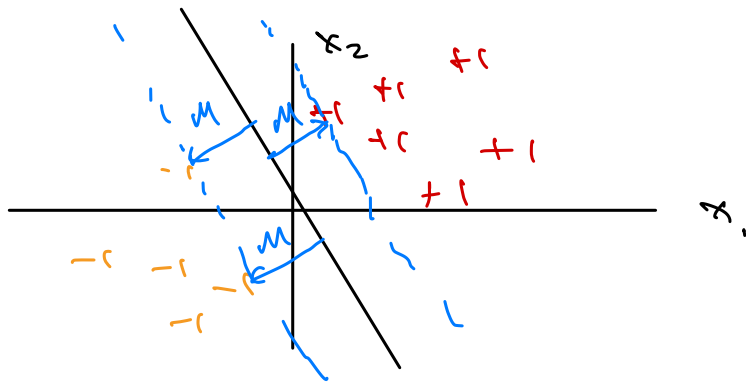


Recall

- Hyperplane is all $\underline{x} \in \mathbb{R}^p \rightarrow \beta_0 + \beta_1 x_1 + \dots + \beta_p x_p = 0$
- Note hyperplane is defined by $\beta_0, \beta_1, \dots, \beta_p$, & splits space into halves
- Data $y_i \in \{-1, +1\} \rightarrow \underline{x}_i = (x_{i1}, x_{i2}, \dots, x_{ip})^T$
 $i = 1, \dots, n$

Suppose linearly separable



Note The maximum margin hyperplane is calculated by:

- Maximize M

$$\beta_0, \beta_1, \dots, \beta_p$$

- Subject to

$$\|\beta\|_2^2 = \sum_{j=1}^p \beta_j^2 = 1 \quad \text{and} \quad y_i (\beta_0 + \beta^T x_i) \geq M \quad \forall i=1, \dots, n$$

$$\beta = \begin{pmatrix} \beta_1 \\ \vdots \\ \beta_p \end{pmatrix} \rightarrow \text{no } \beta_0$$

$p \times 1$



Remark

Restricting β to be length 1 is not a restriction b/c
any \underline{x} that satisfies

$$\beta_0 + \beta^T \underline{x} = 0$$

also satisfies

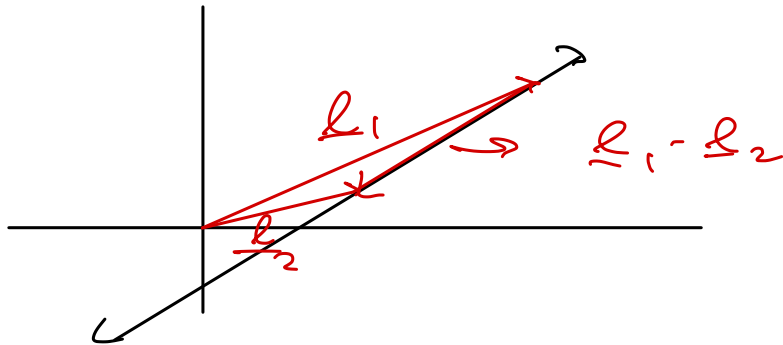
$$k\beta_0 + k\beta^T \underline{x} = 0 \quad \text{for } k \neq 0$$

Claim $y_i (\beta_0 + \beta^T x_i)$ is the distance from x_i to the hyperplane. Thus, M is the margin from the nearest data points to the plane.

Argument in 2 steps:

① Note β is normal to the hyperplane.

Any vector parallel to the hyperplane can be written as a difference of two vectors on the plane



Need to show
 $\beta^T (l_1 - l_2) = 0$

$$\beta^T \underline{x}_1 - \beta^T \underline{x}_2$$

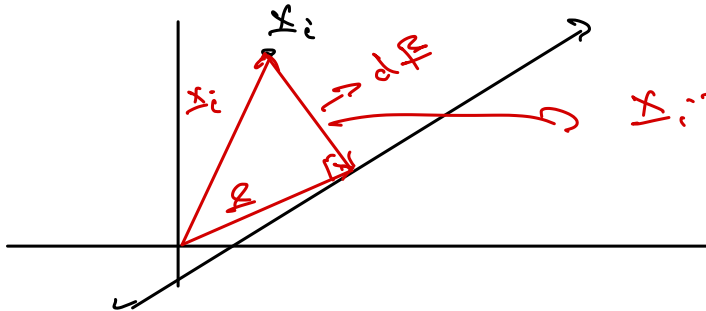
$$[\underline{x}_1, \underline{x}_2 \text{ on plane} \Rightarrow \beta_0 + \beta^T \underline{x}_i = 0 \\ \beta^T \underline{x}_i = -\beta_0]$$

$$= -\beta_0 - (-\beta_0) = 0.$$

(2) Due to the algorithm, $\|\beta\|_2 = 1$ is a unit vector, so the vector from a candidate point \underline{x}_i to the plane can be written

$$\underline{x}_i - \underline{z} = d\beta$$

where d is the signed distance to the plane.



$\underline{x}_i - \underline{z}$ = orthogonal to plane & of length $d = d\beta$

multiply by β^T :

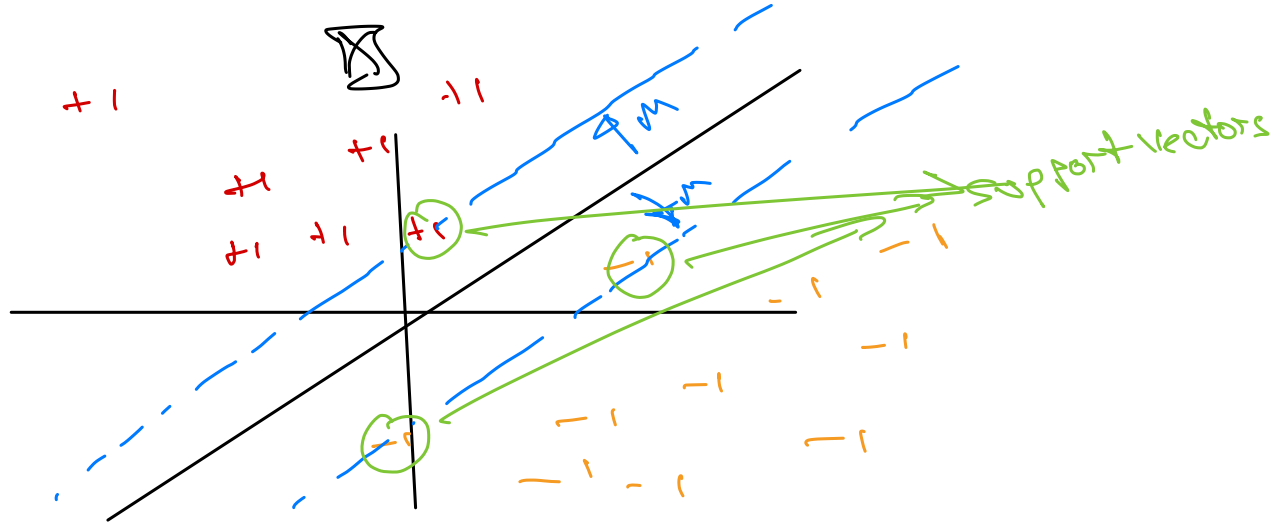
$$\beta^T x_i - \beta^T p_0 = d \beta^T \beta$$

$$\beta^T \beta = \|\beta\|_2^2 = 1$$

$$\Leftrightarrow \beta^T x_i - (-p_0) = d$$

$$\Leftrightarrow p_0 + \beta^T x_i = d$$

$$\Leftrightarrow y_i (p_0 + \beta^T x_i) = y_i d = \text{unsigned distance.}$$



Algorithm

- Maximize M
 β_0, β
 - subject to $\|\beta\|_2 = 1$ + $y_i(\beta_0 + \beta^T x_i) \geq M$ $\forall i$
-

Try to remove $\|\beta\|_2 = 1$ condition

Let β_d be parallel to β but of length d , so $\beta_d = d\beta$

$$y_i(\beta_0 + \beta^T x_i) \geq M$$

$$y_i(\beta_0 + \beta_d^T x_i / d) \geq M$$

$$y_i(d\beta_0 + \beta_d^T x_i) \geq dM$$

The plane defined by $\beta_0 + \beta$ is the same as that defined by $d\beta_0 + \beta_d$, choose $d = \frac{1}{M}$

The algorithm can be rewritten as

$$\left[\begin{array}{ll} \min_{\gamma_0, \beta} & \frac{1}{2} \|\beta\|_2^2 \\ \text{subject to} & \gamma_i (\gamma_0 + \beta^\top x_i) \geq 1 \quad \forall i \end{array} \right.$$

Can convert this to a Lagrangian formulation

[go back $\beta_0 + \beta$ notation]

$$\min_{\beta_0, \beta} \quad \frac{1}{2} \|\beta\|_2^2 - \sum_{j=1}^n \alpha_j (\gamma_j (\beta_0 + \beta^\top x_j) - 1)$$

where $\alpha_1, \dots, \alpha_n$ are Lagrange multipliers

$$\frac{d}{d\beta} \Big| \Rightarrow \beta - \sum_{i=1}^n \gamma_i \alpha_i x_i \stackrel{\text{set}}{=} 0$$

$$\Rightarrow \underline{f} = \sum_{i=1}^n \gamma_i \alpha_i \underline{x}_i$$

$$\frac{d}{d\beta_0} \Rightarrow \sum_{i=1}^n \alpha_i \gamma_i = 0$$

plug back in

$$\sum_{i=1}^n \alpha_i - \frac{1}{2} \sum_{i=1}^n \sum_{j=1}^n \alpha_i \alpha_j \gamma_i \gamma_j \underline{x}_i^T \underline{x}_j$$

maximize over $\underline{\alpha}$ subject to $\alpha_i \geq 0$,

At solution it turns out that:

- $\alpha_i = 0$ happens if $\gamma_i (\beta_0 + \underline{f}^T \underline{x}_i) > 1$
- $\alpha_i > 0$ " " $\gamma_i (\beta_0 + \underline{f}^T \underline{x}_i) = 1$

→ these points are the support vectors, it is only these points that define the hyperplane,

The classification function is still

$$f(\underline{x}) = \beta_0 + \beta^T \underline{x} = \beta_0 + \sum_{i=1}^n \alpha_i \gamma_i \underline{x}_i^T \underline{x}$$

+ classification rule is still

$$\hat{y} = \begin{cases} +1 & f(\underline{x}) > 0 \\ -1 & f(\underline{x}) < 0 \end{cases}$$

only nonzero terms
are support vectors

