

$$f(x) = \beta_0 + \beta^T x$$

(margin)

$$L = \{x: \hat{\beta}_0 + \hat{\beta}^T x = 0\}$$

Optimal separating hyperplane
Maximize the distance
of closest point to
hyperplane.

(Thickness should be same for both hyperplanes)

$\beta_0, \beta, \max M$

($\|\beta\|=1$)

subject to $\gamma_i (x_i^T \beta + \beta_0) \geq M$ for all $i=1, \dots, N$

β , should n't be
larger.

We need to maximise
 M , so that there is
a maximum separation
between class.

distance of any x_i from hyperplane

$$\frac{f(x)}{\|f(x)\|} = \frac{(x_i^T \beta + \beta_0)}{\|\beta\|} \rightarrow \text{signed distance}$$

Need n't worry about $\|\beta\|=1$

multiply by γ_i to make
positive

$$\frac{\gamma_i (x_i^T \beta + \beta_0)}{\|\beta\|} \geq M \Rightarrow \gamma_i (x_i^T \beta + \beta_0) \geq M \|\beta\|$$

~~Just~~ Just
make sure

$$\gamma_i (x_i^T \beta + \beta_0) \geq 1$$

$$\|\beta\| = 1/M$$

Max M

\downarrow
Min β

only change
in direction

$$\left(\begin{array}{l} \min \|\beta\| \\ \text{s.t. } \gamma_i (x_i^T \beta + \beta_0) \geq 1 \end{array} \right)$$

[1]

SVM - interpretation & Analysis

$$s.t. y_i(x_i^T \beta + \beta_0) > 1 \quad \min_{\beta} \|\beta\| \Rightarrow \left\{ \min \frac{1}{2} \|\beta\|^2 \right\} \text{ by our choice}$$

So that it is easily differentiable

Lagrangian:

$$L_P = \frac{1}{2} \|\beta\|^2 - \sum_{i=1}^N \alpha_i [y_i(x_i^T \beta + \beta_0) - 1]$$

Setting derivatives to '0'.

$$\frac{\partial L_P}{\partial \beta} = 0 \Rightarrow \frac{1}{2} \cdot 2 \cdot \beta - \sum_{i=1}^N \alpha_i [y_i x_i] = 0 \quad \text{--- (1)}$$

$$\beta = \sum_{i=1}^N \alpha_i x_i y_i$$

$$\frac{\partial L_P}{\partial \beta_0} = 0 - \sum_{i=1}^N \alpha_i y_i = 0 \quad \text{--- (2)}$$

$$\text{Dual: } L_D := \sum_{i=1}^N \alpha_i - \frac{1}{2} \sum_{i=1}^N \sum_{k=1}^N \alpha_i \alpha_k y_i y_k x_i^T x_k \quad \text{--- (3)}$$

Subject to $\alpha_i \geq 0, \forall i$

$$\alpha_i [y_i (x_i^T \beta + \beta_0) - 1] = 0 \quad \forall i \quad \text{--- (4)}$$

$\alpha_i = 0$ For non-marginal points \rightarrow which are very far from hyperplane.

\downarrow will not contribute to calculating β .

$\alpha_i \neq 0$, x_i is on the margin. \rightarrow (support vectors)

\rightarrow only points on the margin plane contribute to β .

\downarrow points which lies on margin plane.

Calculate β_0 from here.

$$f(x) = x^T \hat{\beta} + \hat{\beta}_0$$

$$= \sum \alpha_i y_i x_i^T x + \beta_0$$

SVM is stable

\downarrow less variance

Support vectors will affect SVM

points outside the margin not affect SVM. only

Lecture-23.

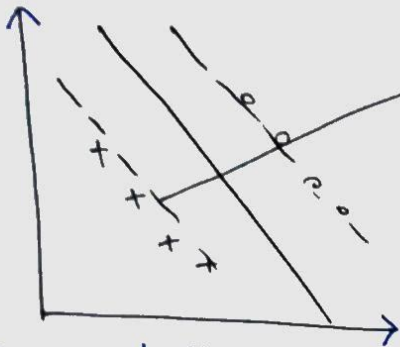
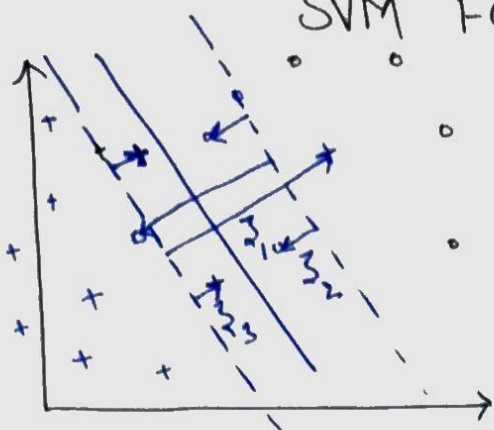
SVM FOR linear Non-seperable data

1. Minimize, these distances.

$\xi_1, \xi_2, \xi_3 \rightarrow$ misclassified samples

Qe? why not minimize the maximum distance rather than minimizing sum of distances.

Ans! classifier will try to minimize the noise, and hyperplane will get shifted towards Noise. like.



Noise

Formulation!

$$y_i (x_i^T \beta + \beta_0) \geq M(1 - \xi_i) \quad \begin{matrix} \nearrow \text{slack} \\ \text{variable} \end{matrix}$$

1. Ideally we want to maximise M , and $\xi_i = 0$

2. but I want some sort of relaxation in maximising M .

why not $(M - \xi_i)$ Ans! Non-Convex optimization loss problem.

ξ_i : what fraction of margin (M)

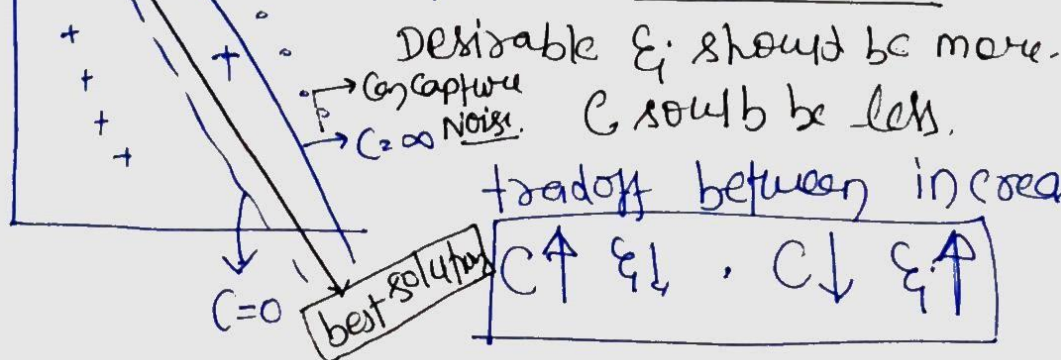
① $\forall_i \xi_i \geq 0$, ② $\sum_{i=1}^N \xi_i \leq \text{constant}$
 \uparrow we don't want ξ_i to be very large.

$$\min \frac{1}{2} \|\beta\|^2 + c \sum_{i=1}^N \xi_i, \text{ sub. to } y_i (x_i^T \beta + \beta_0) \geq (1 - \xi_i)$$

Since $M=1$
 $\|\beta\|$

$\xi_i \downarrow \Rightarrow$ smaller margin. $\xi_i \geq 0$
 ξ_i is very small. \rightarrow less robust

like the previous formulation.

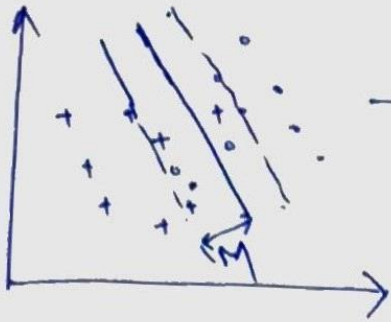


tradeoff between increasing C and ξ_i .

$C \uparrow \xi_i \downarrow, C \downarrow \xi_i \uparrow$

NOTE:

Introducing ξ makes SVM more stable and less prone to noisy data. Suppose: we put $\xi=0$, then we will try to maximize M as shown in figure.



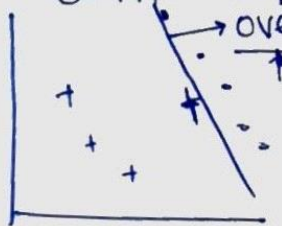
→ we should provide some slackness, because M is calculated via support vectors which lie on the boundary. It doesn't care about the

points which lie within the margin. So, we need to provide some slackness in variable M .

by introducing ξ which is distance of points which are within the margin from margin plane. If ξ ~~max~~ we need sum of ξ_i ~~to~~ to be constant so it shouldn't overshoot.

If we make $C=\infty$, it means all ξ_i are very small, since classifier will try to minimize loss, in doing so, it has to make all $\xi_i=0$, \Rightarrow

the margin will be very small. Overfitting will capture noise also as shown below



If $C=0$, ξ_i will be larger, more biased SVM classifier.

Underfitting

Lecture-24 SVM KERNELS

$$f(x) = x^T \beta + \beta_0 \Rightarrow \sum_{i=1}^N \alpha_i y_i \underbrace{x_i^T x}_{\text{inner product}} + \beta_0$$

is there any good way to compute inner product.

$x \rightarrow h(x)$ "Transformation on x "

$$L_D = \frac{1}{N} \sum_{i=1}^N \alpha_i - \frac{1}{2} \sum_{i=1}^N \sum_{j=1}^N \alpha_i \alpha_j y_i y_j \langle h(x_i), h(x_j) \rangle$$

Similarly

$$f(x) = \sum_{i=1}^N \alpha_i y_i \underbrace{\langle h(x), h(x_i) \rangle}_{\text{inner product}} + \beta_0$$

$$\langle h(x), h(x_i) \rangle = \underbrace{k(x, x_i)}_{\text{kernel}} - \text{similarity measures}$$

$$\underbrace{\text{Kernel}}_{\text{(+ve) semi definite}} - \underbrace{\text{Symmetric}}_{\text{semi definite}} \quad x^T A x \geq 0 \quad \downarrow \text{semi definite}$$

"popular choice of k "

Poly : $(1 + \langle x, x' \rangle)^d$

RBF : $\exp(-\gamma \|x - x'\|^2)$

ANN : $\tanh(k_1 \langle x, x' \rangle + k_2)$
 \downarrow Constant \downarrow Constant

Poly

$$(1 + \langle x_i, x \rangle)^d = (1 + x_1 x'_1 + x_2 x'_2)^2, \quad d=2$$

$$\boxed{X = (x_1, x_2)} = \frac{1}{1} \frac{1 + 2x_1 x'_1}{2} + \frac{2x_2 x'_2}{3} + \frac{(x_1 x'_1)^2}{4} + \frac{(x_2 x'_2)^2}{5} + \frac{2x_1 x'_1 x_2 x'_2}{6}$$

$$h_1(x) = 1, \quad h_2(x) = \sqrt{2} x_1, \quad h_3(x) = \sqrt{2} x_2$$

$$h_4(x) = x_1^2, \quad h_5(x) = x_2^2, \quad h_6(x) = \sqrt{2} x_1 x_2$$

(6-dimension space)

Primal

Lecture-23

$$L_p = \frac{1}{2} \|\beta\|^2 + C \sum_{i=1}^N \xi_i - \sum_{i=1}^N (\alpha_i [\gamma_i (x_i^T \beta + \beta_0) - (1 - \xi_i)]) - \sum_{i=1}^N \mu_i \xi_i$$

Setting derivatives to 0.

$$\beta = \sum \alpha_i \gamma_i x_i \quad (1) \quad 0 = \sum_{i=1}^N \alpha_i \gamma_i \quad (2) \quad \text{we don't need}$$

$$\alpha_i = C - \mu_i \quad (3)$$

$$\sum \xi_i < \text{Constant}$$

It's taken care by 'C' and optimization equation.

Dual:

$$L_D = \sum_{i=1}^N \alpha_i - \frac{1}{2} \sum_{i=1}^N \sum_{j=1}^N \alpha_i \alpha_j \gamma_i \gamma_j x_i^T x_j$$

Subject to

$$0 \leq \alpha_i \leq C < \sum_{i=1}^N \alpha_i \gamma_i = 0$$

KKT Condition

$$C \xi_i = 0 \quad (4) \quad \alpha_i [\gamma_i (x_i^T \beta + \beta_0) - (1 - \xi_i)] = 0 \quad (5) \quad \mu_i \xi_i = 0 \quad (6)$$

$\gamma_i (x_i^T \beta + \beta_0) - (1 - \xi_i) \geq 0$ (1) If $\alpha_i = 0 \Rightarrow x_i$ far away since $\xi_i = 0$ for far points.

if $\gamma_i (x_i^T \beta + \beta_0) \geq 1 \Rightarrow \xi_i = 0 \Rightarrow \alpha_i = 0$

if $\gamma_i (x_i^T \beta + \beta_0) = 1 \Rightarrow 0 \leq \alpha_i \leq C$ (2) $0 < \alpha_i < \infty, \xi_i = 0$

if $\gamma_i (x_i^T \beta + \beta_0) < 1 \Rightarrow$

$\xi_i > 0, \mu_i = 0, \alpha_i = C$

(3) $\gamma_i (x_i^T \beta + \beta_0) - (1 - \xi_i) < 1$ support vector if $\xi_i \neq 0$

Misclassified

in both cases $\alpha_i \neq 0 \Rightarrow$ those x_i 's are support vectors.

$$\alpha_i = C$$

Lecture-24 SVM kernel

Using kernels, we don't need to compute the basis function to transform vectors.

like in previous example, we don't need to compute $h(x)$.

C-SVM → we call it C-SVM - whatever we learnt so far.

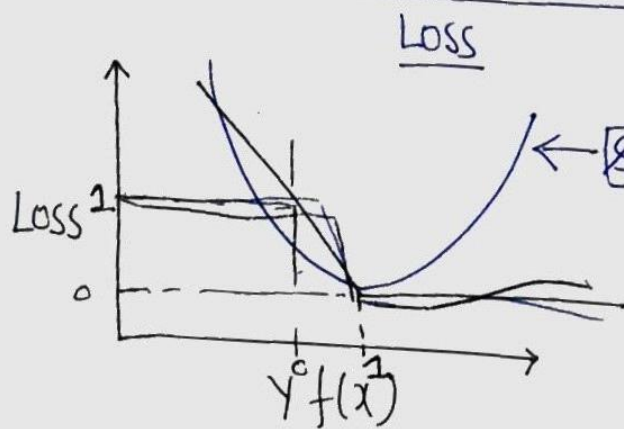
Different type of SVM's are possible where people use different optimisation constraints.

Lecture-25 - Hinge Loss Formulation

$$L_p = \cancel{\sum_{i=1}^N \log \frac{1}{1 + \exp(-y_i(x_i^T \beta + \beta_0))}} = \frac{1}{2} \|\beta\|^2 - \sum_{i=1}^N \left[y_i (x_i^T \beta + \beta_0) \right]$$

$$\min_{\beta, \beta_0} \sum_{i=1}^N \left[1 - y_i f(x_i) \right] + \frac{1}{2} \|\beta\|^2$$

↓ count only when it's positive.



Penalty

← **squared error**

$$\frac{(1 - y f(x))^2}{2}$$