# ② NEURAL NETWORKS

## INTRODUCTION TO NEURAL NETWORKS (LESSON 1)

### CLASSIFICATION & PERCEPTRONS



$X_2$ (grades)    rejected

accepted

classification model

$w_1 x_1 + w_2 x_2 + b = 0$

University entrance
- rejected
- accepted

$X_1$ (test score)

all are summed

activation function
- step function
- sigmoid
- relu...

$X_1, w_1$

$X_2$  $w_2$   $\Sigma$  $\sigma$

bias b

$\geq 0$: reject
$< 0$: accept

PERCEPTRON:
building block of neural networks: it's basically a classification model!
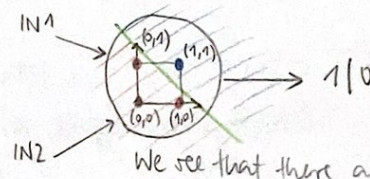
<mark>See CVND p.16 - Perceptrons</mark>

### Perceptrons as logical operators

- logical operators can be modeled with perceptrons!

Example:   AND

| IN 1 | IN 2 | ouT |
|------|------|-----|
| 1 | 1 | 1 • |
| 1 | 0 | 0 • |
| 0 | 1 | 0 • |
| 0 | 0 | 0 • |

IN1

(0,1)  (1,1)
(0,0)  (1,0)

IN2

$\rightarrow$ 1|0

We see that there are actually many possible models...

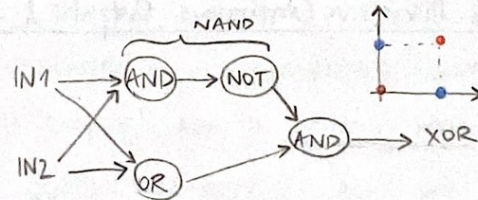| $w_1$ | $w_2$ | b |
|-------|-------|------|
| 1 | 1 | -1.5 |
| 1 | 1 | -1.75 |
| ... | | |

- It's easy to create models for OR, NOT, AND

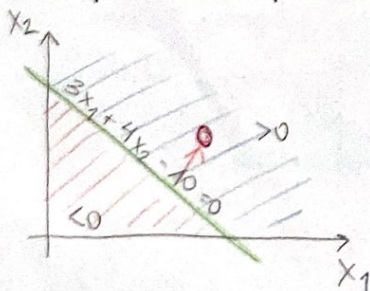- For more complex operators, eg XOR, we combine other logical operators / perceptron classification models

Example:  XOR

| IN 1 | IN2 | ouT |
|------|-----|-----|
| 1 | 1 | 0 |
| 1 | 0 | 1 |
| 0 | 1 | 1 |
| 0 | 0 | 0 |

only true if one is different

NAND

IN1 $\rightarrow$ AND $\rightarrow$ NOT

IN2 $\rightarrow$ OR

$\rightarrow$ AND $\rightarrow$ XOR



### Perceptron model optimization



$X_2$

$3x_1 + 4x_2 - 10$
$>0$
$<0$

$X_1$

$3x_1 + 4x_2 - 10 = 0$ model line

- miss classified point O :  (4,5)
↳ a fast method to move the model line towards the missclassified point consists in subtracting the point coords to the model params:

| | | | |
|---|---|---|---|
| - | 3 | 4 | -10 | model |
| | 4 | 5 | 1 | point (1 for bias) |
| | -1 | -1 | -11 | new model |

if move to ⊕ area, then ⊖, else ⊕

1

- But usually a <u>learning rate</u> is used for the point values to move slowly, eg:

```
   3    4   -10    model
  0.4  0.5  0.1    point × 0.1 learning rate
  2.6  3.5 -10.1 ] ——→ new model, closer to missclassified point!
```

$2.6x_1 + 3.5x_2 - 10.1 = 0$
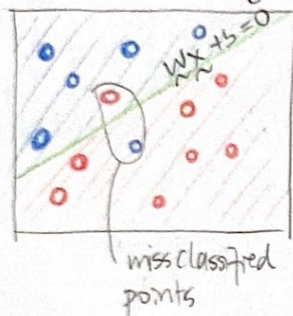
- If we want to move to negative area, we add instead of substracting!

—▷ PERCEPTRON ALGORITHM: it's based on the previous trick

   Given: - m samples of dimension n each
         - Start with random model weights $W_1, ..., W_n$ and $b$
         - learning rate $\alpha$



miss classified points

for every missclassified point $x$

  if $y = Wx + b < 0$ / prediction was $1/\oplus$ and we
                                  want to move to $0/\ominus$
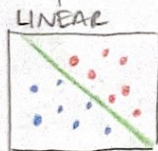    $W_i = W_i + \alpha x_i$
    $b_i = b_i + \alpha$

  else if $y = Wx + b < 0$ / prediction was $0/\ominus$ and we
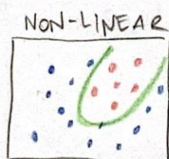                                  want to move to $1/\oplus$
    $W_i = W_i - \alpha x_i$
    $b_i = b_i - \alpha$

—▷ IMPORTANT: a perceptron is a linear classifier; if we want more complex
    classification regions, we need to use more perceptrons!
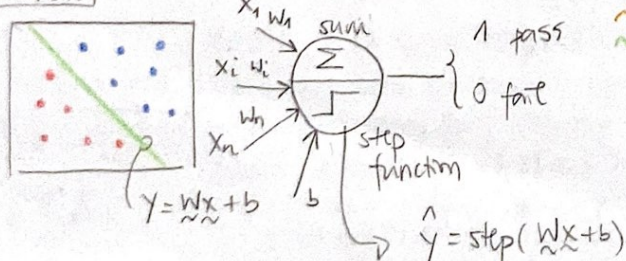


single neuron / perceptron

several neurons

▷ $\sigma : \mathbb{R} \longrightarrow [0, 1]$
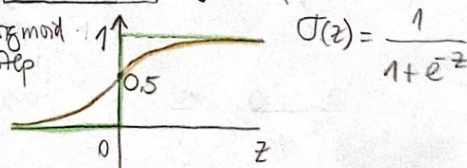all real numbers mapped to $[0, 1]$

<u>Error & Discrete vs Continuous Outputs & sigmoid function</u>

- Although classification is a DISCRETE problem (yes/no, dog/cat/rat),
the <u>error function</u> is used to find the model parameters with <u>gradient descend</u>,
and for that, we need CONTINUOUS outputs and error function: the <u>sigmoid</u>
function achieves that.

DISCRETE



$y = Wx + b$

$\hat{y} = step(Wx + b)$

1 pass
0 fail

step function

CONTINUOUS | sigmoid function



~ sigmoid
~ step

$\sigma(z) = \dfrac{1}{1 + e^{-z}}$

$\hat{y} = \sigma(Wx + b)$ ; $\hat{y}$ tells us the probability of pass (close to 1) or fail

2

## Softmax function for 2+ classes

→ SOFTMAX converts the outputs in PROBABILITIES

$$\hat{y}$$

class 1 $\begin{bmatrix} 2.6 \\ 0.5 \\ -3.7 \end{bmatrix}$  $\xrightarrow{\text{SOFTMAX}}$  $P$ $\begin{bmatrix} 0.87 \\ 0.10 \\ 0.03 \end{bmatrix}$

class 2

class 3

$\dfrac{e^{z_i}}{\sum e^{z_i}}$ : soft max

Softmax for 2 classes is equivalent to the sigmoid!

↳ sum = 1

Note: usually the one-hot encoding is used for multi-class problems:
for an input, the output has the size of the # classes.

## Maximum Likelyhood

A

$P(r) = 0.7$
$P(b) = 0.3$

$P(r) = 0.5$
$P(b) = 0.5$

$P(r) = 0.8$
$P(b) = 0.2$

$P(r) = 0.3$ $P(b) = 0.7$

B

$P(r) = 0.2$  $P(r) = 0.3$
$P(b) = 0.8$  $P(b) = 0.7$

$P(r) = 0.6$  $P(r) = 0.7$
$P(b) = 0.3$  $P(b) = 0.3$

Model A is clearly better.
We can check that numerically with the maximum likelyhood method

$\downarrow$

We computing a model, we're finding the one that has the biggest maximum likelyhood.

$P(red) = 1 - P(blue)$

MODEL A: maximum likelihood of all points & probability of **correct** values

$P(all) = 0.7 \times 0.8 \times 0.5 \times 0.7$

$= 0.196$

MODEL B:

$P(all) = 0.2 \times 0.6 \times 0.7 \times 0.3$

$= 0.025$

→ the maximum likelihood of A is much higher ⟹ better model!

- The main issue with the product of probabilities is that they become too small
- Solution: we use the logarithm:

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

$\ln(ab) = \ln(a) + \ln(b)$ ← the natural logarithm is usually used, but we could take $\log_{10}$ too

$\ln(0.7 \times 0.9) = -0.36 - 0.1$

logarithm is always $< 0$ in $[0,1]$

## Cross Entropy

Negative sum of the logarithms of the probabilities:

$C.E.: -\sum_{\text{samples}}^{m} \ln(p_j) \rightarrow \begin{cases} \uparrow\uparrow \text{ BAD model} \\ \downarrow\downarrow \text{ GOOD model} \end{cases}$

→ Now, our goal has become to minimize the Cross Entropy instead of maximizing the P(all)

$C.E. = -\sum_{\text{samples}}^{m} y_j \ln(p_j) + (1-y_j) \ln(1-p_j)$ ; $y = \{^1_0\}$ →

Note that for binary classification one term cancels!

Event  P(Event)  complementary event  P(complementary) event

3

**\* Multi-class cross-entropy:**

$$C.E. = -\sum_i^n \sum_j^m y_{ij} \ln(p_{ij})$$

probability of event $y_{ij}$

classes    samples

event: $y_{ij} = 1$ if sample $j$ belongs to class $i$
                  $\emptyset$ otherwise

this formula is equivalent to the previous!

<u>Error Function.</u>  — in the previous (binary classification), a term is always cancelled!

— The error function is basically the averaged cross entropy:

Binary: $E = -\left(\frac{1}{m}\right)\sum_j^m y_j \ln(p_j) + (1-y_j)\ln(1-p_j)$

Multiclass: $E = -\left(\frac{1}{m}\right)\sum_i^n \sum_j^m y_{ij} \ln(p_{ij})$

Since: $p = \sigma(\underset{\sim}{W}x + b) \rightarrow E(\underset{\sim}{W},b) = -\frac{1}{m}\sum_j^m y_j \ln(\sigma(\underset{\sim}{W}x+b)) + (1-y_j)\ln(1-\sigma(\underset{\sim}{W}x+b))$

— We are going to find the $\underset{\sim}{W}$ & $b$ values that minimize this error with gradient descend

## Gradient Descend

$f \longrightarrow \underset{\sim}{\nabla}f = \left[\frac{\partial f}{\partial x_1}, \frac{\partial f}{\partial x_2}, -, \frac{\partial f}{\partial x_n}\right]$ : Gradient of $f$: direction of maximum increase!

$E(\underset{\sim}{W}, b) = E(W_1, W_2, ..., W_n, b)$

$$\underset{\vee}{\nabla}E = \left[\frac{\partial E}{\partial W_1}, \frac{\partial E}{\partial W_2}, -, \frac{\partial E}{\partial W_n}, \frac{\partial E}{\partial b}\right]$$

We update:

$$W_i' = W_i - \alpha \cdot \frac{\partial E}{\partial W_i}$$

$$b' = b - \alpha \frac{\partial E}{\partial b}$$

Gradient descend: we initialize the $\underset{\sim}{W}$ & $b$ with random values and iteratively change them with the gradient until the error is small enough / # iterations max.

learning rate: small to avoid large uncontrolled steps

In a binary classification problem with sigmoids and cross-entropy error:

$$\left.\frac{\partial E}{\partial w_j}\right| = ... = -\sum_i^m (y_i - \hat{y}_i) \cdot x_{j} \cdot \frac{1}{m}$$

$$x = [x_1, ..., x_j, ..., x_n]$$

$$\frac{\partial E}{\partial b} = ... = -\sum_i^m (y_i - \hat{y}_i) \cdot \frac{1}{m}$$

$\rightarrow$ sum over all $m$ samples

scalar: $d$?
truth − prediction

$$\nabla E = -\sum_i (y_i - \hat{y}_i) \cdot [x_1, x_2, -, x_j, -, x_n, 1] \cdot \frac{1}{m}$$

bias    average

ground    prediction
truth

$X_i$: point with $n$ features

$\hat{y} = \sigma(\underset{\sim}{W}x + b)$

$i$: samples
$j$: features of each sample

Perceptron Algorithm vs Gradient Descend Algorithm: They are basically the same, but gradient descend repels the classification boundary if correct classification

Gradient Descend

$$w_j \leftarrow w_j + \alpha \frac{1}{m} \sum_i^m (y_i - \hat{y}_i) x_{ij}$$

If correct classification, NOTHING is done in that term!

Perceptron

$$w_j \leftarrow w_j + \alpha \frac{1}{m} \sum_i^m x_{ij}$$ if missclassified, nothing otherwise

the only thing missing is

$(y_i - \hat{y}_i)$

which:

• $y_i - \hat{y}_i = 0$ if correct

• $y_i - \hat{y}_i = 1 | (-1)$

→ that's the only difference!

| | Gradient descend | Perceptron |
|---|---|---|
| Correct classification | Repel further boundary | Do nothing |
| Miss classification | Attract boundary | Attract boundary |

this makes sense: error becomes smaller and our region increases

See Jupyter Notebooks: Perceptrons, Gradient Descend, Student Admissions

Nonlinear Models: Multi-layer Perceptrons

See CIND p-16    If we combine several neurons (Multi-layer Perceptrons) we obtain NEURAL NETWORKS, which can handle NONLINEAR models.

We have
- input layer
- hidden layers: If we have several hidden layers, we have deep neural networks.
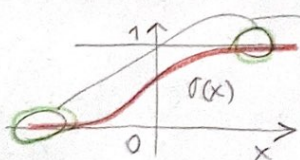- output layer

Feed forward
Back propagation

see p. 16 - 18 of the CIND.
19, 20

↓

IMPORTANT: keep in mind we are differentiating functions and applying the chain rule, thus, if a function has derivative almost $\emptyset$, the total error derivative will be $\emptyset$!
Therefore, we can easily run into numerical issues...

→ derivative is almost $\emptyset$ for large $|x|$ values!

This is called the VANISHING GRADIENT problem

$\sigma(x)$

# IMPLEMENTING GRADIENT DESCEND (LESSON 2)

## Different Error Functions

$$E_{cross-entropy} = -\sum_i^n \sum_j^m y_{ij} \ln(\hat{y}_{ij})$$

classes   samples

event: TRUE

output/probability of event: PREDICTION

$\hat{y}_{ij} = P_{ij} = \sigma(\cdot) \in [0,1]$

mean    Sum of Squared Errors

$$E_{MSE} = \frac{1}{m} \cdot \frac{1}{2} \sum_i^n \sum_j^m (y_{ij} - \hat{y}_{ij})^2$$

‖
Mean Squared Error

classes samples

event: TRUE

prediction; probability PREDICTION

$\hat{y} = f(\underline{w,b,x})$

we are looking for the w & b values that minimize the error E

| sample | | | class | | |
|---|---|---|---|---|---|
| 0 | 1 | | 1 | | |
| | 0 | | 2 | | |
| | 0 | | 3 | | |
| 1 | 0 | | 1 | | |
| | 1 | | 2 | | |
| | 0 | | 3 | | |
| 2 | ⋮ | | 1 | | |
| | | | 2 | | |
| | | | 3 | | |

$\hat{y}$
| 0.9 |
| 0.1 |
| 0 |
| 0.2 |
| 0.6 |
| 0.2 |

## Gradient Descend (notes on the notebook for Student Admission)

$$W_k \leftarrow W_k \boxed{- \eta \cdot \frac{\partial E}{\partial w_k}}$$

$$b \leftarrow b - \eta \frac{\partial E}{\partial b}$$

$$\Delta w_k = \eta \, \delta \, x_k \quad : \quad W_k \leftarrow W_k + \Delta w_k \qquad \propto -\frac{\partial E}{\partial w_k} \text{ prop.}$$
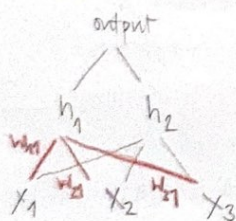
→ value of feature k

learning rate, aka $\alpha$

error term: $\delta = (y - \hat{y}) f'(\sum w_k x_k)$ : Note: it contains already the '−' sign !!!

$f(h) = \sigma(h), \quad h = \sum w_k x_k$

activation function



$\sigma(w_1 x_1 + w_2 x_2 + b) = $ output

input  b

## Multi-Layer Perceptron  see CVND 16-19

output



$h_1 \quad h_2$

$$\begin{array}{c} x_1 \\ x_2 \\ x_3 \end{array} \begin{bmatrix} w_{11} & w_{12} \\ w_{21} & w_{22} \\ w_{31} & w_{32} \end{bmatrix}$$

$f_j = \sigma(\sum_i w_{ij} x_i)$
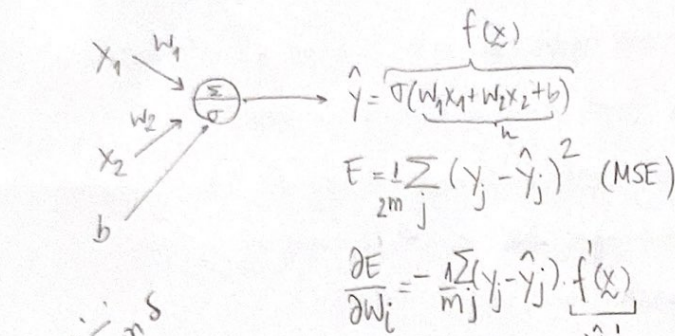
$\sigma(W^T x) = f$

$$\sigma\left(\begin{bmatrix} w_{11} & w_{21} & w_{31} \\ w_{12} & w_{22} & w_{32} \end{bmatrix} \begin{bmatrix} x_1 \\ x_2 \\ x_3 \end{bmatrix}\right) = \sigma\left(\begin{bmatrix} h_1 \\ h_2 \end{bmatrix}\right) = \begin{bmatrix} f_1 \\ f_2 \end{bmatrix}$$

Note: Numpy
- Transpose of a matrix (defined as array): arr.T
- BUT, for 1-D arrays, we arr[:, None]

6

# Backpropagation : a nice and intuitive explanation is intended here

$$\hat{y} = \sigma(\underbrace{w_1 x_1 + w_2 x_2 + b}_{h})$$

$$E = \frac{1}{2m} \sum_j (y_j - \hat{y}_j)^2 \quad (MSE)$$

$$\frac{\partial E}{\partial w_i} = -\frac{1}{m}\sum_j (y_j - \hat{y}_j) \cdot \underbrace{f'(x)}_{\hat{y}'}$$

because $\hat{y}'$ is already in $\delta$

ERROR TERM: $\delta = (y - \hat{y}) \cdot \boxed{f'(\underline{x})}$

WEIGHT UPDATE: $w_i \leftarrow w_i + \boxed{\frac{\alpha \sum_j \delta}{m}}$ → average across all samples $j$; ↗ all error terms

NOTE:
we are with a single perceptron / neuron here...

$$f'(\underline{x}) = \frac{\partial h}{\partial w_i} = \frac{\partial \hat{y}}{\partial w_i} = \frac{\partial \sigma(h(w_i \underline{x}))}{\partial w_i} =$$

$$= \frac{\partial \sigma}{\partial h} \cdot \frac{\partial h}{\partial w_i} = \sigma(h)(1 - \sigma(h)) \cdot \frac{\partial h}{\partial w_i} = \boxed{\sigma(\underline{x})(1 - \sigma(\underline{x})) \cdot x_i}$$

$$h(\underline{w}, \underline{x}) = w_1 x_1 + w_2 x_2 + b = \underline{w}^T \underline{x} + b$$

① Example with a unique perceptron

→ Error term = difference prediction output multiplied by the derivative of the output function

(note: we're using MSE and sigmoid)

learning rate, aka $\eta$

‖ value of $x$ only in dimension $i$

$\sigma'$: if we use another activation function, then this term is different.

② Example with a Multi-Layer Perceptron network:

Feedforward:

$$\hat{y} = f(h_1, h_2) = \sigma(w_1 \sigma(h_1) + w_2 \sigma(h_2))$$

$$= \sigma\left( w_1 \sigma(w_{11} x_1 + w_{21} x_2 + w_{31} x_3) + w_2 \sigma(w_{12} x_1 + w_{22} x_2 + w_{32} x_3) \right)$$

$$\begin{cases} h_1 = (w_{11} x_1 + w_{21} x_2 + w_{31} x_3) \longrightarrow \sigma(h_1) \\ h_2 = (w_{12} x_1 + w_{22} x_2 + w_{32} x_3) \longrightarrow \sigma(h_2) \end{cases}$$



Errors: $\delta = (y - \hat{y}) \cdot f'(\cdot)$ , $f' \equiv$ output function derivation of the node wrt weight we want to update
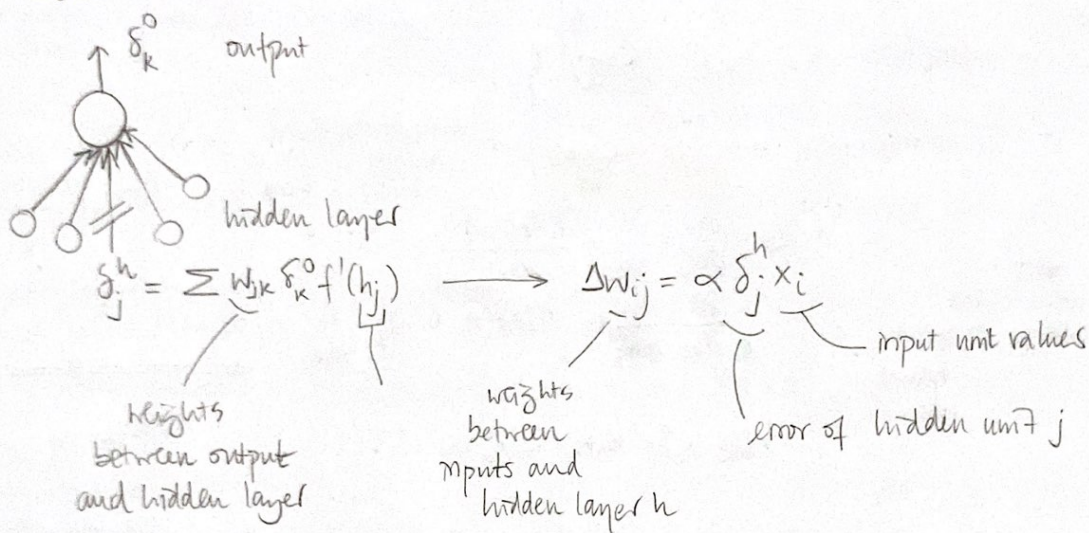
$$\delta_1 = \delta \cdot w_1 \cdot f'(h_1) \qquad \delta_2 = \delta \cdot w_2 \cdot f'(h_2)$$

previous error | weight | derivative

Therefore:
1) Bigger weights, bigger errors
2) Derivative must be considered carefully: not → ∞ not → 0

Backpropagation propagates the error backwards in the network! (analogous to propagating inputs forwards); Now the error is the input!

→ the oldest errors are propagated by being multiplied by: (1) weights (2) function derivative

- For training/learning, we want to update the weights so that the error is minimized
- Backpropagation is gradient descend with the chain rule and it turns out it's like applying the weights backwards to the final error — hence, the name.
- We take the final error and basically scale it through the network multiplying it by the weights along the way: this gives the portion/contribution of each weight to the final error ⇒ the update each weight needs.



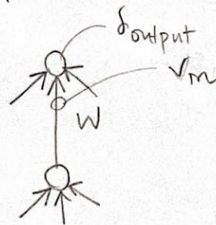$\uparrow \delta_k^o$  output

hidden layer

$$\delta_j^h = \sum W_{jk}\, \delta_k^o\, f'(h_j) \longrightarrow \Delta w_{ij} = \alpha\, \delta_j^h\, x_i$$

weights between output and hidden layer

weights between inputs and hidden layer h

input unit values

error of hidden unit $j$

- In general, for any level/layer:

$$\boxed{\Delta w = \alpha\, \delta_{output}\, V_{in}}$$

$\delta_{output}$

$V_m$

W

output error in node w points to

value through ling with weight W

③ Example exercise: compute $\Delta w$ for a simple network manually:

$\hat{y}$: output: should be 1.

$x_0 \leftarrow$

$0.1 = W_0$

$w_1 = 0.4$

$-0.2 = W_2$

$x_1 = (0.1)$  $x_2 = (0.3)$  INPUT

output:
$$h = \sum W_i x_i = 0.4 \cdot 0.1 - 0.2 \cdot 0.3 = -0.02$$
$$f(h) = \sigma(h) = \frac{1}{1 + e^{+0.02}} = 0.495 = x_0$$
$$\hat{y} = f(W_0 x_0) = \sigma(W_0 x_0) = \frac{1}{1 + e^{-0.0495}} = 0.512$$

- - - - -

Error in output:
$$\delta^o = (y - \hat{y}) \cdot f'(W_0 x_0)$$
$$= (y - \hat{y}) \cdot f(W_0 x_0)(1 - f(W_0 x_0)) =$$
$$= (1 - 0.512) \cdot 0.512 (1 - 0.512) = 0.122$$

8

Error in first layer:

$$\delta_j^h = \sum_K W_{jk} \delta_k^o f'(h_j) = W_0 \delta^o f'(h) = 0.1 \cdot 0.122 \cdot 0.495(1-0.495) = 0.003 = \delta^h$$

$$\Delta W_0 = \alpha \delta^o \cdot X_0 = 0.5 \times 0.122 \cdot 0.495 = 0.0302$$

$$\Delta W_1 = \alpha \boxed{\delta^h} \cdot X_1 = 0.5 \times 0.003 \times 0.1 = 0.00015$$

$$\Delta W_2 = \alpha \boxed{\delta^h} \cdot X_2 = 0.5 \times 0.003 \times 0.3 = 0.00045$$
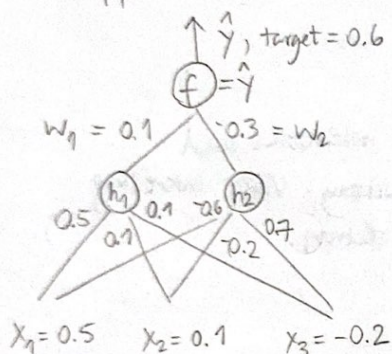
we see we have **VERY SMALL** steps! That's related to the **VANISHING GRADIENT** problem, here fueled by the sigmoid function w.r.t. f

weights in same layer have same error term!

other activation functions (eg. ReLU) alleviate this issue

④ Numpy exercise



$$\hat{y}, \text{ target} = 0.6$$

$W_1 = 0.1$      $0.3 = W_2$

$0.5 \quad 0.1 \quad 0.6$
$0.1 \quad 0.7$
$0.2$

$X_1 = 0.5 \qquad X_2 = 0.1 \qquad X_3 = -0.2$

$$\delta^o = (y-\hat{y}) \cdot f', \quad f' = f(1-f) = \hat{y}(1-\hat{y})$$

$$\delta^h = W_1 \cdot \delta^o \cdot f'(h_1) + W_2 \delta^o f'(h_2) =$$

$$= W_1 \delta^o f(h_1)(1-f(h_1)) + W_2 \delta^o f(h_2)(1-f(h_2))$$

⑤ **Backpropagation algorithm**

- Error term in output layer: $\delta_k^o = (y_k - \hat{y}_k) f'(a_k)$
- Error term in hidden layer: $\delta_j = \sum (W_{ij} \delta_k) \cdot f'(h_j)$

Set
. input to hidden weights: $\Delta W_{ij} = 0$
. hidden to output weights: $\Delta \overline{W}_j = 0$

for each training sample (a total of m samples):

$\hat{y} = \text{forward}(x)$

error term of output unit: $\delta^o = (y-\hat{y}) f'(z)$, $z = \sum_j \overline{W}_j a_j$ ← input of the output unit

propagate errors to hidden layers:

$$\delta_j^h = \delta^o \cdot \overline{W}_j \cdot f'(h_j)$$

update weight steps

$$\Delta \overline{W}_j = \Delta \overline{W}_j + \delta^o a_j$$

$$\Delta W_{ij} = \Delta W_{ij} + \delta_j^h a_i$$

then:

$$\overline{W}_j = \overline{W}_j + \frac{\alpha}{m} \Delta \overline{W}_j$$

$$W_{ij} = W_{ij} + (\alpha/m) \cdot \Delta W_{ij}$$

repeat for e epochs



$$\delta^o \qquad \hat{y} = f(z) = \sigma\left(\sum_j \overline{W}_j a_j\right)$$

$W_1 \qquad \overline{W}_j$
$W_2$
$a_1 \quad a_2 \quad a_j$

$$\delta_j^h$$

$W_{ij}$
$a_i$

9

TRAINING NEURAL NETWORKS ⊢▻ it's 1:1 the same as in the CVND: p. 20-24

SENTIMENT ANALYSIS by Andrew Trask: NLP PhD student at Oxford
author of Grokking deep learning
— We get a human generated text and want to know whether the content +/—.

⟶▻ Go to the repository and work there: deep-learning-v2-pytorch/sentiment-analysis
The section is divided in 6 mini-projects
↳ -network

25,000 movie reviews
labelled as positive/
negative sentiment

See Jupyter notebooks and repository!

Very interesting lesson: a neural network class is implemented
from-the-scratch using numpy for sentiment analysis.
A basic workflow for NLP/ text processing is shown.

DEEP LEARNING WITH PYTORCH ⊢▻ see the notebooks and
the repository. VERY interesting
lesson/section.

PROJECT: PREDICTING BIKE·SHARING PATTERNS