# Math 390.4 Lec# 3/7/18

Supervised Learning $\mathbb{D}, \mathcal{H}, \mathcal{A}$  then  $g = \mathcal{A}(\mathbb{D}, \mathcal{H})$

How to see a prediction for $\vec{x} \in \mathbb{D}$?  must have

$\hat{y}_i := g(\vec{x}_i)$   $i \in \{1, \dots, n\}$   same features as $\vec{x}_i$'s in $\mathbb{D}$, otherwise it is not in the domain of $g$!!

Then the next question is mathematically obvious but conceptually quite a leap!!

How to predict for new person? $\vec{x}^*$ ... let $\hat{y}^* = g(\vec{x}^*)$. Hopefully $\hat{y}^* \approx y^*$ most of the time

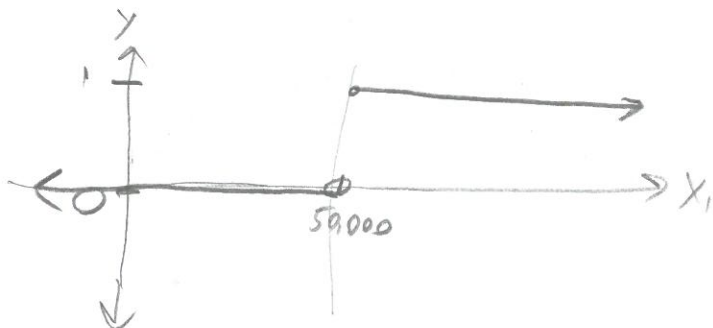real, unobserved, as of $n$

Back to our problem...

$y \in \{0, 1\}$,  $\vec{x} \in \mathcal{X}$  where $p = 3$. To make it even simpler, let's say we only have $x_1$, salary.

let $\mathcal{H} = \left\{ \mathbb{1}_{x \geq x_T} \right\}$

Indicator function is   $\mathbb{1}_{x \geq x_T} = \begin{cases} 1 & \text{if } x \geq x_T \\ 0 & \text{if not} \end{cases}$

What would a $g$ look like?



Threshold Model
Is it a good model??  [NO]

Parameter notation

$\theta, \beta, w,$ any other symbols used (autonomously)

What do all possible choices of $\mathcal{H}$ look like? Elements in $\mathcal{H}$ are indexed by $x_T$. Every single value of $x_T$, the threshold creates a new model. $x_T$ is called a "parameter".

What does $A$ do? It selects the best candidate $g \in \mathcal{H}$.
That is, is selects the "best" $x^*$. How??

Consider the following... there are $n$ $x_i$'s. They $x_T = x_1, x_T = x_2, \ldots x_T = x_n$. pick the "best". How to know the best?

$\hat{y}_i := h(\vec{x}_i)$

$$SSE(h) := \sum_{i=1}^{n} \left( \hat{y}_i - y_i \right)^2$$

↑ sum of squared error

i.e. the $\#$ of times the hypothesis function disagrees with the response (this is to be small)

Next to determine a way of measuring error. This will be a bit diff when we get to start models

$$SAE = \sum_{i=1}^{n} |\hat{y}_i - y_i| = \sum \mathbb{1}_{\hat{y} \neq y_i}$$

# of misclassifications

$$g = \underset{h \in \mathcal{H}}{argmin} \left\{ SSE(h) \right\} \iff x$$

MSE, MAE, misclassification rate. Any monotonic increasing function

But we never used $x_2, x_3 \ldots$ wouldn't we do better? Yes... SSE always gets smaller as you increase $p$.

Why don't we try a more complicated model?

How about $x_1, x_2, x_3$.

we can do a model like

$$\mathcal{H} := \left\{ \mathbb{1}_{x_1 \geq x_1^*} \mathbb{1}_{x_2 \geq x_2^*} \mathbb{1}_{x_3 \geq x_3^*} : \langle x_1^*, x_2^*, x_3^* \rangle \in \mathcal{X} \right\}$$

Why would this be bad?



You are creating a very inflexible model!
Why not a linear model? Imagine $x_1, x_2$ (both continuous)



$a + bx_1$

$$h(x_1, x_2) = \mathbb{1}_{x_2 > a + bx_1} = \mathbb{1}_{a + bx_1 - x_2 < 0}$$

$$= \mathbb{1}_{-a - bx_1 + x_2 > 0} = \mathbb{1}_{w_0 + w_1 x_1 + w_2 x_2 > 0}$$

$$\mathcal{H} := \left\{ \mathbb{1}_{w_0 + w_1 x_1 + w_2 x_2 > 0} : w_0, w_1, w_2 \in \mathbb{R} \right\}$$

$$= \left\{ \mathbb{1}_{w_0 + \vec{w} \cdot \vec{x} > 0} : w_0 \in \mathbb{R}, \vec{w} \in \mathbb{R}^2 \right\} \quad \text{where} \quad \begin{matrix} \vec{x} = [x_1, x_2] \\ \vec{w} = [w_1, w_2] \end{matrix}$$

if we augment $\vec{x}$ to include an $x_0 = 1 \implies \vec{x} = [1 \; x_1 \; x_2]$

, , , , $\vec{w}$ , , , , , , , , , $w_0 \implies \vec{w} = [w_0 \; w_1 \; w_2]$

$$= \left\{ \mathbb{1}_{\vec{w} \cdot \vec{x} > 0} : \vec{w} \in \mathbb{R}^3 \right\} \quad \text{Note that} \quad \overset{2}{\overset{..}{p}} + 1 = 3$$

Here, a 1 is prepended to the feature vector so the intercept $w_0$ can be part of the dot product. This is purely a convenience for linear models. This is sometimes implied without explicit telling you!

---

# What is $A(\mathcal{H}, \mathcal{D})$? How to pick $g \in \mathcal{H}$?

$$g = \underset{h \in \mathcal{H}}{\text{argmin}} \left\{ SSE(h) \right\} = \mathbb{1}_{\vec{w}^* \cdot \vec{x} > 0}$$

Since $\mathcal{H}$ is parameterized by $\vec{w}$, this is the same problem as:
$$\vec{w}^* = \underset{\vec{w} \in \mathbb{R}^3}{\text{argmin}} \left\{ \sum_{i=1}^{n} \mathbb{1}_{\vec{w} \cdot \vec{x}_i > 0} \neq y_i \right\}$$

- You can't check all $\vec{w} \in \mathbb{R}^3$ like $x^* \in \{x_1, ..., x_s\}$ like last time!

- You can't take the derivative wrt to $\vec{w}$ since the # errors is continuous...

- 1957 "Perceptron Learning Algorithm", an iterative, imperfect algorithm