

# Supervised learning (3 ingredients)

$$y = t(z) = F(\vec{x}) + j$$

$j$ : error due to ignorance

## 1) Training Data

$$D = \langle X, Y \rangle$$

$$X = \begin{bmatrix} x_1 \\ \vdots \\ x_n \end{bmatrix}, Y = \begin{bmatrix} y_1 \\ \vdots \\ y_n \end{bmatrix}, \vec{x}_i \in X, y_i \in Y$$

## 2) $\mathcal{H}$ : a candidate set of functions

## 3) $A$ : an algorithm which takes in data $D$ and set $\mathcal{H}$ and produces a model

$$g = A(D, \mathcal{H})$$

Is  $F \in \mathcal{H}$  generally? No

however  $h^* \in \mathcal{H}$  which is the closest possible model to  $F$  (function)

$$y = t(z) = F(x) + j = h^*(x) + \epsilon$$

$$h^*(\vec{x}) + \underbrace{(F(\vec{x}) - h^*(\vec{x}))}_{\text{mispecification error (model too simple for a complex zero)}} + \underbrace{(t(\vec{z}) - F(\vec{x}))}_{\epsilon \text{ "epsilon error"}}$$

mispecification error (model too simple for a complex zero)

$\epsilon$  epsilon error

$$y = g(x) + \underbrace{(h^*(x) - g(x))}_{\text{estimation error}} + \underbrace{(F(x) - h(x))}_{\epsilon} + \underbrace{(t(z) - F(x))}_{j}$$

model

$\epsilon$  residual

## How to reduce errors:

1)  $j$ , ignorance error can be reduced by: measuring more  $x_j$ 's (features) of the units that contain information

2) Misspecification error:

expand set  $\mathcal{H}$  to make more complicated functions  $\} \rightarrow$  machine learning

3) Estimation error can be reduced by: increase sample size

## example of supervised learning:

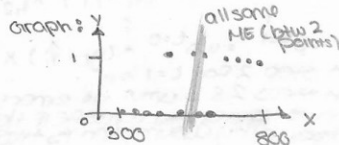
$$1) y = 20.113$$

# of variables  $\uparrow$  paid back mortgage  
 $\hookrightarrow p=1$   $x_1$ : credit score

$$D = \langle X, Y \rangle = \left\langle \begin{bmatrix} 810 \\ 370 \\ 750 \\ \vdots \end{bmatrix}, \begin{bmatrix} 0 \\ 1 \\ \vdots \end{bmatrix} \right\rangle$$

$$X = [300, 850]$$

$$n = 100$$



$$2) \mathcal{H} = \left\{ \begin{array}{l} 1 \text{ if } x \geq \theta \\ 0 \text{ if } x < \theta \end{array} \right\}$$

threshold model  $\uparrow$  indicator function  $\uparrow$  parameter  $\theta$  (also denoted  $\beta_1, w, \dots$ )

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

## 3) Algorithm A

$$g(x) = A(\mathcal{H}, D)$$

A estimates  $\theta$  (parameter)

"misspecification error"

$$ME := \frac{1}{n} \sum_{i=1}^n |g(x_i) - y_i| \neq y_i = \frac{1}{n} \sum_{i=1}^n |e_i|$$

$$\text{accuracy} := 1 - ME$$

A: minimize ME over  $\theta \in \Theta$

$\theta \in \Theta$  unique  $x$ 's

supervised absolute error (SAE)

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

mean absolute error (MAE)

sum squared error

mean squared error

\* Just because  $h^* \in \mathcal{H}$  does NOT mean  $A$  will locate it.  $A$  will not be perfect and the values of  $\epsilon$  will confuse  $A$ . then  $g \neq h^*$ ;  $g$  is the best  $A$  can do

$\hat{y} = g(\vec{x})$   
prediction of  $y$  in setting  $\vec{x}$

$e = y - \hat{y}$   
residual if  $\vec{x} \in D$  otherwise they're unknown

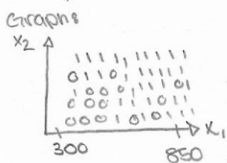
ex 2)

$$y \in \{0, 1\}$$

$x_1$ : credit score  $\in [300, 850]$

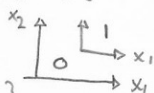
$x_2$ : salary  $\in \mathbb{R}$   
(\$1000s)

$$D = \left\langle \begin{bmatrix} 810 & 63.1 \\ 770 & 58.7 \\ 750 & 137.6 \\ \vdots & \vdots \end{bmatrix}, \begin{bmatrix} 1 \\ 0 \\ 1 \end{bmatrix} \right\rangle$$



$$H = \{ \mathbb{1}_{x_1 \geq \theta_1} \& \mathbb{1}_{x_2 \geq \theta_2} : [\theta_1, \theta_2] \in \Theta \}$$

$$f(x) = \mathbb{1}_{x_1 \geq 650 \& x_2 \geq 100}$$



misclassification error

$$H = \{ \mathbb{1}_{x_0 \geq a + bx_1} : a, b \in \mathbb{R} \}$$

parameter space has dimension 2 (degrees of freedom)



$$x_i \geq a + bx_1 \Rightarrow \underbrace{-a}_{w_0} + \underbrace{-b}_{w_1} x_1 + \underbrace{(1)}_{w_2} x_2 \geq 0$$

bias weights

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

$$X = \begin{bmatrix} 1 & x_1 & x_2 \end{bmatrix}$$

redesign the matrix by appending a column of 1s on the left

$$\tilde{X} = [1, x_1, x_2]$$

$$H = \{ \mathbb{1}_{\tilde{w} \cdot \tilde{X} \geq 0} : \tilde{w} \in \mathbb{R}^3 \}$$

this is overparameterized model each line has infinite  $\tilde{w}$ s that specify it.

$P+1 = 3$  (# of cols in  $X$ )  
default scores/salary for the bias term

need algorithm A:

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

assume 0's and 1's are linearly separable



Perceptron Learning Algorithm (1957)

1) initialize  $\tilde{w}^t = 0 = \vec{0}$  or random (compute  $\hat{y}$ )

2) for  $j=0, 1, 2, \dots, P$  let:

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

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

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

$$X = \begin{bmatrix} x_{11} & x_{12} & \dots & x_{1P} \\ x_{21} & x_{22} & \dots & x_{2P} \\ \vdots & \vdots & \ddots & \vdots \\ x_{n1} & x_{n2} & \dots & x_{nP} \end{bmatrix}$$

\* update weights until right decision



$$w_0^{t+1} = w_0^t + (y_i - \hat{y}_i) x_{i,0}$$

3) Repeat step 2 for  $i=1 \rightarrow n$

4) Repeat steps 2 & 3 until no errors

\* Perceptron is proven to converge if the linear separability assumption is true

recompute  $\hat{y}$