

Math 309 Lec 3 2/5/19

$$y = g(\hat{x}) + \underbrace{h(\hat{x}) - g(\hat{x})}_{\text{estimation error}} + \underbrace{(\hat{y}) - g(\hat{x})}_{\text{misspecified error}} + \underbrace{t(\hat{x}) - \hat{y}}_{\text{error due to ignore}} \\ \underbrace{\hspace{10em}}_{e \text{ (the "residual")}} \quad \hat{y} = g(\hat{x})$$

$g = A(D, \mathcal{H})$
 $\uparrow \quad \uparrow$
 \hookrightarrow training examples set of candidate models / model space / allowable
 Let's see a full example...

Back to the story about loans

Want to build model g .
 "Binary classification model"
 "classification model"

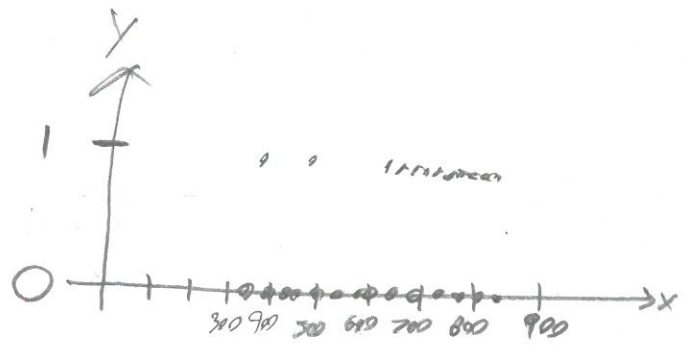
$y = \{0, 1\}$ binary response
 $\uparrow \quad \uparrow$
 will not pay back loan (not ...)
 will pay back loan (credworthy)

Imagine we have one x : credit score. $\mathcal{X} = [300, 850]$

$D = \langle X, y \rangle = \left[\begin{array}{c|c} \begin{matrix} 810 \\ 390 \\ 750 \\ \vdots \end{matrix} & \begin{matrix} 1 \\ 0 \\ 1 \\ \vdots \end{matrix} \end{array} \right] \leftarrow$ historical data for 4 different people.

Conveniently...
 also...

$$= \begin{bmatrix} X & y \end{bmatrix} = \begin{bmatrix} 810 & 1 \\ 390 & 0 \\ 750 & 1 \\ \vdots & \vdots \end{bmatrix}$$

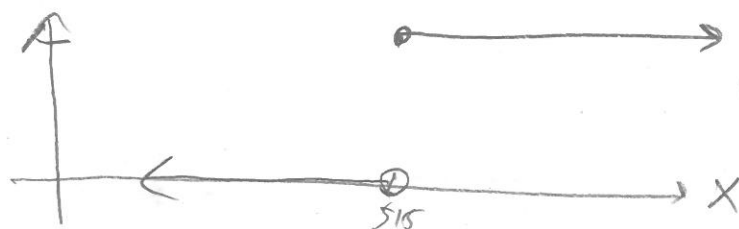


Null Model. No features! $\hat{y} = \text{Model}[x]$, $\mathcal{H} = \mathcal{Y} = \{0, 1\}$ [2]

Let us pick the simplest model space

Boolean function $\mathbb{1}_{x \geq \theta} = \begin{cases} 1 & \text{if } x \geq \theta \\ 0 & \text{otherwise} \end{cases}$

$\mathcal{H} = \left\{ \mathbb{1}_{x \geq \theta} : \theta \in \Theta \right\}$ eg $g(x) = \mathbb{1}_{x \geq 5.151} \in \mathcal{H}$,
 "threshold model" parameter θ, β, w, \dots parameter space
 $g(x) = \mathbb{1}_{x \geq 611.3} \in \mathcal{H}$, etc...



$g(x) = \hat{y}$ these are the predictions this model produces

the algorithm

A produces g . Since g is fully specified by θ , A just needs to select a θ . How about selecting the θ that gives the least prediction errors on \mathcal{D} ?

error is when $y = 0$ but $\hat{y} = 1$ or when $y = 1$ but $\hat{y} = 0$

misclassification error
 $ME := \frac{1}{n} \sum_{i=1}^n \mathbb{1}_{g(\tilde{x}_i) \neq y_i}$
 accuracy

$ACC := \frac{1}{n} \sum_{i=1}^n \mathbb{1}_{g(\tilde{x}_i) = y_i} = 1 - ME$

minimizing error \equiv maximizing accuracy

A: Find θ by considering all $\theta \in \Theta$ and finding the one that

minimizes ME. ME is called an objective function, target function, fitness function, something to be optimized.

A does $\mathcal{O}_g = \argmin_{\theta \in \mathcal{H}} \left\{ \frac{1}{n} \sum_{i=1}^n \mathbb{I}_{x_i \geq \theta \neq y_i} \right\}$ corner but highly! (3)

Not the common $\frac{1}{n}$ will affect the argmin's argmin

errors. Are they residuals? The $y - \hat{y}$ would be the residuals

$$y = g(x) + e$$

$$y = \hat{y} + (y - \hat{y})$$

Can we rewrite A 's objective function to be a function of e_1, \dots, e_n ?

$$ME = \frac{1}{n} \sum_{i=1}^n |y_i - \hat{y}_i| = \frac{1}{n} \sum_{i=1}^n |e_i| = \frac{1}{n} \sum_{i=1}^n e_i^2$$

$\underbrace{\quad}_{MAE} \quad \underbrace{\quad}_{SAE} \quad \underbrace{\quad}_{SSE} \quad \underbrace{\quad}_{MSE}$

Minimizing any of these is equivalent in classification.

A: $g = \argmin_{h \in \mathcal{H}} \{SSE(h)\} \Leftrightarrow \mathcal{O}_g = \argmin_{\theta \in \mathcal{H}} \{SSE(\theta)\}$ okay to run in computer?

We have to specify \mathcal{H} . Ideas? $\mathcal{H} = \mathcal{X}$. All suitable values in \mathcal{D} .

Let's now consider x_1, x_2 two parameters $\dim(\mathcal{H}) = 2$

$$\mathcal{H} = \left\{ \mathbb{I}_{x_1 \geq \theta_1, x_2 \geq \theta_2} : \theta_1, \theta_2 \in \mathcal{H} \right\}$$

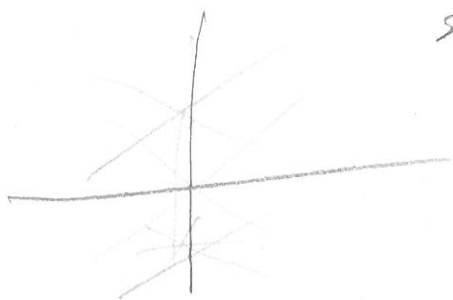


\mathcal{H} is very restrictive! Still could work well... but how about something less restrictive but still very simple? Linear?

$$\mathcal{H} = \{ \mathbb{I}_{x_2 \geq a + b x_1} : a, b \in \mathbb{R} \}$$

↑ ↑

slope and intercept gives enough "degrees of freedom" for any separating line



we need an A to find a & b . we find representation ...

$$x_2 \geq a + b x_1 \Rightarrow -a + b x_1 + x_2 \geq 0 \Rightarrow \begin{matrix} \text{"intercept" or "bias"} & \text{weight } x_1 & \text{weight } x_2 \\ \downarrow & \downarrow & \downarrow \\ w_0 + w_1 x_1 + w_2 x_2 \geq 0 \end{matrix}$$

Now we have 3 parameters, but we only needed 2. This is an "overparameterized" model. For each b , there are ∞ equivalent models... all its comb of \vec{w} i.e. $c\vec{w}$ s.t. $c \in \mathbb{R}$

add a to each observation

$$\text{If we also set } \vec{x} = \begin{bmatrix} 1 & x_1 & x_2 \end{bmatrix} \Rightarrow \vec{w} \cdot \vec{x} \geq 0$$

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

$$p+1 = 3$$

↑
for bias or "intercept" term

Need A to find f . $f = \argmin_{f \in \mathcal{F}} \{SSE(f)\}$ equivalent to

$$\vec{w} = \argmin_{\vec{w} \in \mathbb{R}^3} \left\{ \sum_{i=1}^n \mathbb{1}_{\vec{w} \cdot \vec{x}_i \neq y_i} \right\}$$

Last time we checked only all thresholds in x . We can't check all $\vec{w} \in \mathbb{R}^3$!
Can you solve it independently? No... introduce linear ions difference.

Consider the following algorithm called the perceptron learning algorithm (1957): For the general case of p classes:

- ① Introduce $\vec{w}_{t=0} = \vec{0}$ or random
- ② Compute $\hat{y}_i = \mathbb{1}_{\vec{w}_{t=0} \cdot \vec{x}_i}$
- ③ For $j = 0, 1, \dots, p$ e_i

$$e_i = 1$$

Lead to
make $\vec{w} \cdot \vec{x}$
bigger

if $e_i = 0$...
... smaller

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

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

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

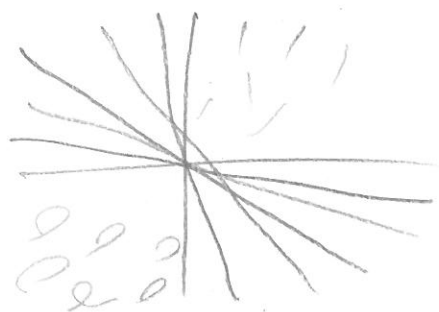
④ Repeat steps 2 & 3 for $i = 1, \dots, n$

⑤ ' ' ' ' 2, 3, 4 until no errors or until a pre-specified max # of iterations



The perceptron is proven to converge if the y 's are linearly separable. If not, it will fail!!

Besides solutions being equal to a multiplicative constant, there are infinitely many solutions that give linear separability e.g.



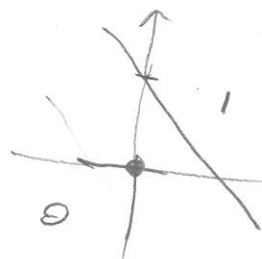
Which one is best? The one that divides the margin evenly.

Weaknesses

- ① Needs $1/n$ steps
- ② Cannot find "best" $1/n$ sep. model

Simple

i	x_1	x_2	y
1	-1	-1	0
2	1	1	1



$$\vec{w}_{t=0} = \begin{pmatrix} 0 \\ 0 \end{pmatrix}$$

$$\hat{y}_1 = \text{sgn}(\vec{w} \cdot \vec{x}) \geq 0 = 1$$

$$t=1, i=1$$

$$w_0 = (0) + (-1)(1) = -1$$

$$w_1 = (0) + (-1)(-1) = 1$$

$$w_2 = (0) + (-1)(1) = -1$$

$$t=1, i=2$$

$$w_0 = (1) + (0)$$

$$w_1 = (1) + (0)$$

$$w_2 = (1) + (0)$$

$$t=2, i=1$$

$$w_0 = (1) + (0)$$

$$w_1 = (1) + (0)$$

$$w_2 = (-1) + (0)$$

$$\hat{y} = 1 - 1 + x_1 - x_2 \geq 0$$

$$\Rightarrow x_2 \geq 1 - x_1$$

$$\hat{y}_2 = 1 + (1) \cdot (1) \geq 0 = 1$$

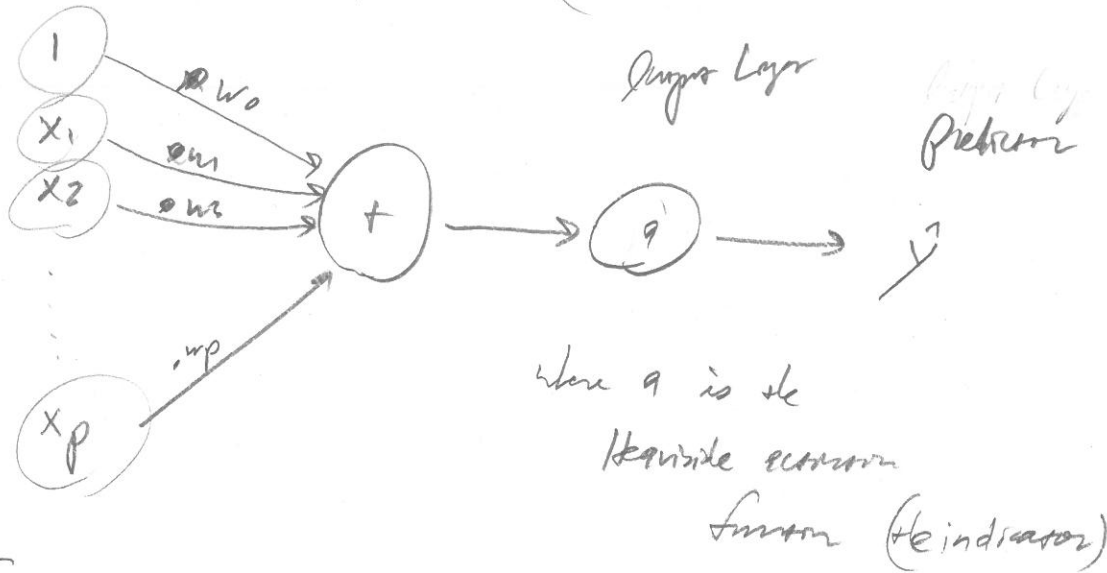
$$\hat{y}_1 = 1 - 1 + (1)(1) + (1)(1) \geq 0$$

$$= 0$$



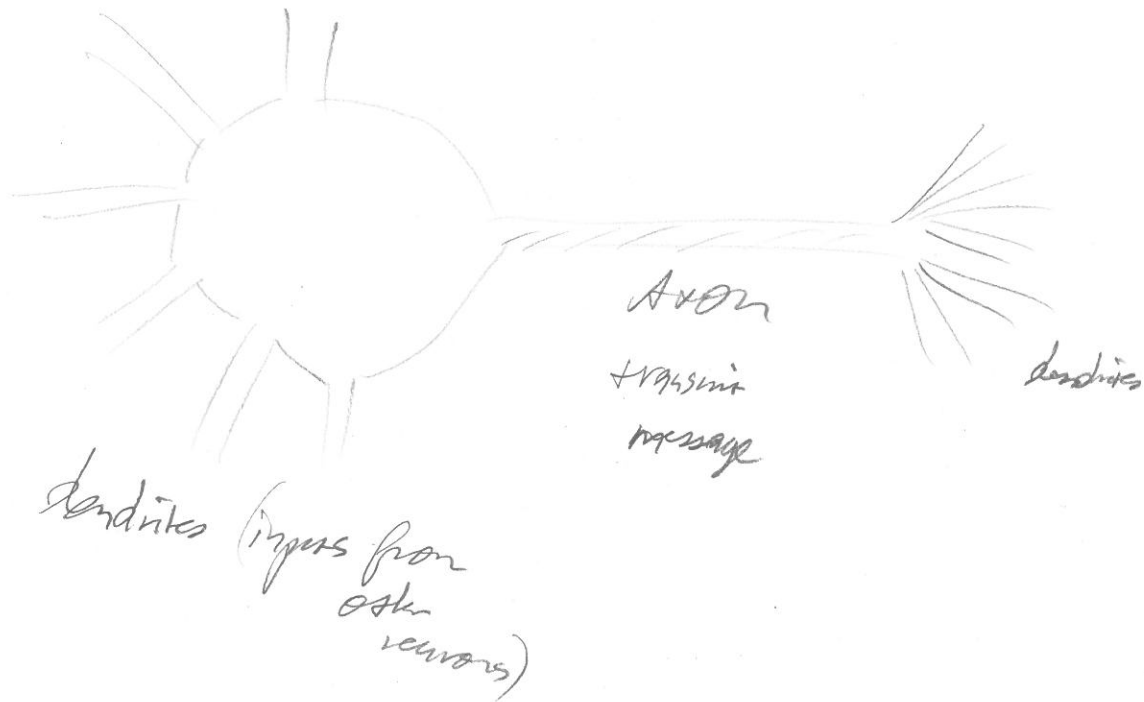
converged

Diagram of perceptron (a single layer neural network)

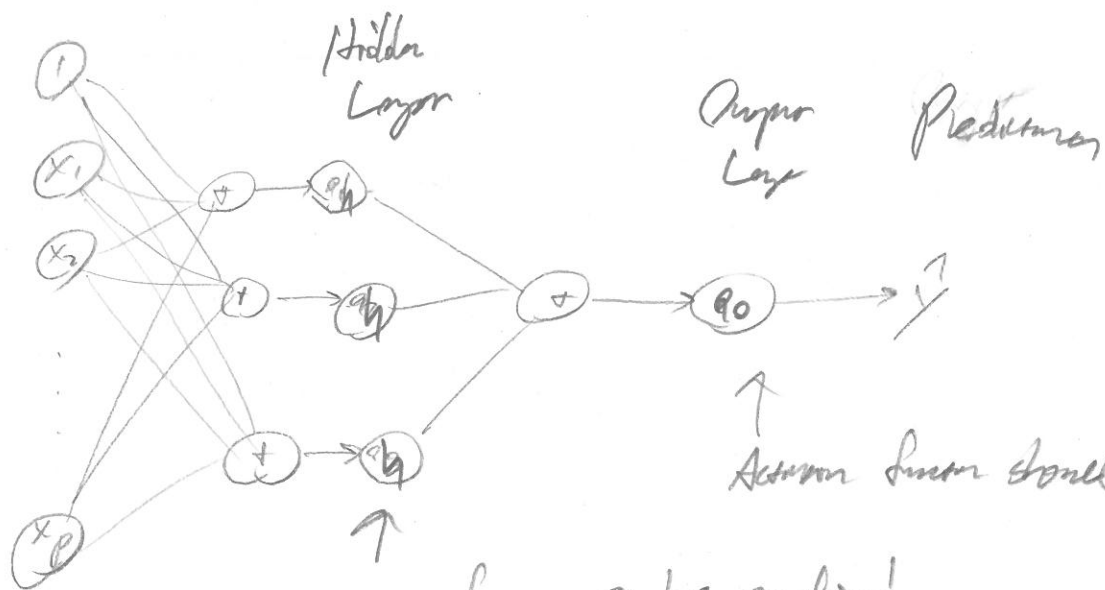


Input
layer

Why called neural network? Inspired by neuron..



Deep Learning



Input
Layer

Let's attempt to solve problems with the perceptron

① Finds "best" model out of all models if lin. sep

$$\mathcal{H} = \{ \perp \vec{w} \cdot \vec{x} \geq 0 : \vec{w} \in \mathbb{R}^{p+1} \}$$