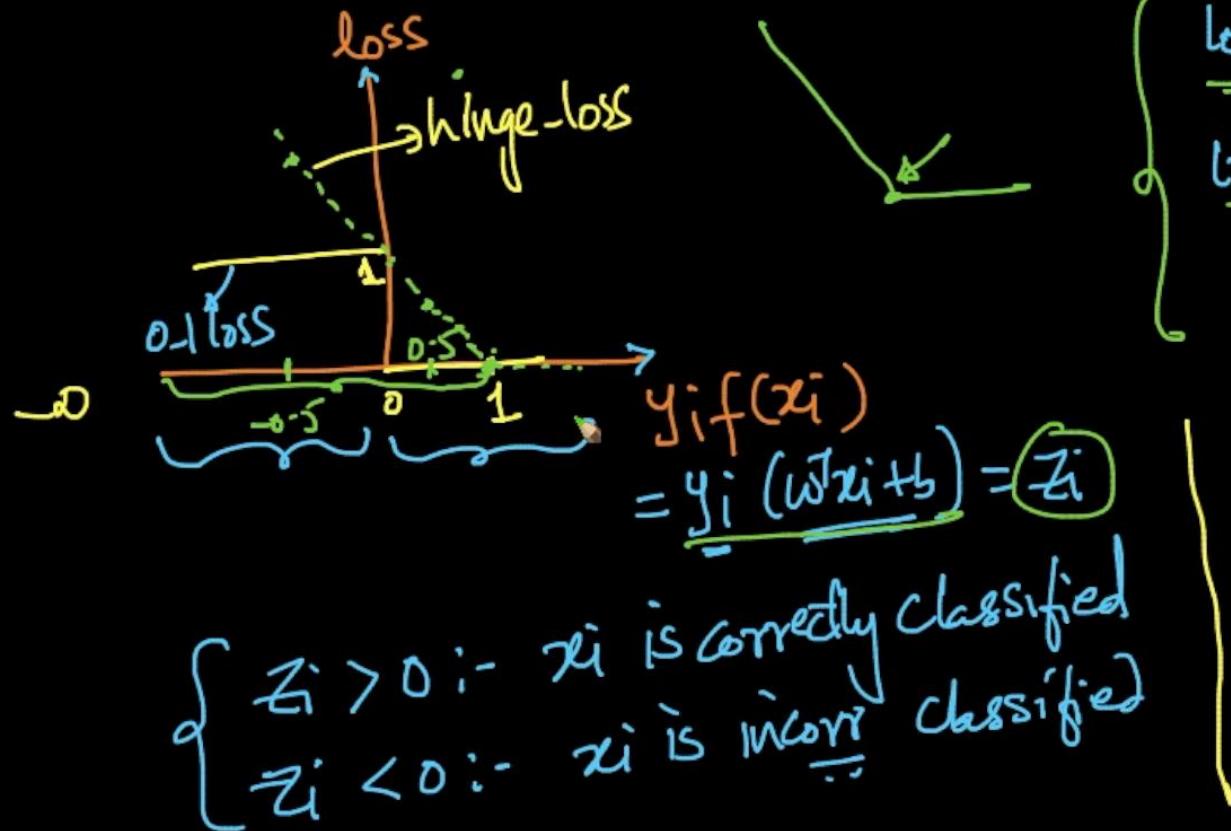
internet based

K is large \Rightarrow $\#SVS$

$f(x_1)$

low-latency is not possible

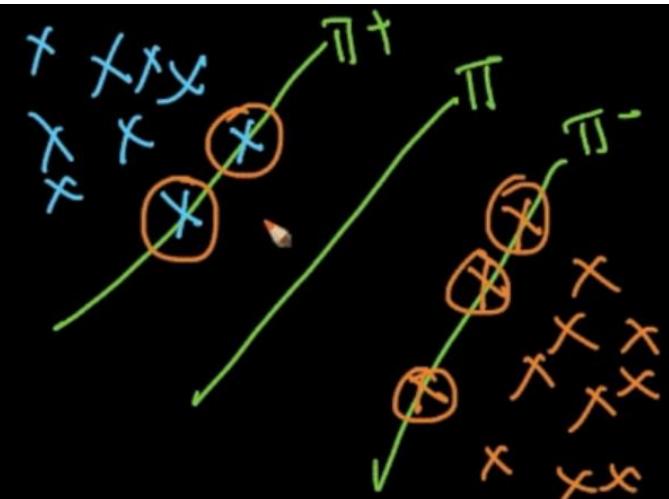
Loss-minimization = hinge-loss



{

- log-reg:- logistic loss + reg
- lr-reg:- lr. loss + reg
- SUM:- hinge-loss + reg

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

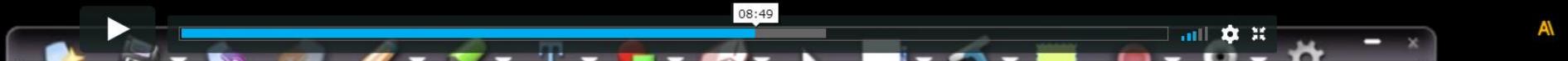


SVs :- $\alpha_i > 0$

non SVs :- $\alpha_i = 0$

$f(x_0)$:- only pts that matter are SVs

SUM



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

Kernel-fn

Sim(x_i, x_j)

s.t. $\alpha_i > 0 \rightarrow \begin{cases} \alpha_i = 0 \text{ f/svs} \\ \alpha_i > 0 \text{ f/nonsvs} \end{cases}$

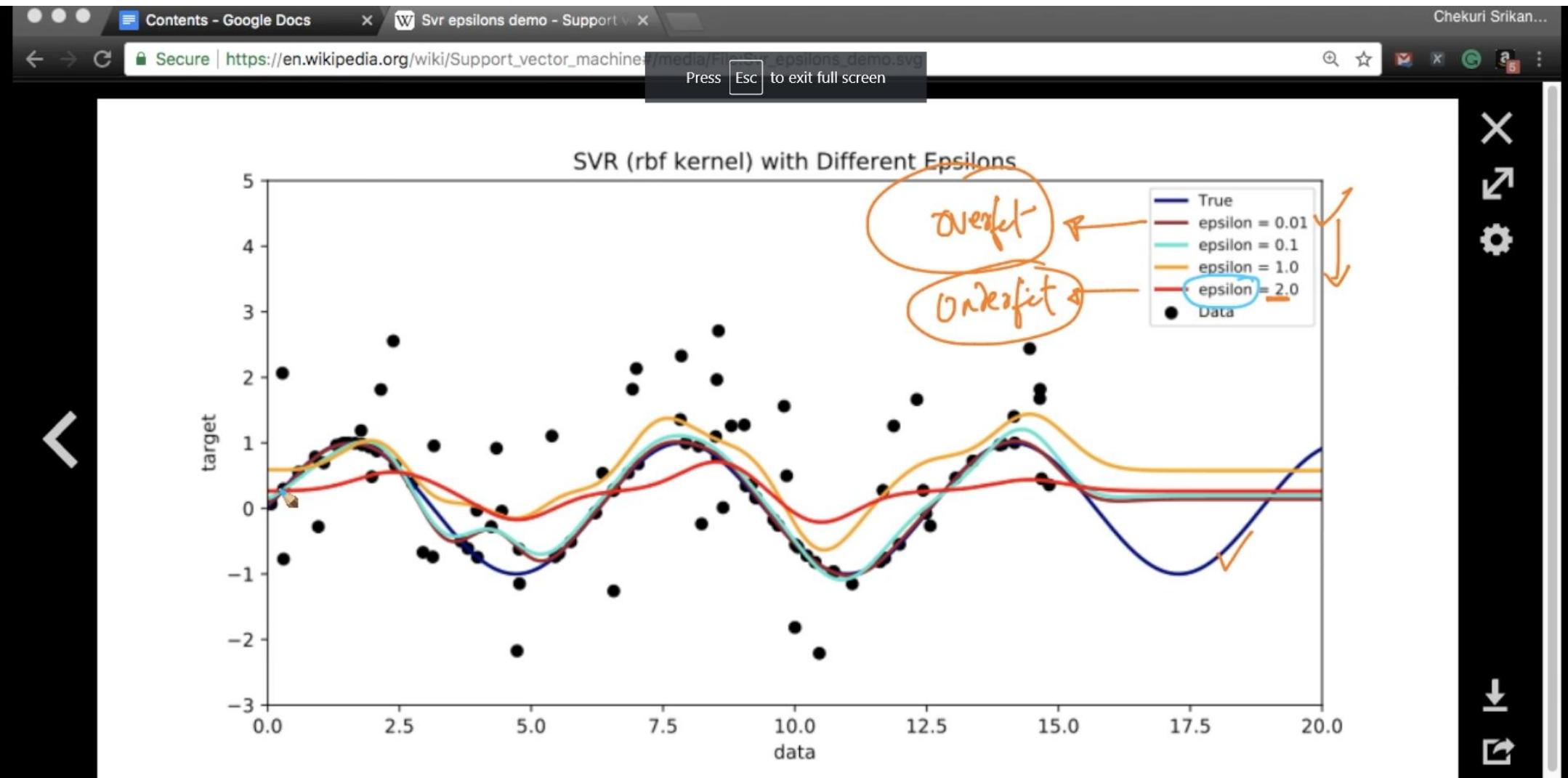
$$\sum_{i=1}^n \alpha_i y_i = 0$$

$$\{ x_i^T x_j = x_i \cdot x_j = \overbrace{\text{Cosine-Sim}(x_i, x_j)}^{\text{if } \|x_i\| = 1, \|x_j\| = 1}$$

Run-time

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





Support Vector Regression (prediction) with different thresholds ε . As ε increases, the prediction becomes less sensitive to errors.

More details

