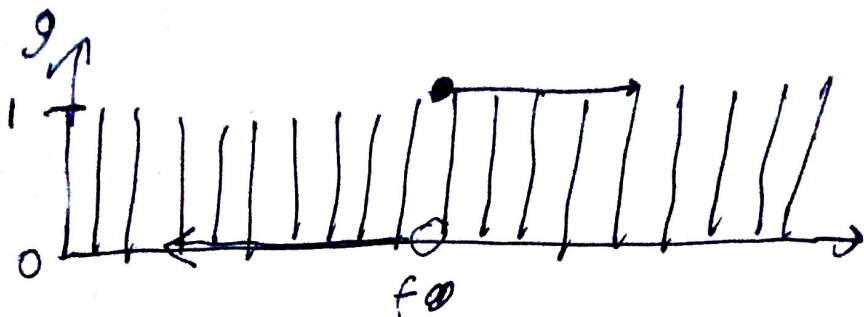


# Lecture 04

$$H = \{ \mathbb{I}_{x \geq \theta} : \theta \in \Theta(H) \}$$

model  $\nearrow$   
Parameter  $\nwarrow$  Parameter Space



Prediction

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

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

The algorithm  $A$  produces  $g$ . Since  $g$  is fully specified by  $\theta$ , the algorithm selects / estimates / optimizes / fits a  $\theta$ .

Let's create an algorithm. A bad algorithm will have high estimation error.

$$y \begin{matrix} & \begin{matrix} 0 & 1 \end{matrix} \\ \begin{matrix} 0 \\ 1 \end{matrix} & \begin{bmatrix} 0 & -1 \\ +1 & 0 \end{bmatrix} \end{matrix} e$$

Let's define an overall error function / objective function called "misclassification error" (ME)

$$ME = \frac{1}{n} \sum_{i=1}^n \mathbb{1}_{g(\vec{x}_i) \neq y_i} = \frac{1}{n} \sum_{i=1}^n |e_i|$$

or accuracy (ACC) as

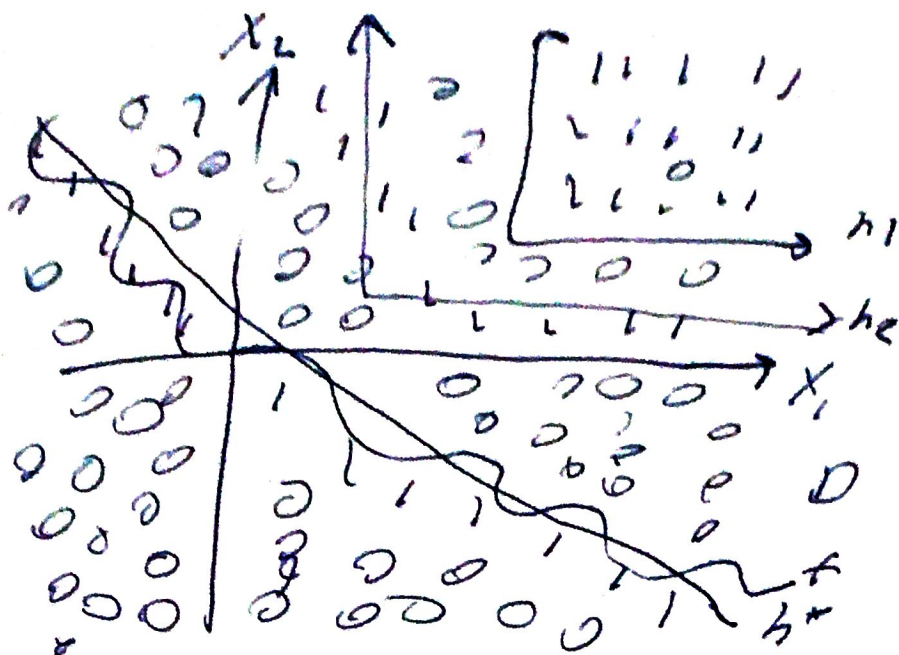
$$ACC = \frac{1}{n} \sum_{i=1}^n \mathbb{1}_{g(\vec{x}_i) = y_i} = 1 - ME$$

Goal of the algorithm is to minimize ME (or maximize ACC). To do so, we check every possible  $\theta \in \Theta$  and keep track of the ME( $\theta$ ) and then return the model with the lowest ME

How to define parameter space? It must be finite because we need to check (i.e. compute ME) each element. Gabriel says grid up  $[300, 850]$  e.g.  $\{351, 352, \dots, 849, 850\}$ . That's fine, but it's more convenient to only check the unique values of  $x$ .

$$\text{A produce } g(x) = \mathbb{1}_{x \geq \text{argmin} \left\{ \frac{1}{n} \sum_{i=1}^n \mathbb{1}_{x_i \neq y_i} \right\}}$$

Let's make a loan model with two continuous  $x$ 's i.e.  $x_1, x_2$  ( $p=2$ )  $\dim[\theta] = 2 = p$



A two dimensional threshold model  
 extending what we have before has  
 Candidate set

$$H = \{ \mathbb{1}_{x_1 \geq \theta_1} \mathbb{1}_{x_2 \geq \theta_2} : \begin{bmatrix} \theta_1 \\ \theta_2 \end{bmatrix} \in \mathbb{R}^2 \}$$

The candidate set of a "angle bracket"

-looking thing is very restrictive! which  
 means we will probably have high misspecification  
 error. Let's use another hypothesis set:  
 all lines

$$H = \{ \mathbb{1}_{x_2 \geq a + bx} : a \in \mathbb{R}, b \in \mathbb{R} \}$$

intercept
slope

The Slope and Intercept provide you with enough "degree of freedom" to specify any separating line. We need an algorithm to find  $g$  i.e. to specify  $a$  and  $b$ . This is a hard problem so we will study it with different conditions.

We will first reparameterize the hypothesis space to be:

$$\vec{w}, \vec{x} \geq 0$$

$$H = \left\{ \underset{\substack{\uparrow \\ \text{Intercept term} \\ \text{of "bias" }}}}{1} w_0 + \underset{\substack{\uparrow \\ \text{weight of} \\ \text{the first} \\ \text{feature}}} {w_1} x_1 + \underset{\substack{\uparrow \\ \text{weight of second feature}}} {w_2} x_2 \geq 0 : w_0 \in \mathbb{R}, w_1, w_2 \in \mathbb{R} \right\}$$

In order to fit this model, we "add" a dummy value of 1 to each data record

$$\vec{x} = [750 \ 850000] \rightarrow \vec{x} = [1 \ 750 \ 850000]$$

So we append the  $n$ -dim Counter vector to  $x$ , the matrix of features in  $D$

We only need 2 parameters  $(a, b)$  but here we have three  $(w_0, w_1, w_2)$  and hence we are "over-parameterized" meaning we have infinite solutions seen here.

$$\mathbb{I}_{\vec{w} \cdot \vec{x} \geq 0} = \mathbb{I}_{\vec{w} \cdot \vec{x} \leq 0} \quad \forall \vec{w} \neq 0 \quad \begin{matrix} x_1 + x_2 = 1 \\ p \\ p_n \end{matrix}$$

At find  $w_0, w_1, w_2$  to minimize ME  
i.e.

$$\vec{w}_* := \arg \min_{\vec{w} \in \mathbb{R}^3} \left\{ \sum_{i=1}^n \mathbb{I}_{\vec{w} \cdot \vec{x}_i \geq 0} = y_i \right\}$$

$$= \arg \min \{ME\}$$

We have a problem here. There is no analytic solution since the indicator function is non-differentiable

We need a way to search over all possible lines. So (1) we need to reduce the number of lines like before, (2) use an iterative algorithm to find a local solution (not the best but hopefully pretty good) or (3) change our objective function.

In the setting of perfect linear separability, e.g. where ME of that linear discrimination model is zero (i.e. no errors), consider the 1957 perceptron iterative algorithm for  $p$  features



Step 1 Initialize  $\vec{w}^{t=0} = \vec{0}_{p+1}$  or to a random vector value

Step 2: Compute  $z_i = \vec{w}^{t=0} \cdot \vec{x}_i \geq 0$

Step 3: For  $i = 0, 1, \dots, p$  set

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

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

$\vdots$

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

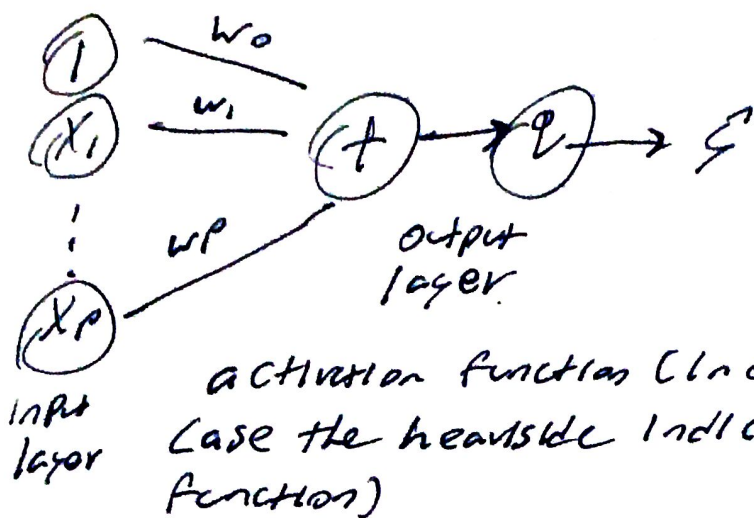
Step 4 Repeat steps 2 and 3 for  $i = 1, \dots, n$  (all the observations).

Step 5: Repeat steps 2, 3, and 4 until  $ME = 0$  i.e. all  $e_i$ 's are 0 or until a pre-specified (large) number of iterations

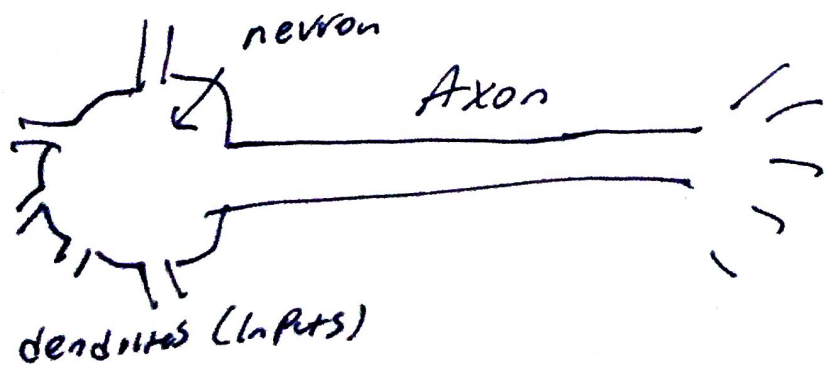
Note:  $t$  is the iteration number. It starts at 0,  $t=1$  is first iteration

The perceptron is proved to converge for linearly separable datasets but for non-linearly separable datasets, anything can happen so it may fail

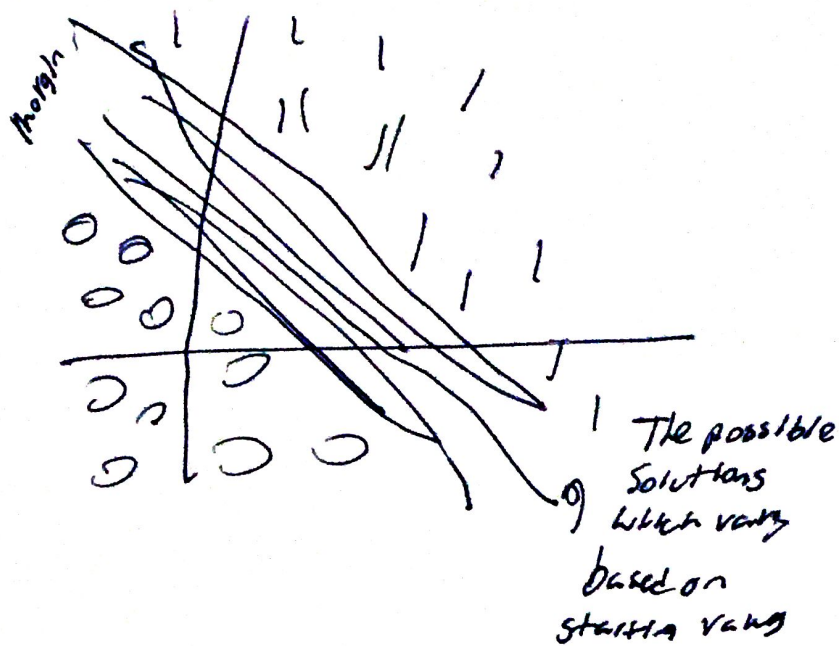
## Diagram of perception



The perceptron is a type of "neural network model". So are deep learning models. They're called neurons since they kind of act like neurons;



The perceptron has infinitely many solutions



But you kinda see there's a best model. This best model divides the margin (ANA wedge) evenly. This "best" model is called the "maximum margin hyperplane" and it was proven in 1998 to be the optimal linear classifier.