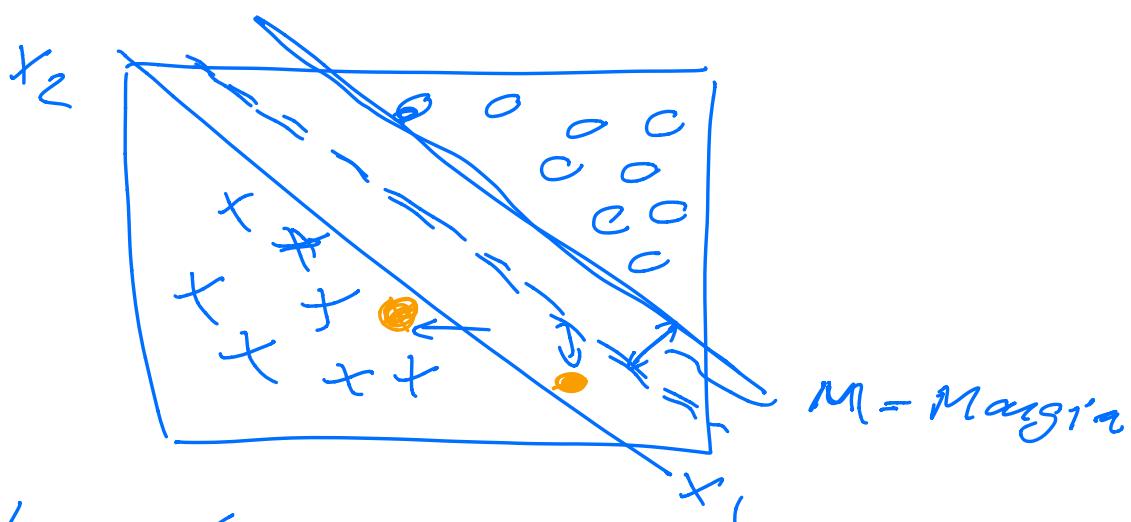
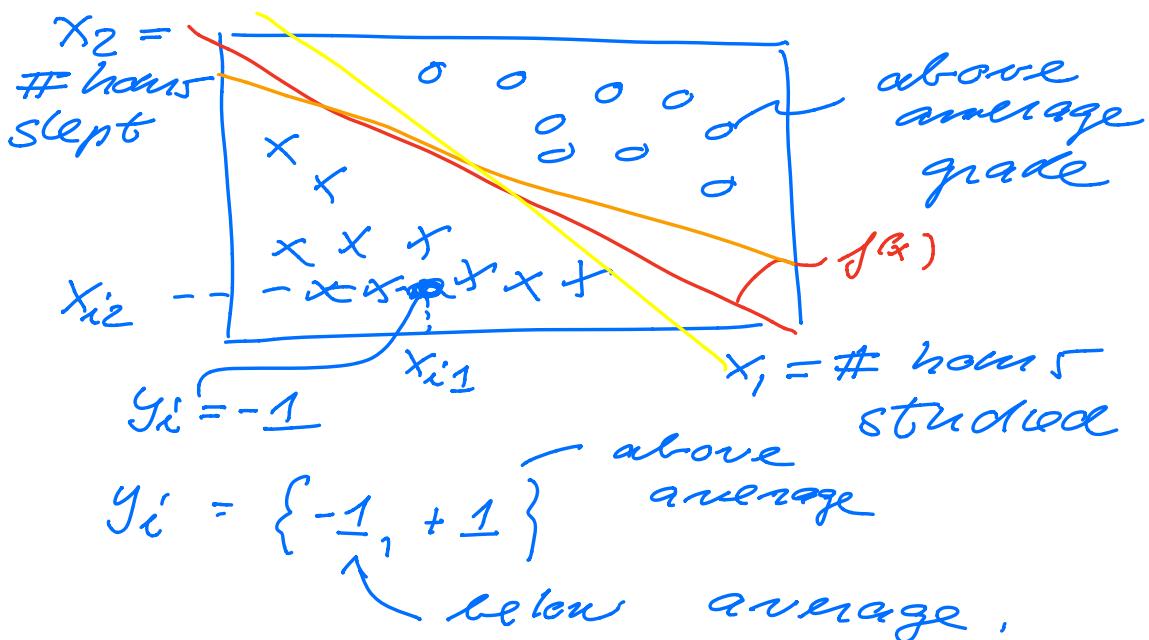


Lecture November 19

- Linear Classification.



in 2-dim a hyperplane is
a straight line

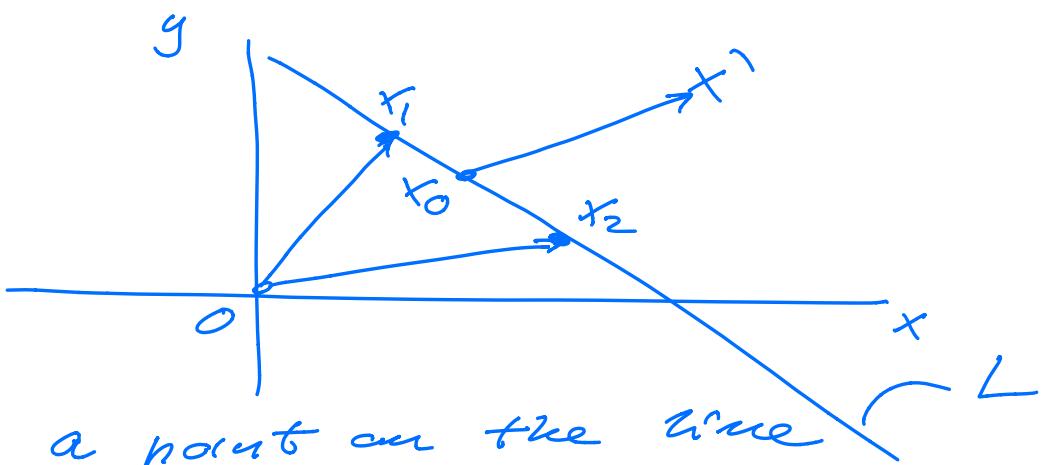
$$b + w_1 x_1 + w_2 x_2 = 0$$

$$\begin{aligned} \mathbf{x}^T &= [x_1, x_2] & \mathbf{w}^T &= [w_1, w_2] \\ (\beta_0 + \beta_1 x_1 + \beta_2 x_2) \end{aligned}$$

$$f(\mathbf{x}) = b + \mathbf{x}^T \mathbf{w}$$

- in 3-dim a hyperplane is a 2-dim plane,
- in p-dim dimension space, $p-1$ hyperplanes,

vector algebra



$$X : b + w_1 x_1 + w_2 x_2 = 0$$

$$\mathbf{x} = [x_1 \ x_2]$$

$$f(\mathbf{x}) = \mathbf{x}^T \mathbf{w} + b = 0$$

any two points defined by
the vectors x_1 & x_2 in L

$$\underline{(x_1 - x_2) w = 0}$$

vector on L

For this to be zero, w must be perpendicular to L

Define a vector $w^* = \frac{w}{\|w\|}$

$$\|w\| = w^T w$$

For any point x_0 in L

$$\underline{x_0^T w = -b}$$

The signed distance of any point x' to L is given

$$(x'^T - \underline{x_0^T}) w = \frac{1}{\|w\|} \frac{(b + x^T w)}{f(x)}$$

$$x' \rightarrow x$$

$$(x^T - x_0^T) w = \frac{f(x)}{\|f(x)\|}$$

$$f(x) = b + x^T w$$

$f(x)$ is proportional to the

Signed distance from
any point x to the
hyperplane defined by

$$f(x) = 0$$

$$f(x) = x^T w + b$$

our model, w_1 and w_2
are unknown parameters.

Observations $y_i = \{-1, 1\}$

$$f(x_i) = x_i^T w + b = \{-1, 1\}$$

$$= [x_{i1} \ x_{i2}] \begin{bmatrix} w_1 \\ w_2 \end{bmatrix} + b$$

One possible approach

want to minimize

$$C(b, w) = - \sum_{\substack{i \in \text{MIS} \\ \text{class}}} \frac{y_i (x_i^T w + b)}{\|x_i\|}$$

unknown b and w

$$\frac{\partial C}{\partial w} = - \sum_{\substack{i \in \text{MIS} \\ \text{class}}} y_i x_i$$

$$\frac{\partial C}{\partial w} = - \sum_{\substack{i \in M \\ \text{class}}} y_i'$$

Problem

- when data are well separated in classes can have many solutions that are more pernicious,
- $\begin{bmatrix} b \\ w \end{bmatrix} \leftarrow \begin{bmatrix} b \\ w \end{bmatrix} + \gamma \begin{bmatrix} -\sum_i y_i' \\ -\sum_i y_i' x_i \end{bmatrix}$
learning rate
- May need many GD iterations. Need to tune γ ,
- with no separable data, no convergence,
instead of $f(x) = 0$, let us try

$$\max_{b, w} M \quad \text{subject to}$$

$$y_i(x_i^T w + b) \geq M \quad i = 0, 1, \dots, n-1$$

all points are at a signed distance from the decision boundary defined by b and w and M (largest M)

$$y_i \frac{(x_i^T w + b)}{\|w\|} \geq M$$

$$y_i(x_i^T w + b) \geq M \|w\|$$

$$\text{we can scale } M = \frac{1}{\|w\|}$$

$$y_i(x_i^T w + b) \geq 1$$

- Maximize M or minimize $\|w\| = w^T w \Rightarrow \frac{1}{2} w^T w$
- subject to (s.t.)

$$\underline{y_i(x_i^T w + b) \geq 1} \text{ for}$$

λ_i

Lagrange formulation

$$\mathcal{L}(w, b, x) = \frac{1}{2} w^T w$$

$$- \sum_{i=0}^{n-1} \lambda_i (y_i(x_i^T w + b) - 1)$$

Lagrangian
multiplier

$$\frac{\partial \mathcal{L}}{\partial w} = w - \sum_{i=0}^{n-1} \lambda_i y_i x_i = 0$$

$$\frac{\partial \mathcal{L}}{\partial b} = - \sum_{i=0}^{n-1} \lambda_i y_i = 0$$

$$w = \sum_{i=0}^{n-1} \lambda_i y_i x_i \quad \sum_{i=0}^{n-1} \lambda_i y_i = 0$$

$$\mathcal{L}(\lambda) = \underbrace{\sum_{i=0}^{n-1} \lambda_i}_{} - \frac{1}{2} \sum_{i,j} \lambda_i \lambda_j y_i y_j \\ \times x_i^T x_j \in [x_j, x_{j+1}]$$

subject to

$$\sum \lambda_i y_i = 0$$

$$\lambda_i \geq 0$$

in addition we must have

$$\lambda_i [y_i(x_i^T w + b) - 1] = 0$$

$$H_i'$$

minimize w.r.t to λ
subject to the conditions

$$\min_{\lambda} \frac{1}{2} \lambda^T \begin{bmatrix} y_1 y_1 x_1^T x_1 & y_1 y_2 x_1^T x_2 & \dots \\ y_2 y_1 x_2^T x_1 & x_2^T x_2 y_2 y_2 & \dots \\ \vdots & \vdots & \ddots \\ y_m y_1 x_m^T x_1 & \dots & \ddots \end{bmatrix} \lambda + (-1) \lambda$$

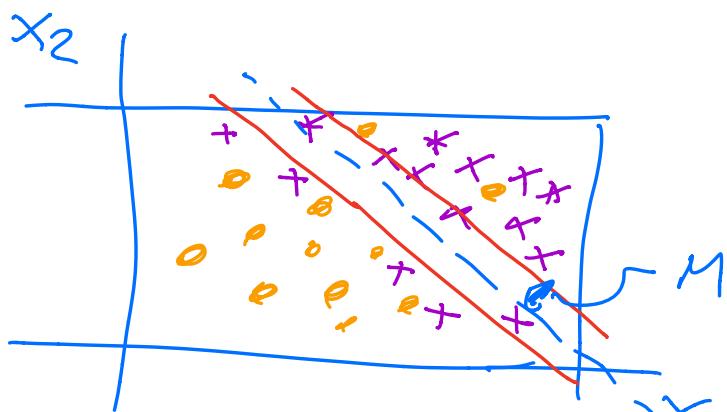
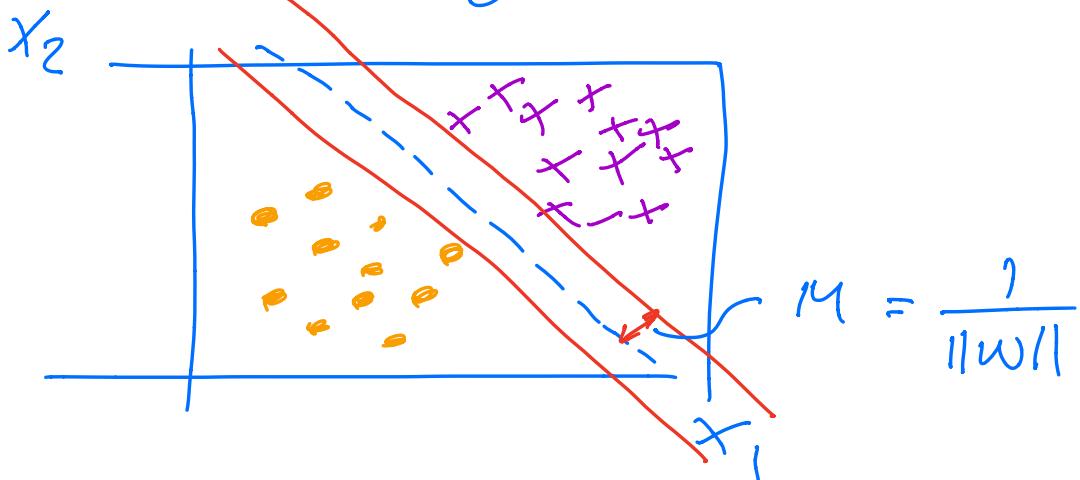
$$\text{subject to } g^T \lambda = 0$$

$$\text{and } 0 \leq \lambda_i \leq \infty$$

The points defined by
 x_i are called the support
vectors,

Convex optimization
problem. Python [CVXOPT]

Hard margin case M



We will still maximize

M, but will allow for some points which are on the wrong side of the margin

Introduce new hyperparameter

$$\xi^T = \xi_+ - \xi_- \leq 1$$

$\rightarrow \text{--- } \dots \text{ --- } \text{ --- }$

$$\frac{1}{\|w\|} y_i (x_i^T w + b) \geq M(1 - \xi_i)$$

$\forall i \quad \xi_i \geq 0 \quad \sum_{i=0}^{n-1} \xi_i \leq \text{constant}$

$$M = \frac{1}{\|w\|}$$

$$y_i (x_i^T w + b) \geq 1 - \xi_i$$

↑
Slack
variable

Need to minimize the

$$\frac{1}{2} w^T w + C \sum_{i=0}^{n-1} \xi_i$$

$S, T,$

$$y_i (x_i^T w + b) \geq 1 - \xi_i \quad \forall i$$

Lagrange formulation

$$\begin{aligned} \mathcal{L}(w, b, \lambda, \xi, x) &= \frac{1}{2} w^T w \\ &+ C \sum_{i=0}^{n-1} \xi_i - \sum_{i=0}^{n-1} (\lambda_i [y_i (x_i^T w + b) \\ &\quad - (1 - \xi_i)]) \end{aligned}$$

$$-\sum_{i=0}^n \gamma_i \underline{\underline{\gamma}}_i$$

$$\frac{\partial L}{\partial w} = w - \sum_i \lambda_i y_i x_i = 0$$

$$\frac{\partial L}{\partial b} = - \sum_i \lambda_i y_i = 0$$

$$\frac{\partial L}{\partial \gamma_i} = C - \lambda_i - \gamma_i = 0 \Rightarrow$$
$$\gamma_i = C - \lambda_i$$

\nearrow

$C = \text{hyperparameter}$.