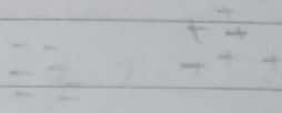


Here we are imposing a strict constraint that no data point is closer than 1 from the D.B. This is called hard margin classifier.

With primal approach, if data is not linearly separable by hyperplane, then this method gives no solution.



Linearly separable

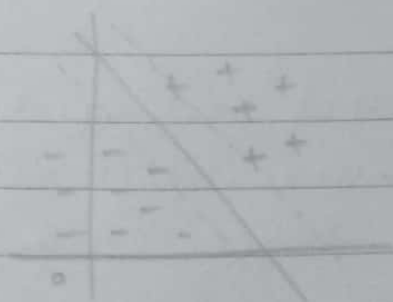


Not linearly separable



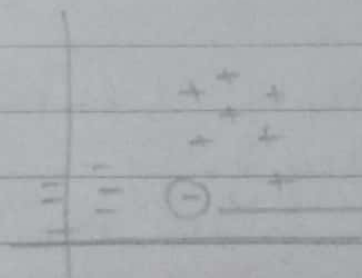
At least for some points, the constraint  $y^{(i)} (\beta^T x^{(i)} + \beta_0) \geq 1$  is not satisfied.

In order to solve this problem when data is not linearly separable, we use soft margin classifier. (SVM with regularisation)

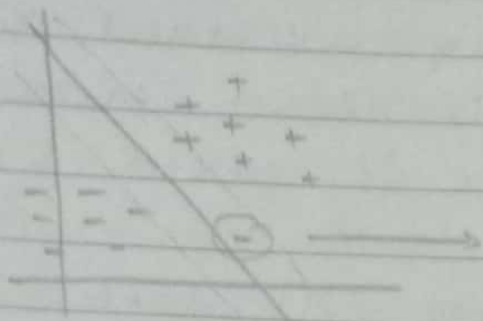


$$\beta^T x + \beta_0 = 0$$

→ when we do not have such a nice case with clear boundary, use soft margin.



outlier: far from its group  
(makes margin very narrow)



we let go and make an error in this case. But in turn we get a decision boundary with a reasonably good margin.

So when no. of outliers and the points lying in completely wrong group are less, we ignore those as errors and still go for a linear D.B.

### Soft margin SVM (SVM with regularisation)

Problem formulation - we still get hyperplane taking care of outliers and non linear data

$$\min_{\underline{P}} \frac{1}{2} \|\underline{P}\|^2 \quad \text{such that} \quad y^{(i)} (\underline{P}^T x^{(i)} + \beta_0) \geq 1 - \epsilon_i$$

for all  $i = 1 \dots m$

$\epsilon_i \geq 0$  is slack parameter

But here we do not get unique  $\underline{P}^*$  and  $\beta_0^*$  because we can have different  $\epsilon_i$ . Hence, we add other term (regularisation)

$$\min_{\underline{P}, \epsilon_i} \left[ \frac{1}{2} \|\underline{P}\|^2 + \underbrace{\left[ \frac{c}{2} \sum_{i=1}^m \epsilon_i \right]}_{\text{Regularisation parameter (chosen, not estimated)}} \right]$$

Penalize when they occur i.e. add weight to a function which is to be minimised

Regularisation parameter (chosen, not estimated)

outliers and misclassified

For softmax : —

→ Objective function :  $\min \frac{\|\beta\|^2}{2} + c \sum_{i=1}^m \epsilon_i$

→ constraints : (1)  $y_i (\beta^T x_i + \beta_0) \geq 1 - \epsilon_i$   
(2)  $\epsilon_i \geq 0, \forall i$

Primal variables —  $\beta, \beta_0, \epsilon_i$

\* Dual variables —  $\alpha_i, \lambda_i$  (2 constraints)

→ Lagrangian is given as follows —

$$L(\beta, \beta_0, \epsilon_i, \alpha_i, \lambda_i) = \frac{1}{2} \|\beta\|^2 + c \sum_{i=1}^m \epsilon_i + \sum_{i=1}^m \alpha_i [1 - \epsilon_i - y_i (\beta^T x_i + \beta_0)] + \sum_{i=1}^m \lambda_i [-\epsilon_i]$$

$$\frac{\partial L}{\partial \epsilon_i} = c - \alpha_i - \lambda_i = 0$$

For dual problem formulation,  
replace the primal variables.

$$L = \frac{1}{2} \sum_{i=1}^m \sum_{j=1}^m \alpha_i \alpha_j y_i y_j \underline{x_i^T x_j} + c \sum_{i=1}^m \epsilon_i$$

$$+ \sum_{i=1}^m \alpha_i - \sum_{i=1}^m \alpha_i \epsilon_i - \sum_{i=1}^m \lambda_i \epsilon_i$$

$$- \sum_{i=1}^m \sum_{j=1}^m \alpha_i \alpha_j y_i y_j \underline{x_i^T x_j}$$

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

$$+ \underbrace{(c - \alpha_i - \lambda_i)}_{=0} \sum_{i=1}^m \epsilon_i$$

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

with  $\alpha_i \geq 0$

$\lambda_i \geq 0$

But no  $\lambda_i$

in this expression so using

$\lambda_i = c - \alpha_i$  we get  $c - \alpha_i \geq 0$

Hence  $\alpha_i \leq c$ .

So  $0 \leq \alpha_i \leq c$

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

(The  $\lambda_i$   
written  
in bold)

## KKT Conditions

$$\frac{\partial L}{\partial \beta} = 0 ; \quad \frac{\partial L}{\partial \beta_0} = 0 ; \quad \frac{\partial L}{\partial \epsilon_i} = 0$$

$$y_i ( \beta^T x_i + \beta_0 ) \geq 1 - \epsilon_i$$
$$\epsilon_i \geq 0$$

Primal feasibility

$$\alpha_i \geq 0$$
$$\lambda_i \geq 0$$

Dual feasibility

$$\alpha_i [ y_i ( \beta^T x_i + \beta_0 ) - 1 + \epsilon_i ] = 0$$
$$\lambda_i \epsilon_i = 0$$

Complementary slackness

$$\alpha_i > 0 \Rightarrow y_i ( \beta^T x_i + \beta_0 ) = 1 - \epsilon_i$$

Two possibilities here

$$A) \epsilon_i > 0 \Rightarrow \lambda_i = 0 \Rightarrow \alpha_i = \gamma$$

$$B) \epsilon_i = 0 \Rightarrow \lambda_i \text{ could be greater than } 0$$
$$\text{so } \gamma - \alpha_i > 0$$
$$\alpha_i < \gamma$$

so all  $\alpha_i > 0$  are support vectors.  
Those in B) are on the margin.  
and those that don't violate margin.  
Those in A) may violate margin.



→ As  $\xi_i$  is close to zero, we approach hard margin.

→ We want to minimize the entire function, so when  $c$  is very large,  $\xi_i$  are small hence we are close to hard margin <sup>close to zero</sup>.

→ As the value of  $c$  decreases, the margin becomes narrower.

For hard margin problem, dual problem formulation is:

$$\max_{\lambda_i} \sum_{i=1}^m \lambda_i - \sum_{i,j} \lambda_i \lambda_j y^{(i)} y^{(j)} x^{(i)\top} x^{(j)}$$

with

$$\lambda_i \geq 0$$

$$\sum \lambda_i y_i = 0$$

→ For soft margin, dual problem is:

$$\max_{\lambda_i} \sum_{i=1}^m \lambda_i - \sum_{i,j} \lambda_i \lambda_j y^{(i)} y^{(j)} x^{(i)\top} x^{(j)}$$

with

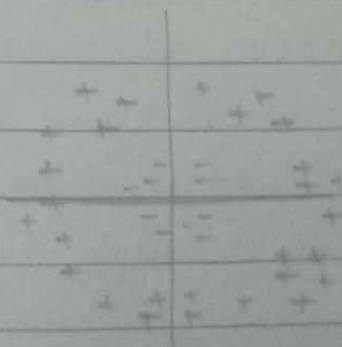
$$0 \leq \lambda_i \leq c$$

$$\sum \lambda_i y_i = 0$$

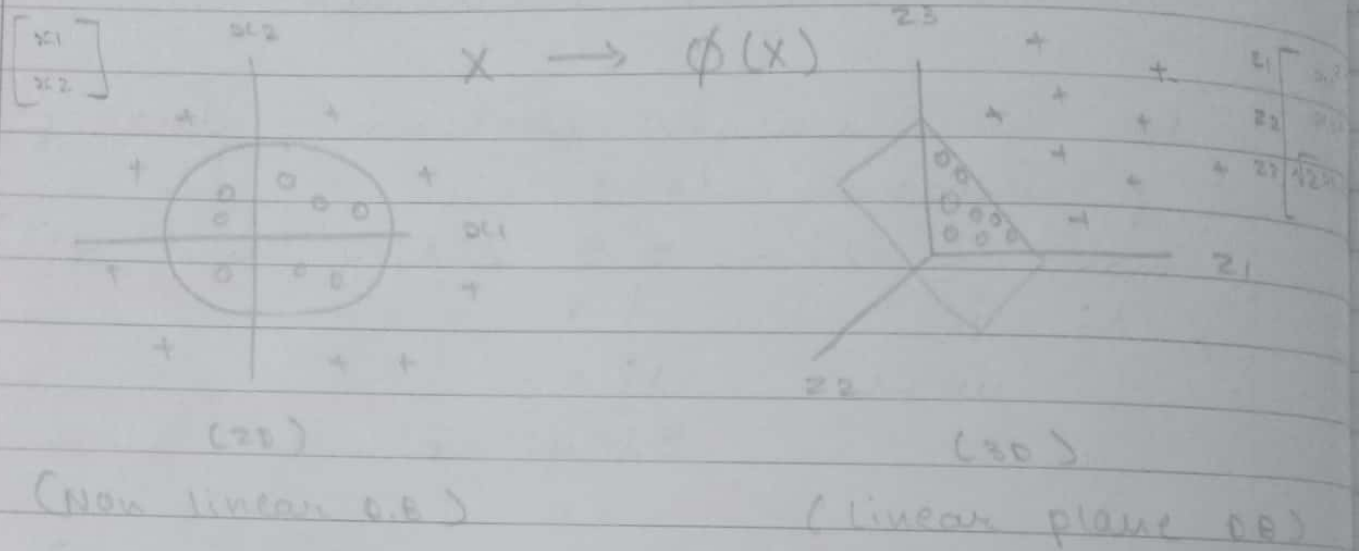
Proof!

### SVM with kernel function

→ This is used for non linearly separable data point. (The best solution - a general method that can be used in ...)



- The trick used here is kernel trick - convert data from lower dimension to higher dimension [Eg. from a vector in 2D to a vector in 4D]
- After this conversion, the data now becomes separable and hence we can fit a decision boundary.



- $\phi$  function need not be known
- Kernel function - Radial basis function (RBF)

Hard margin:

- Data has to be linearly separable

Soft margin:

- Slack variables  $\xi_i$  are used
- Takes care of non linear separation and outliers can be handled well
- Still not the best classifier

Using Kernel function:

- Best svm classifier.
- When data is not linearly separable

onto higher dimensional space where the data becomes linearly separable (Then we can use concepts of hard margin / soft margin classifier for problem formulation). The data now used is  $\phi(x)$ . But  $\phi$  itself is not used due to the use Kernel function

A kernel function is given by:

$$K(\underline{x}, \underline{z}) = [\phi(\underline{x})]^T [\phi(\underline{z})]$$

Thus  $K$  performs dot product of  $\underline{x}$  and  $\underline{z}$  not in lower dimension but in higher dimension. (Here  $\underline{x}$  and  $\underline{z}$  are 2 examples)

Eg.  $\underline{x} = \begin{bmatrix} x_1 \\ x_2 \end{bmatrix}$        $\underline{z} = \begin{bmatrix} z_1 \\ z_2 \end{bmatrix}$

$$\phi(\underline{x}) = \begin{bmatrix} x_1^2 \\ x_2^2 \\ \sqrt{2} x_1 x_2 \end{bmatrix} \quad \phi(\underline{z}) = \begin{bmatrix} z_1^2 \\ z_2^2 \\ \sqrt{2} z_1 z_2 \end{bmatrix}$$

$$\begin{aligned} \phi(\underline{x})^T \phi(\underline{z}) &= x_1^2 z_1^2 + x_2^2 z_2^2 + 2 x_1 x_2 z_1 z_2 \\ &= (x_1 z_1 + x_2 z_2)^2 \end{aligned}$$

Let us consider  $K(\underline{x}, \underline{z}) = (\underline{x}^T \underline{z})^2$

So  $K(\underline{x}, \underline{z}) = (x_1 z_1 + x_2 z_2)^2$

Thus, the kernel operates on lower dimension data, but is actually doing dot product in higher dimension data.

One of the popular kernel  $K(\underline{x}, \underline{z})$  is called radial basis function (RBF) or



$$K(\underline{x}, \underline{z}) = e^{-\|\underline{x} - \underline{z}\|^2 / 2\sigma^2} \quad \text{where } \sigma^2 \text{ must be known}$$

Intuitively, Kernel function is finding out similarity in  $\underline{x}$  and  $\underline{z}$ . When  $\underline{x}$  and  $\underline{z}$  are similar,  $\underline{x} - \underline{z}$  is small close to zero hence  $K(\underline{x}, \underline{z})$  will be close to one (highest possible). This makes sense because when  $\underline{x}$  and  $\underline{z}$  are similar, their dot product should be high.

Similarly,  $K(\underline{x}, \underline{z})$  is small when  $\underline{x}$  and  $\underline{z}$  are not similar.

After converting in higher dimension, we formulate using dual approach in hard margin.

$$\max_{\lambda_i} \sum_{i=1}^m \lambda_i - \sum_{i=1}^m \sum_{j=1}^m \lambda_i \lambda_j y^{(i)} y^{(j)} \underbrace{K(\underline{x}^{(i)}, \underline{x}^{(j)})}_{\downarrow \phi(\underline{x}^{(i)})^T \phi(\underline{x}^{(j)})}$$

( $\underline{x}^{(i)T} \underline{x}^{(j)}$  becomes  $\phi(\underline{x}^{(i)})^T \phi(\underline{x}^{(j)})$  in higher dimension now)

(We do not worry about what the higher dimension is. Kernel function will ultimately give a single value only.)

$$\text{Now, } \underline{B} = \sum_{i=1}^m \lambda_i^* y^{(i)} \phi(\underline{x}^{(i)})$$

But now, there is a problem here because  $\phi$  is unknown. Not a major problem because we actually want the decision boundary not  $\underline{B}$  and  $\underline{B}^*$ .

Hence  $\sum_{i=1}^m \lambda_i y^{(i)} [\phi(x^{(i)})]^T \phi(z) + \beta_0 = 0$

$\downarrow$   
 $K(x^{(i)}, z)$   
 set example

So now if  $\beta \phi^T(z) + \beta_0 > 0$  — classify as 1  
 $\beta \phi^T(z) + \beta_0 < 0$  — classify as 0

We still need to find  $\beta_0$ .

suppose  $x^{(i)}$  is an example in the training set with non zero  $\lambda_i$ , i.e.  $x^{(i)}$  is a support vector.

so  $\beta^T \phi(x^{(i)}) + \beta_0 = 1$

$$\sum_{j=1}^m \lambda_j y^{(j)} [\phi(x^{(i)})]^T \phi(x^{(j)}) + \beta_0 = 1$$

$\downarrow$   
 $K(x^{(i)}, x^{(j)})$

$\downarrow$   
 Find from here

Here we are using Gaussian Kernel.

There are other types of kernel as well:

① Polynomial kernel  $K(x^{(i)}, x^{(j)}) = (1 + x^{(i)T} x^{(j)})^p$   
 where  $p$  is an integer  $> 0$

② Linear kernel  $K(x^{(i)}, x^{(j)}) = x^{(i)T} x^{(j)}$   
 (used in original problem)