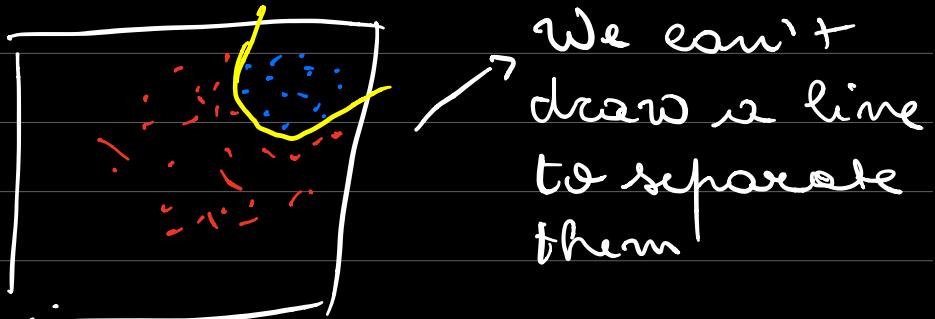
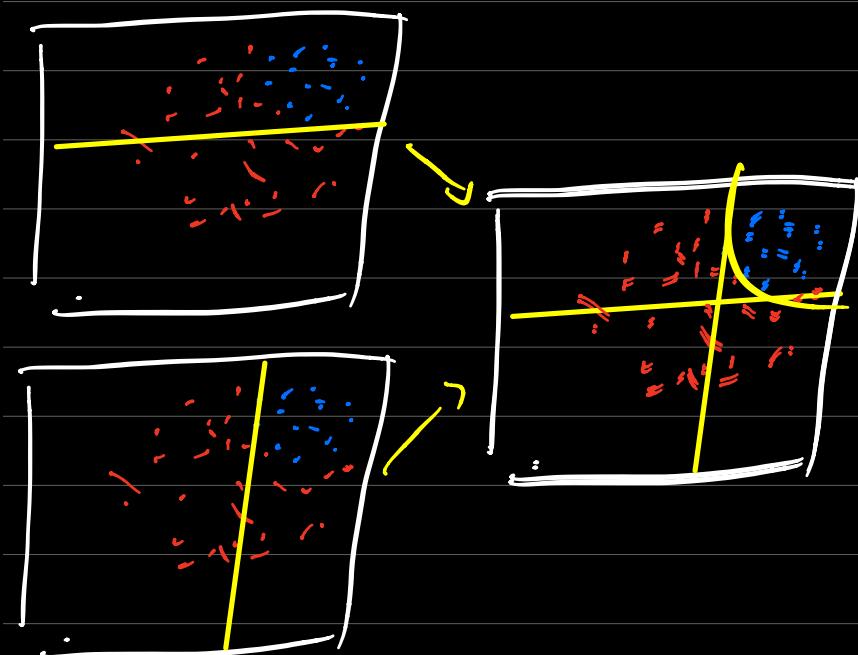


Non-linearly separable data



We use a simple trick:

- We combine two linear models into a non-linear one!



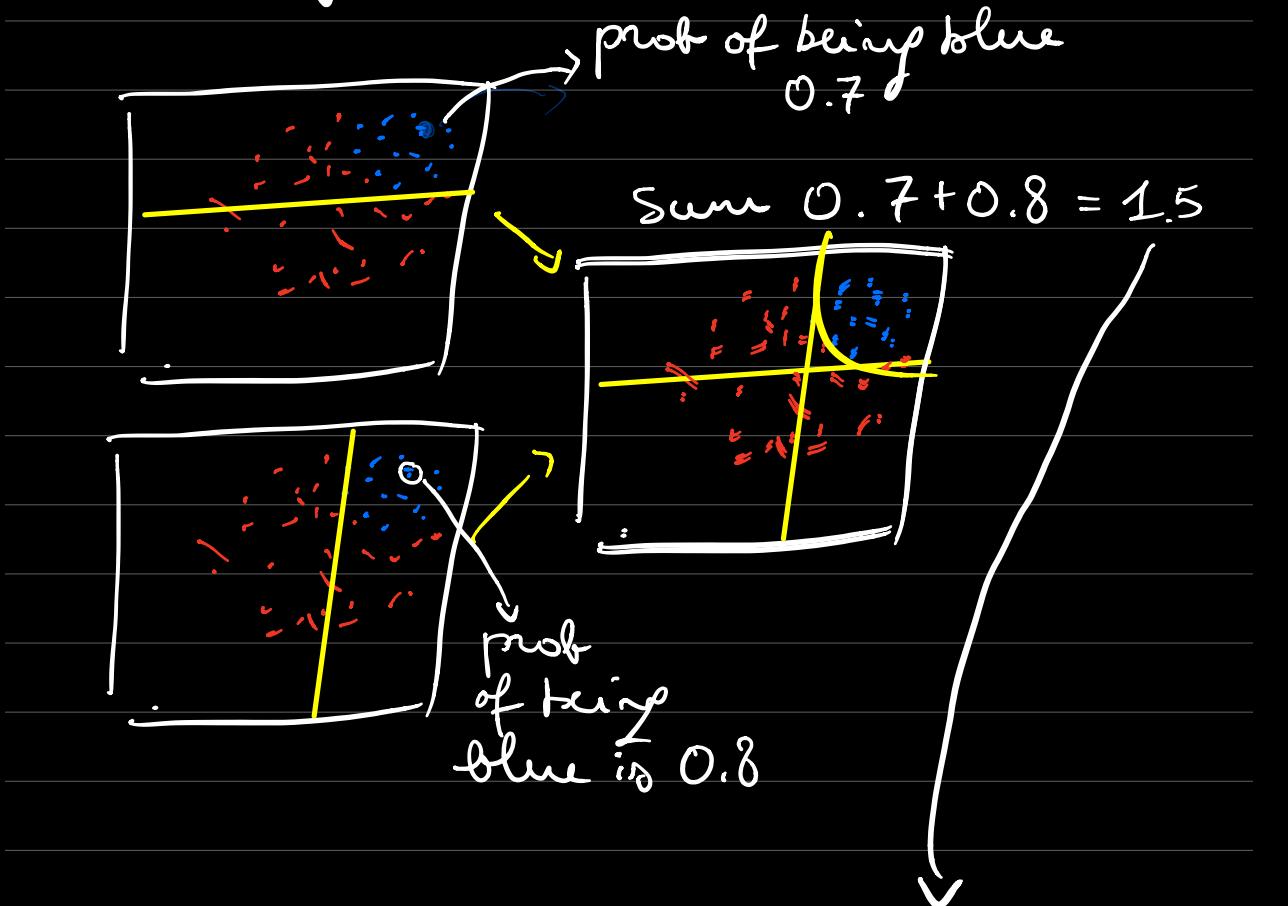
If we overimpose the two linear models we kind of get this the desired

one

if δ

How do we combine them?

We know that our linear models assign a probability to each point of being either blue or red.



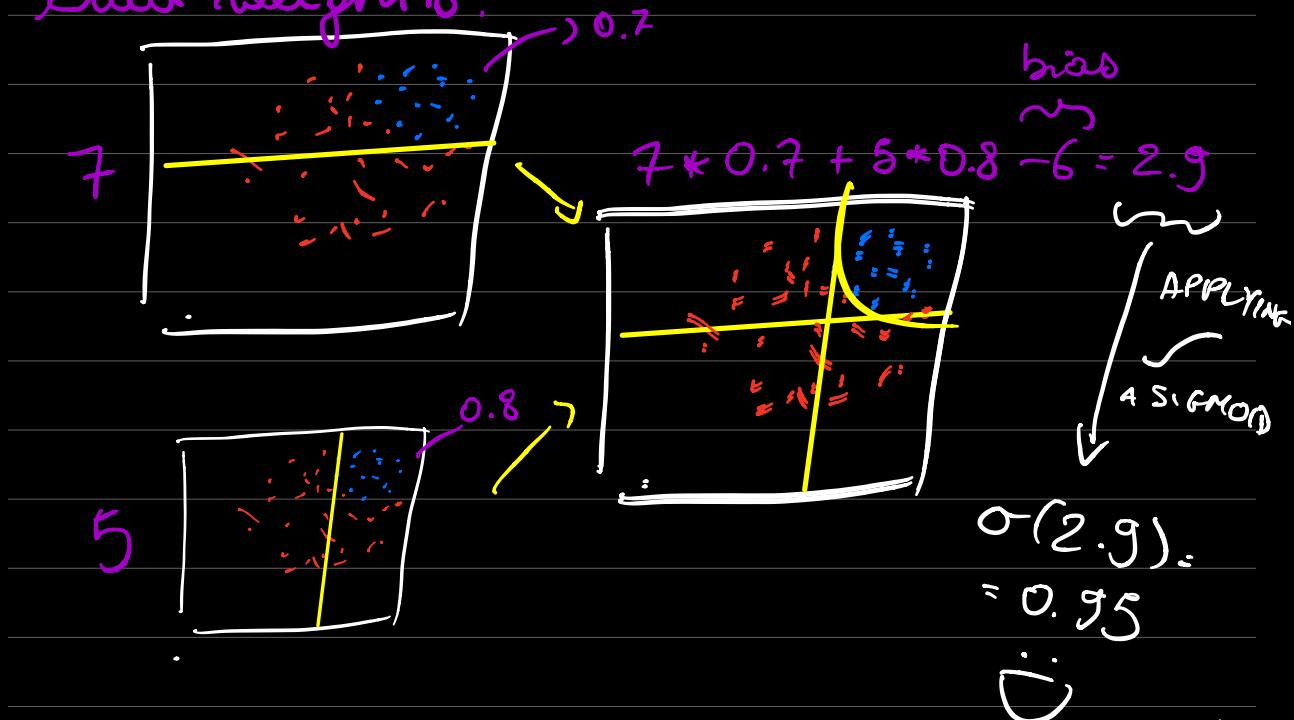


We can use a sigmoid σ

$$\sigma(1.5) = \frac{1}{1 - e^{-1.5}} = 0.82$$

What if we want to give more weights to the upper linear model than the bottom one?

Add weights!

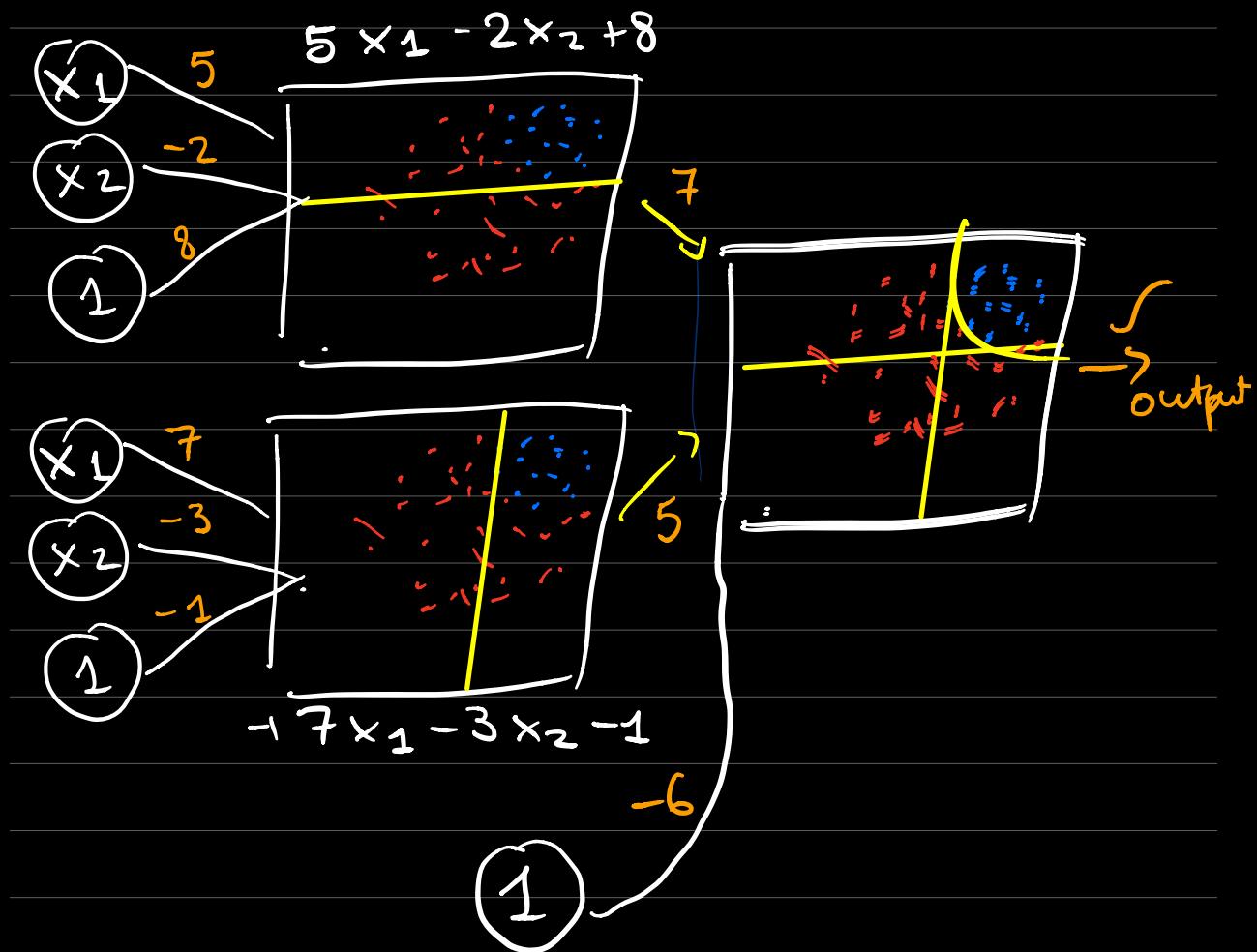


This is where we start to build a Neural Network.

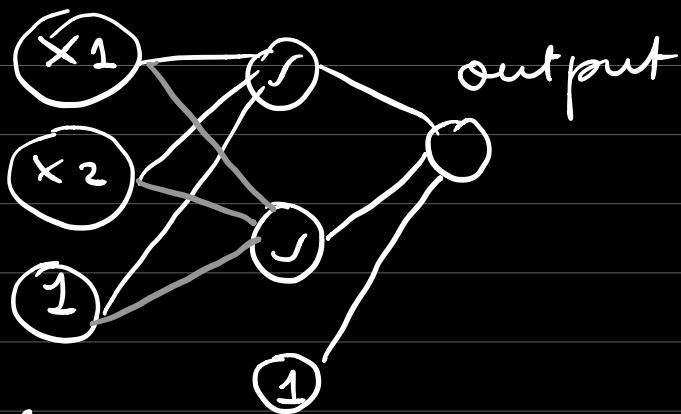
And as you could have noticed, we

did the same thing we did already with perceptrons, but using as input of the new perceptron the output of the previous ones.

Let's go into more Math, using the example of above



We represents neural networks
in the following way

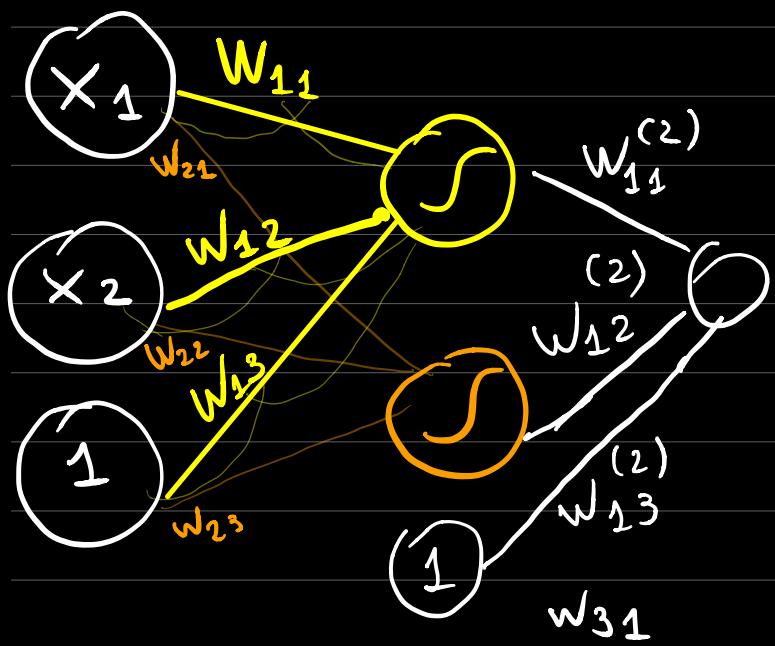


In training a neural network
there are two main steps:
Feed forward, forward pass
Backward, Backpropagation

Feed forward

input → output

⚠ MATH ALERT ⚠



$$W^2 = \begin{pmatrix} W_{11} & W_{12} & W_{13} \\ W_{21} & W_{22} & W_{23} \end{pmatrix}$$

$$\begin{aligned} & \downarrow \\ & \left(\begin{matrix} W_{11}^{(2)} & W_{12}^{(2)} & W_{13}^{(2)} \\ W_{21} & W_{22} & W_{23} \end{matrix} \right) \left(\begin{matrix} x_1 \\ x_2 \\ 1 \end{matrix} \right) \\ & \quad \left(\begin{matrix} W_{11}x_1 + W_{12}x_2 + W_{13} \cdot 1 \\ W_{21}x_1 + W_{22}x_2 + W_{23} \cdot 1 \end{matrix} \right) \end{aligned}$$

The diagram shows the calculation of the output of the second layer. The weight matrix W^2 is multiplied by the input vector $\begin{pmatrix} x_1 \\ x_2 \\ 1 \end{pmatrix}$. The result is a vector of two values: $W_{11}x_1 + W_{12}x_2 + W_{13} \cdot 1$ and $W_{21}x_1 + W_{22}x_2 + W_{23} \cdot 1$. The term $W_{13} \cdot 1$ and $W_{23} \cdot 1$ are labeled "bias".

As for the perceptron, our Neural Network will produce an error too.

How did the ERROR worked for the Perceptron?

So for each data point

$$\begin{array}{c} (x_i, y_i) \\ \downarrow \quad \searrow \\ (x_1, x_2, \dots, x_m) \text{ label (0 or 1)} \end{array}$$

we computed a score by doing

$$0 < \sigma(w^T x + b) < 1$$

$$\underbrace{\hat{y}}$$

$$\hat{y}$$

gives a probability.

What we want is to maximize the probability of being correct, that is minimize the probability of being wrong.

Let's start with an example.

Two models give me the probability of being accepted at the University:

- One says 80% chance
- The other 55% ..

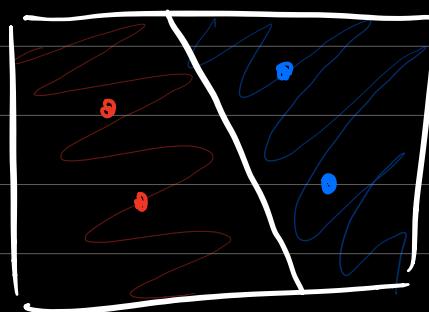
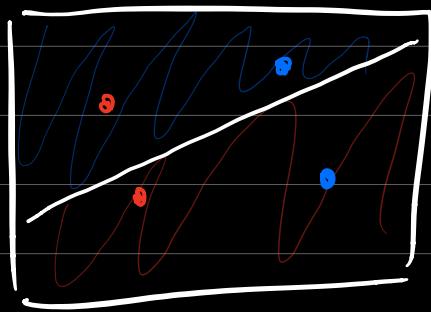
Which is better?

Well, if I actually get accepted, the first one.

So the one it is more likely to predict the correct answer.

This is called

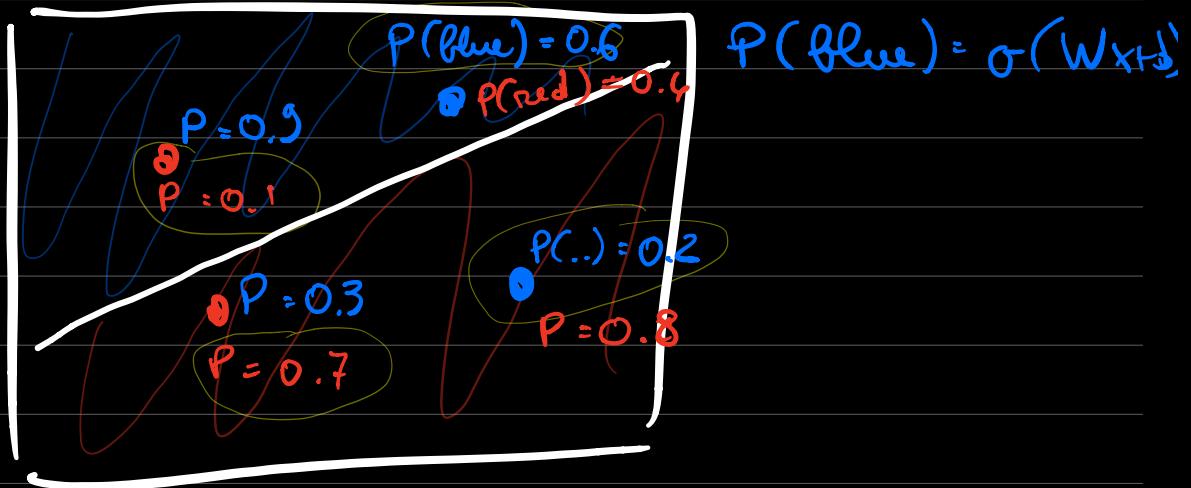
MAXIMUM LIKELIHOOD.



↓
this is more

correct

From a probability perspective

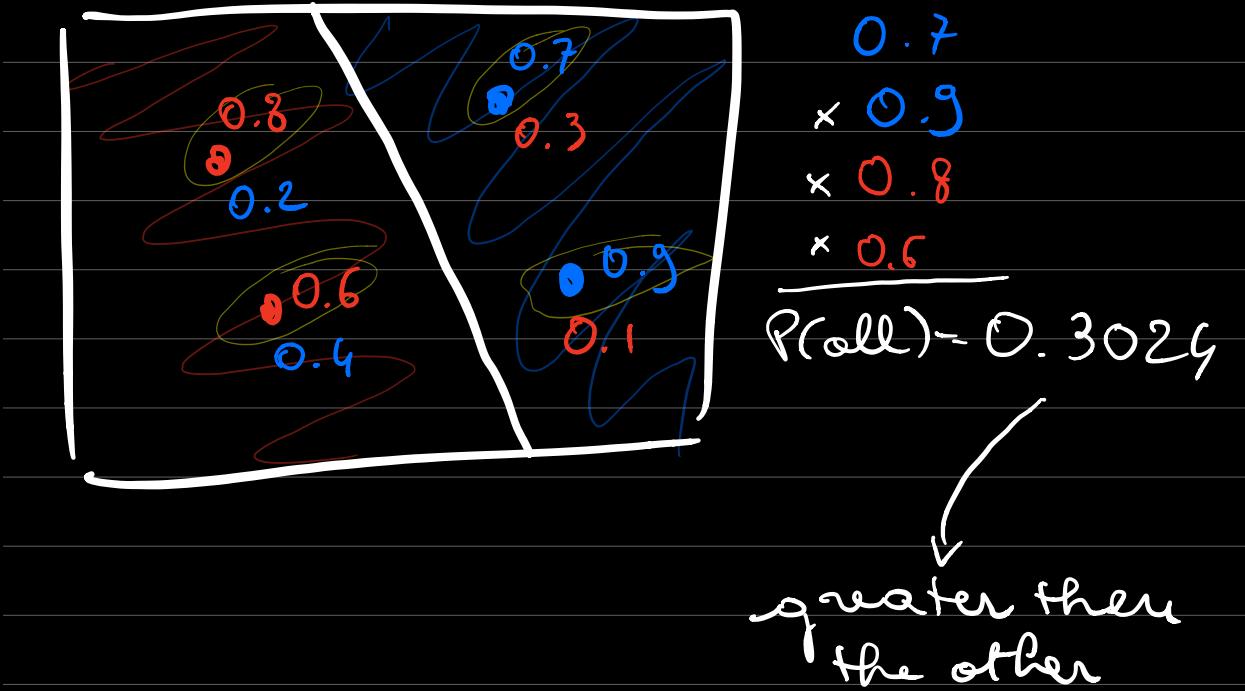


Since the colors of the points are INDEPENDENT EVENTS, then the probability of the model to guess the right label is the product:

$$\begin{aligned} P(\text{red}) &= 0.1 \\ \times P(\text{blue}) &= 0.6 \\ \times P(\text{red}) &= 0.7 \\ \times P(\text{blue}) &= 0.2 \end{aligned}$$

$$\underline{P(\text{all}) = 0.0084}$$

Proof of being correct for all



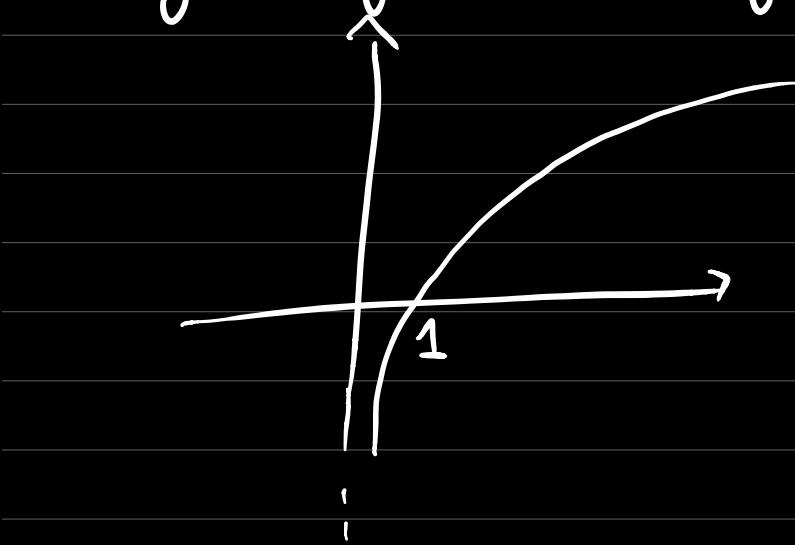
MINIMIZING THE ERROR FUNCTION
IS EQUIVALENT TO MAXIMIZING
THE PROBABILITY OF BEING CORRECT.

Problem: the probability we are computing requires products.

↓
harder
to compute

Transform them into sums by

writing $\log(ab) = \log(a)\log(b)$



so instead of having

$$0.6 \times 0.2 \times 0.1 \times 0.7 = 0.0084$$

\downarrow \downarrow \downarrow \downarrow

$$\log(0.6) + \log(0.2) + \log(0.1) + \log(0.7)$$

\downarrow

$$-0.51 - 1.61 - 2.3 - 0.36 = -4.8$$

The negative of this sum is called

CROSS-ENTROPY!

so a good model will give a

LOWER CROSS ENTROPY

(the other example gives 1.2)

So the triad is:

INPUT x , if label y is blue

$$P(x \text{ is blue}) = \hat{y}$$

$$\text{Error} = -\ln(\hat{y})$$

INPUT x , if label y is red

$$P(x \text{ is red}) = 1 - P(x \text{ is blue}) = 1 - \hat{y}$$

$$\text{Error} = -\ln(1 - \hat{y})$$

We can combine these things

$$\text{Error} = (\underbrace{(1 - \hat{y})(\ln(1 - \hat{y}))}_{\text{if } y \text{ is 1,}} - \underbrace{y \ln(\hat{y})}_{\text{if the label is 0, this is not used}})$$

if y is 1,
 $1 - 1 = 0$

so we use

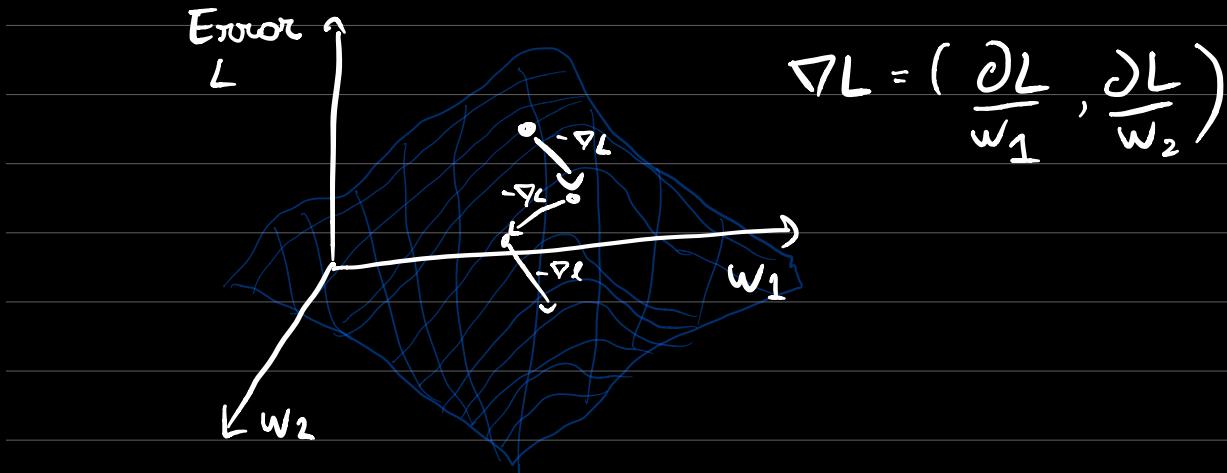
if the
label is 0,
this is not
used

The total error will be the mean for each of the data points

$$\text{Error} = -\frac{1}{m} \sum_{i=1}^m (1-y_i)(\ln(1-\hat{y}_i)) + y_i \ln(\hat{y}_i)$$

How do we minimize the ERROR?

Gradient Descent



∇L tells us the direction in which the function increases

$-\nabla L$ tells us the direction in which the function decreases

To avoid making too big updates, we will use a learning rate α

So we

- Compute probability $\hat{y} = \sigma(Wx)$

where $W = [w_1, w_2, b]$ in our example

- Compute the error

$$L = (1-y)(\ln(1-\hat{y})) - y \ln(\hat{y}) =$$

$$= (1-y)(\ln(1-\sigma(w_1x_1 + w_2x_2 + b))) - y \ln(\sigma(w_1x_1 + w_2x_2 + b))$$

- Compute the gradient

$$\nabla L = \left(\frac{\partial L}{\partial w_1}, \frac{\partial L}{\partial w_2}, \frac{\partial L}{\partial b} \right) =$$

$$\left((1-y) \frac{1}{1-\sigma(\dots)} \cdot -\sigma'(w_1x_1 + w_2x_2 + b) \cdot x_1 \right) -$$

$$- y \frac{1}{\sigma(w_1x_1 + \dots)} \cdot \sigma'(w_1x_1 + w_2x_2 + b) \cdot x_1$$

we know that $\sigma(x) = \frac{1}{1+e^{-x}} = (1+e^{-x})^{-1}$

$$\Rightarrow \sigma'(x) = -(1+e^{-x})^{-2} \cdot e^{-x} \cdot (-1) = \frac{1}{(1+e^{-x})^2} \cdot e^{-x}$$

$$\Rightarrow e^{-x} \sigma'(x) \text{ or } \sigma(x)(1-\sigma(x))$$

$$\Rightarrow - (1-y) \frac{1}{\underline{1-\sigma(\dots)}} \cdot \sigma'(1-\sigma(\dots)) \cdot x_1 -$$

$$y(1-\sigma(\dots)) \cdot x_1 = x_1 (\sigma(\dots) - \cancel{\sigma(\dots)} y -$$

$$- y + y \cancel{\sigma(\dots)}) = -x_1 (y - \sigma(\dots))$$

Indeed our update algorithm is what we did yesterday:

$$w_i = w_i - \alpha \frac{\partial L}{\partial w_i} = w_i + \alpha (y - \hat{y}) x_i$$

Clap Clap Clap!

Let's code now!

The problem in computing the error on a neural network is that we don't

have the "target" of the internal nodes.

Chain Rule

If you have a composition of functions:

$$g(f(x))$$

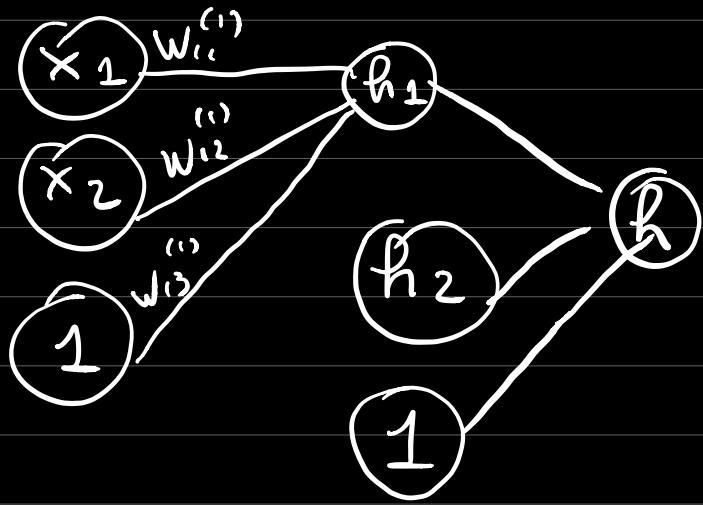
↳ That means

$$\begin{array}{ccc} x & \xrightarrow{f} & A \xrightarrow{g} B \\ \text{input} & A = f(x) & \hookrightarrow B = g(f(x)) \end{array}$$

the derivative is $\frac{\partial g(f(x))}{\partial x} = \frac{\partial g}{\partial f(x)} \frac{\partial f}{\partial x}$

so the derivative of compositions is just a multiplication.

In a Neural Network, what we are doing is composing multiple functions



$$h_1 = W_{11}^{(1)} x_1 + W_{12}^{(1)} x_2 + W_{13}^{(1)}$$

$$h_2 = W_{21}^{(1)} x_1 + W_{22}^{(1)} x_2 + W_{23}^{(1)}$$

Then $y = W_{11}^{(2)} \sigma(h_1) + W_{12}^{(2)} \sigma(h_2)$
 $+ W_{23}^{(2)}$

$$\Rightarrow \hat{y} = \sigma(y)$$

We condense it into a matrix notation

$$\hat{y} = \sigma(W^{(2)}(\sigma(W^{(1)}(x))))$$

We will compute the error

$$L(w) = -\frac{1}{m} \sum_{i=1}^m (y_i \ln(\hat{y}_i) + (1-y_i) \ln(1-\hat{y}_i))$$

$$\nabla L = \left(\frac{\partial L}{\partial w_{11}^{(1)}}, \dots, \frac{\partial L}{\partial w_{31}^{(2)}} \right)$$

$$\frac{\partial E}{\partial w_{11}^{(2)}} = \frac{\partial E}{\partial \hat{y}} \frac{\partial \hat{y}}{\partial h} \frac{\partial h}{\partial h_1} \frac{\partial h_1}{\partial w_{11}^{(1)}}$$

And if you will notice a nice pattern, that is:

Remember $h = \underbrace{w_{11}^{(2)} \sigma(h_1) + w_{12}^{(2)} \sigma(h_2)}_{+ w_{21}^{(2)}} + \dots$

$$\frac{\partial h}{\partial h_1} = w_{11}^{(2)} \sigma(h_1)(1-\sigma(h_1))$$

↓
we have
it from
the forward
pass!

In practice, by doing the forward

pass we provide already the ingredient for computing the backword pass, that will give us the partial derivative of the error w.r.t. whatever weight!