

* If ~~not~~ data is not linearly separable

- Cannot go for zero training error.
- Still learn a maximum margin hyperplane
- 1) Allow some examples to be misclassified.
- 2) Allow some examples to fall inside the margin

* Introducing Slack Variables.

→ Separable case: To ensure zero training loss, constraint was

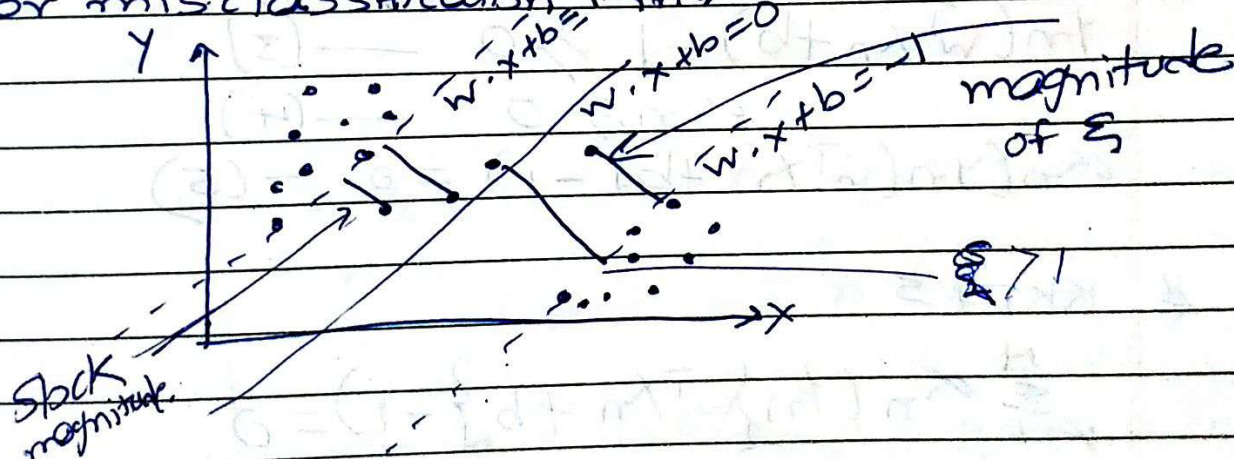
$$y_n (w^T x_n + b) \geq 1 \quad \forall n = 1 \dots N$$

→ Non-Separable case: Relax the constraint

$$y_n (w^T x_n + b) \geq 1 - \xi_n \quad \forall n = 1 \dots N$$

→ ξ_n is called slack variable ($\xi_n \geq 0$) is introduced for non separable

→ for misclassification, $\xi_n \geq 1$



* Relaxing the Constraint.

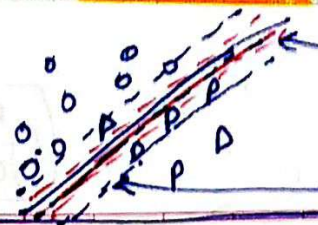
→ It is ok to have some misclassified training example

→ Some ξ 's will be non-zero.

→ minimize the number of such example

→ minimize $\sum_{n=1}^N \xi_n$

C controls
The two
If we want
margin larger
but slack
is less



margin is small $\frac{2}{\|w\|^2}$ small
 $\sum \xi_n$ - slack penalty - small
 $\frac{2}{\|w\|^2}$ is larger $\sum \xi_n$ is large

Optimization Problem for Non-Separable Case

$$\begin{aligned} & \underset{w, b}{\text{maximize}} \quad f(w, b) = \|w\|^2 + C \sum_{n=1}^N \xi_n \\ & \text{Subject to } \gamma_n (w^T x_n + b) \geq 1 - \xi_n, \quad \xi_n \geq 0 \quad n=1, \dots, N \end{aligned}$$

C controls the impact of margin + the margin error

* The Karhuhn-Kuhn-Tucker Conditions (5)

$$\frac{\partial}{\partial w} L_p(w, b, \alpha) = w - \sum_{n=1}^N \alpha_n \gamma_n x_n = 0 \quad (1)$$

$$\frac{\partial}{\partial b} L_p(w, b, \alpha) = - \sum_{n=1}^N \alpha_n \gamma_n = 0 \quad (2)$$

$$\gamma_n (w^T x_n + b) - 1 \geq 0 \quad (3)$$

$$\alpha_n \geq 0 \quad (4)$$

$$\alpha_n (\gamma_n (w^T x_n + b) - 1) = 0 \quad (5)$$

* KKT #5

$$\sum_{n=1}^N \alpha_n (\gamma_n \{w^T x_n + b\} - 1) = 0$$

$\& \alpha_n \geq 0$

$$\alpha_1 (\gamma_1 \{w^T x_1 + b\} - 1) + \alpha_2 (\gamma_2 \{w^T x_2 + b\} - 1) + \dots = 0$$

$$\forall n \quad \alpha_n \geq 0$$

$$\forall n \quad (\gamma_n (w^T x_n + b) - 1) \geq 0$$

$$\Rightarrow \alpha_n (\gamma_n \{w^T x_n + b\} - 1) = 0$$

1) If $\alpha_n = 0$ then $(y_n(w^T x_n + b) - 1)$ can be ≥ 0

2) If $\alpha_n > 0$ then $(y_n(w^T x_n + b) - 1)$ will be 0

x_n is on the margin.

* Key Observation from Dual Formulation

→ KKT conditions # 5

$$\alpha_n (y_n (w^T x_n + b) - 1) = 0$$

→ If x_n not on margin

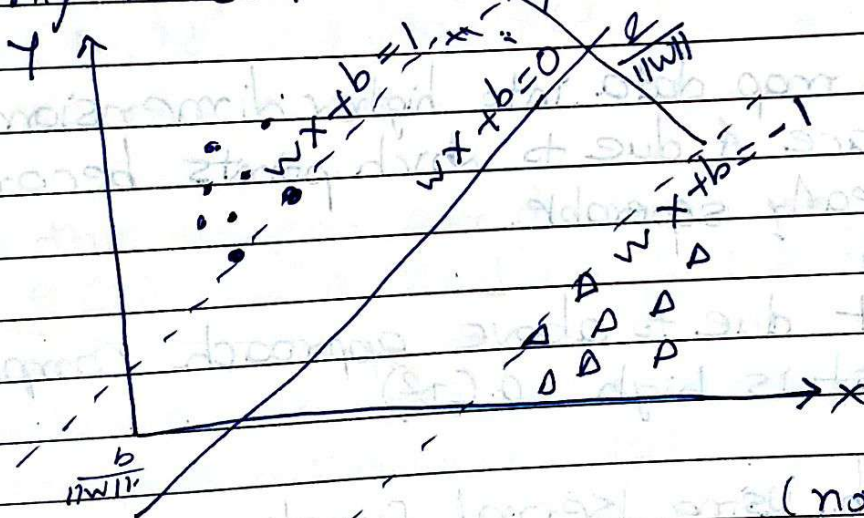
$$y_n (w^T x_n + b) > 1$$

$$\Rightarrow \alpha_n = 0$$

→ $\alpha_n \neq 0$ only for x_n on margin

→ These are the Support Vectors

→ Only need these for prediction.



* In SVM training need to save $\alpha_n + b$ (not w)
 But not α_n only the $\alpha > 0$ means the
 < for support vector

$$= \sum_{\alpha_n > 0} \alpha_n^* (x_n^T x^*) + b$$

↑ for only support vector

Estimating weights

$$w = \sum_{n=1}^N \alpha_n x_n y_n$$

Page:

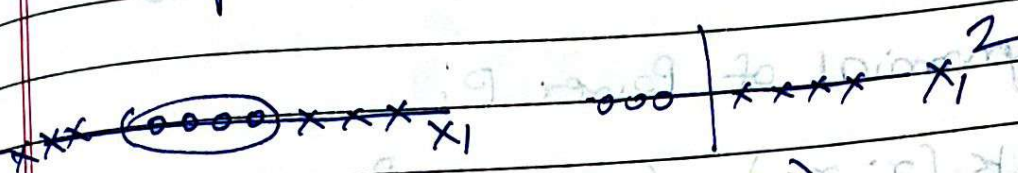
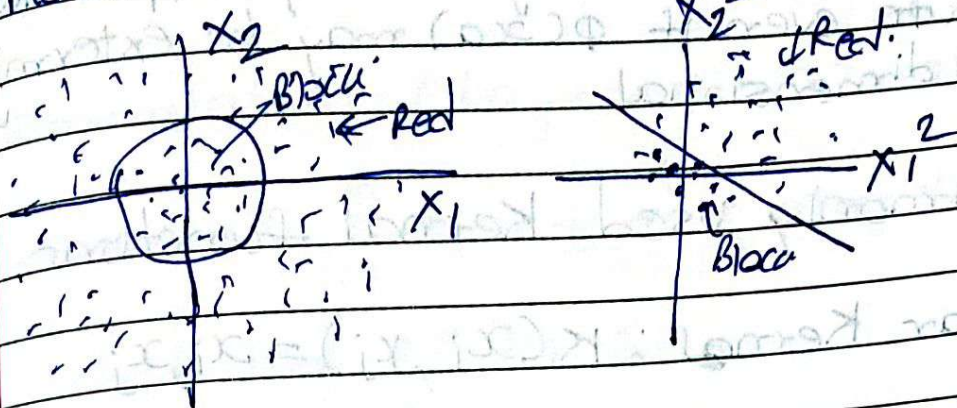
Date: / /

- 1) Points on the margin ($\xi_n = 0$)
- 2) Inside the margin but on the correct side ($0 < \xi_n < 1$)
- 3) On the wrong side of the hyperplane ($\xi_n > 1$)

Nonlinear SVM & Kernel Function SVM

- 1) x is input feature attribute which are truly non separable.
- 2) Suppose x is our I/P features, we can use mapping to map x to $\phi(x)$
 $x \rightarrow \phi(x)$ means transform into new feature space.
- 3) Due to this it is possible I/P feature are linearly separable in the transformed feature space.
- 4) So map data into higher dimensional feature space & due to which points become linearly separable.
- 5) But due to above approach computational cost is high $O(D^2)$
- 6) So by using Kernel function we can achieve this transformation without any major implication on computational cost

* Nonlinear SVM - Feature Space



$$\phi: x \rightarrow \phi(x)$$

So here we can observe if convert feature into high dimensional feature space, the points become linearly separable.

* Kernel Function.

$$K(x_a, x_b) = \phi(x_a) \cdot \phi(x_b)$$

The Kernel Trick

1) With this mapping, our discriminant function

$$g(x) = w^T \phi(x) + b = \sum_{i \in S^+} \alpha_i \phi(x_i)^T \phi(x) + b$$

2) We only use the dot product of feature vectors in both the training & test.

3) A Kernel function is defⁿ as a function that correspond to a dot product of two feature vectors in some expanded feature space

$$K(x_a, x_b) = \phi(x_a) \cdot \phi(x_b)$$

often $K(x_a, x_b)$ may be very inexpensive to compute even if $\phi(x_a)$ may be extremely high dimensional

* Commonly Used Kernel functions

1) Linear Kernel: $K(x_i, x_j) = x_i \cdot x_j$

2) Polynomial of Power P :

$$K(x_i, x_j) = (1 + x_i \cdot x_j)^P$$

3) Gaussian Radial-Basis Function

$$K(x_i, x_j) = e^{-\frac{\|x_i - x_j\|^2}{2\sigma^2}}$$

4) Sigmoid

$$K(x_i, x_j) = \tanh(\beta_0 x_i \cdot x_j + \beta_1)$$

* In general, functions that qualify Mercer's condⁿ is kernel funⁿ.

* Kernel function measures some similarity betⁿ instances x_i & x_j & this kernel functions use.

* Mercer's condⁿ states that any positive semi-definite kernel $f(x, y)$ i.e.

$$\sum_{i,j} K(x_i, x_j) c_i c_j \geq 0$$

can be expressed as a dot product in a high dimensional space

* Nonlinear SVM - Optimization.

→ Formulation (Lagrangian Dual Problem)

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

such that $0 \leq \alpha_i \leq C$

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

The solution of the discriminant funⁿ is

$$g(x) = \sum_{i \in S_V} \alpha_i K(x, x_i) + b$$

* Multi - Class Classification

→ SVM can only handle two - class outputs.

→ Learn M SVM's.

— SVM 1 learns Class 1 vs Rest

— SVM 2 learns Class 2 vs Rest

— ...

— SVM M learns Class M vs Rest.

Then to predict o/p for a new i/p, just predict with each SVM & find out which one puts the prediction the furthest into the positive region.