

Lecture 3:

(2) Supervised learning^o has 3 ingredients

(1) Training Data $D = \langle X, Y \rangle$

$$X = \begin{bmatrix} \vec{x}_1 \\ \vec{x}_2 \\ \vdots \\ \vec{x}_n \end{bmatrix}, Y = \begin{bmatrix} y_1 \\ y_2 \\ \vdots \\ y_n \end{bmatrix}, \vec{x}_i \in \mathcal{X}, y_i \in \mathcal{Y}$$

(1) $Y = t(\vec{x}) = f(\vec{x}) \rightarrow \delta = h^*(\vec{x}) + \underbrace{\epsilon}_{\text{epsilon/error}}$

δ : error due to ignorance.

Misspecification error.

Note: the space of F is very large.

We need to constrain it.

(2) \mathcal{H} : a candidate set of function

③ A , an algorithm which takes in data

ID and set H and produces a model g .

$$g = A(ID, H)$$

Question: is $f: H$ generally? NO

However, there is $h^* \in H$ which is the closest possible model to f .

(function)

Just because $h^* \in H$ does not mean A will locate it. A will not be perfect and the values of ϵ will confuse A .

Arina

Rameesa

Thus $g \neq h^*$ g is the best A can do.

$$y = g(x) + (h^*(x) - g(x)) + (f(x) - h^*(x)) + (1)$$

estimation error.

$$y = g(x) + \underbrace{(h^*(\vec{x}) - g(\vec{x}))}_{\text{model}} + \underbrace{(f(\vec{x}) - h^*(\vec{x}))}_{\delta} + \underbrace{(t(\vec{z}) - f(\vec{x}))}_{\epsilon}$$

ϵ "residuals"

$\hat{y} \leftarrow$ Prediction of y in setting \vec{x}

$$\hat{y} = g(\vec{x})$$

$$e = y - \hat{y} \text{ residual if } \vec{x} \in D$$

Otherwise they're unknown

How to reduce errors:

- (1) δ ignorance error can be reduced by measuring more x_j 's (features) of the units that contain

information about \vec{z}

* (2) Misspecification error can be reduced by ...

expand H to include more complicated functions

(3) Estimation error can be reduced by increase sample size n .

$$SE[\hat{\theta}] = \frac{\sigma}{\sqrt{n}}$$

$$y = \begin{cases} 0, 1 \end{cases}$$

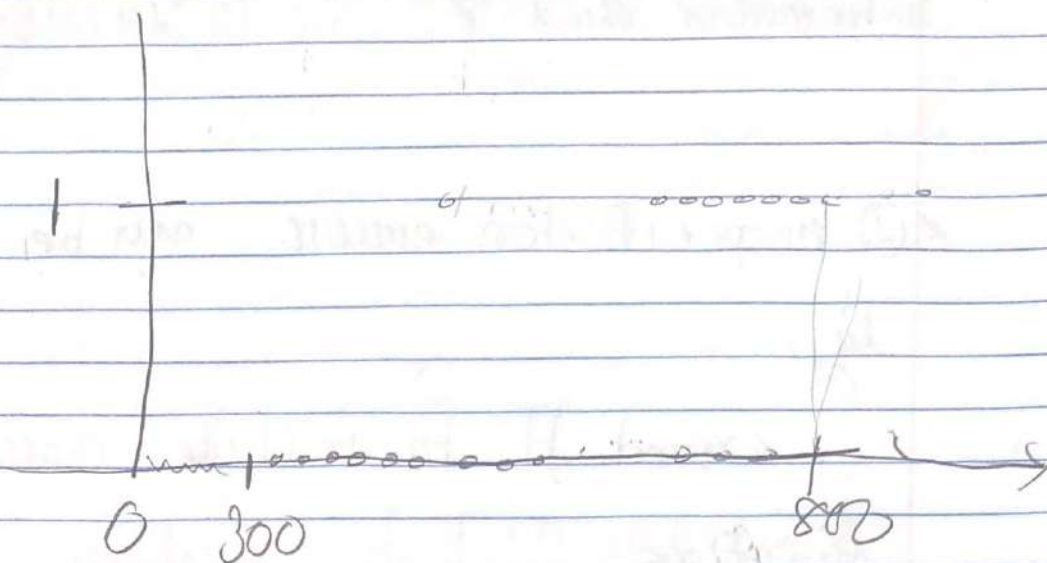
↑ ↑
did not pay back paid back mortgage

(1)
 $n = 100$
 $\mathcal{D} = \langle X, Y \rangle =$ $y \in \{0, 1\}$ (classification)

$$= \left(\begin{bmatrix} 810 \\ 390 \\ 750 \\ \vdots \end{bmatrix}, \begin{bmatrix} 1 \\ 0 \\ 1 \\ \vdots \end{bmatrix} \right)$$

$p = 1$, X_1 is credit score

$$x = [300, 850]$$



threshold model.

$$I = \left\{ \begin{array}{l} 1 \\ 0 \end{array} \right\} \quad x \geq \theta, \quad \theta \in \Theta$$

indicator
function

Parameter

Sometimes
denoted

β, w, \dots , others.

capital θ

in the greek
alphabet

$$I_A(w) = \begin{cases} 1 & \text{if } w \in A \\ 0 & \text{if } w \notin A \end{cases}$$



e.g

$$g(x) = 1_{x \geq 515.3}$$

$$g(x) = 1_{x \geq 407.9}$$

(3) Algorithm A .

Recall,

$$g(x) = A(Z, D)$$

Let Define

"Misclassification error" -

$$ME := \frac{1}{n} \sum_{i=1}^n \mathbb{I}_{g(x_i) \neq y_i}$$

$$= \frac{1}{n} \sum_{i=1}^n |e_i| \quad \rightarrow \text{residuals.}$$

accuracy:

$$ACC := 1 - ME$$

A. minimize ME over

$\theta \in$ unique x 's

Sum absolute error (SAE)

$$\frac{1}{n} \sum_{i=1}^n |e_i| = \frac{1}{n} \sum_{i=1}^n |y_i - \hat{y}_i|$$

mean absolute error (MAE)

minimizing
be (cost)

$$= \frac{1}{n} \sum_{i=1}^n e_i^2 \quad \text{Sum squared error (SSE)}$$

mean squared error

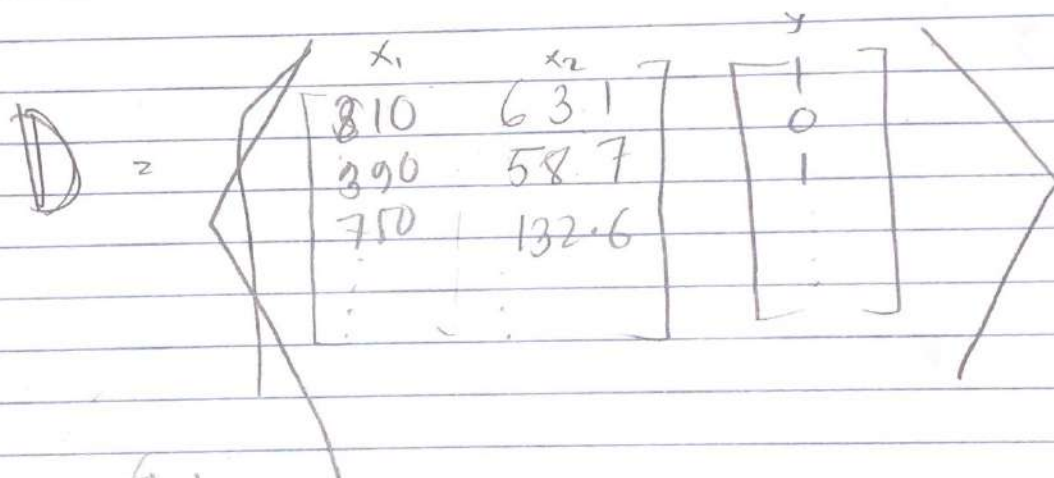
(MSE)

Next model

$$Y \in \mathcal{Y} = \{0, 1\}$$

X_1 : credit score $\in [300, 850]$

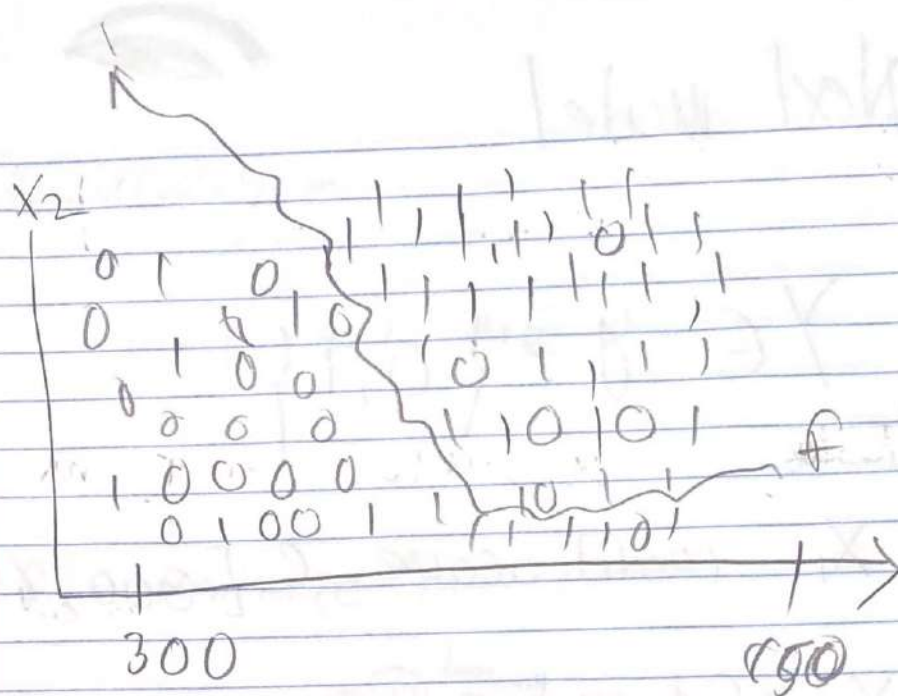
X_2 : Salary^{thousand} $\in \mathbb{R}$



Picture source

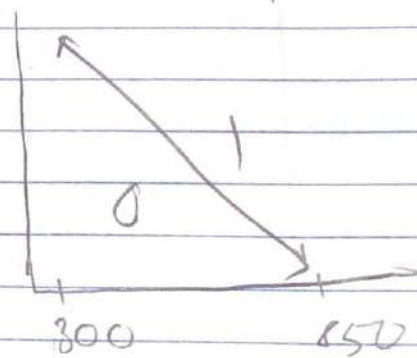
$$\mathcal{H} = \left\{ \underset{\text{[Not good]}}{1_{X_1 > \theta_1 \wedge X_2 \geq \theta_2}} \mid \begin{bmatrix} \theta_1 \\ \theta_2 \end{bmatrix} \in (\mathbb{R})^2 \right\}$$

e.g. $g(x) = 1_{X_1 > 650 \wedge X_2 \geq 100}$



$$H = \left\{ \mathbb{I}_{x_2 \geq a + bx_1} : a, b \in \mathbb{R} \right\}$$

Parameter Space has
dimension 2.



Also known as

2 degrees of freedom.

$$H = \left\{ x_2 > a + bx_1 \mid a, b \in \mathbb{R} \right\}$$

~~$x_2 > a$~~

$$x_2 > a + bx_1 = \underbrace{-a}_{w_0} + \underbrace{-bx_1}_{w_1} + \underbrace{(1)x_2}_{w_2} \geq 0$$

"weights"

$$\Rightarrow w_0 + w_1 x_1 + w_2 x_2 \geq 0$$

↑ credit score, salary

$$P+1 = 3 \text{ (}\# \text{ of cols)}$$

↳ have the bias term.

$$X = \begin{bmatrix} \vdots & x \\ 1 & \end{bmatrix}$$

$$\vec{X} = [1 \ x_1 \ x_2]$$

2nd notation

Redefine the matrix

X by appending

a col of 1's on the left

$$H = \{ \vec{w} \cdot \vec{x} \geq 0 : \vec{w} \in \mathbb{R}^3 \}$$

Thin is an Overparameterized

Model

Each line has infinite \vec{w} 's

that speeding it.

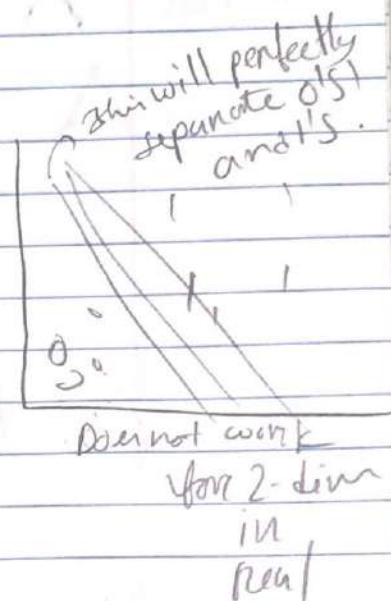
Need algorithm A

$$J \cong A(\mathbb{D}, \#)$$

How many intercepts
1 in slopes

Assume 0's and 1's are
linearly separable.

$\exists \vec{w} \cdot s.t. g(\vec{X})$ has no
error.



Perceptron Learning Algorithm (1957).

- (1) Initialize $\vec{w}^{t=0} = \vec{0}$ or random,
compute \hat{y} .
- (2) Let - -

$$w_0^{t+1} = w_0^{t=0} + \underbrace{(y_i - \hat{y}_i)}_{e_i} (1)$$

$$w_1^{t+1} = w_1^{t=0} + \underbrace{(y_i - \hat{y}_i)}_{e_i} x_{i,1}$$

$$w_2^{t+1} = w_2^{t=0} + (y_i - \hat{y}_i) x_{i,2}$$

$$w_p^{t+1} = w_p^{t=0} + (y_i - \hat{y}_i) x_{i,p}$$

Recompute \hat{y}

$$X_z = \begin{bmatrix} 1 & x_{1,1} & x_{1,2} & \dots & x_{1,p} \\ \vdots & x_{2,1} & x_{2,2} & \dots & x_{2,p} \\ \vdots & \vdots & \vdots & \dots & \vdots \\ 1 & x_{n,1} & x_{n,2} & \dots & x_{n,p} \end{bmatrix}$$

③ Repeat 2 for $i=1 \dots n$

④ Repeat steps 2,3 until no errors

Perceptron is proven to converge if linear separability assumption is true.