



POLITECNICO
MILANO 1863



Artificial Neural Networks and Deep Learning

- From Perceptrons to Feed Forward Neural Networks -

Matteo Matteucci, PhD (matteo.matteucci@polimi.it)

*Artificial Intelligence and Robotics Laboratory
Politecnico di Milano*

«Deep Learning is not AI, nor Machine Learning»

ARTIFICIAL INTELLIGENCE

Early artificial intelligence stirs excitement.



Neural Networks are as old as Artificial Intelligence



Source: Michael Copeland, *Deep Learnig Explained: What it is, and how it can deliver business value to your organization*



The inception of AI

A PROPOSAL FOR THE DARTMOUTH SUMMER RESEARCH PROJECT ON ARTIFICIAL INTELLIGENCE

J. McCarthy, Dartmouth College
M. L. Minsky, Harvard University
N. Rochester, I.B.M. Corp
C. E. Shannon, Bell Telephone

August 31, 1955

A Proposal for the DARTMOUTH SUMMER RESEARCH PROJECT ON ARTIFICIAL INTELLIGENCE

We propose that a 2 month, 10 man study of artificial intelligence be carried out during the summer of 1956 at Dartmouth College in Hanover, New Hampshire. The study is to proceed on the basis of the conjecture that every aspect of learning or any other feature of intelligence can in principle be so precisely described that a machine can be made to simulate it. An attempt will be made to find how to make machines use language, form abstractions and concepts, solve kinds of problems now reserved for humans, and improve themselves. We think that a significant advance can be made in one or more of these problems if a carefully selected group of scientists work on it together for a summer.

The following are some aspects of the artificial intelligence problem:

1) Automatic Computers

If a machine can do a job, then an automatic calculator can be programmed to simulate the machine. The speeds and memory capacities of present computers may be insufficient to simulate many of the higher functions of the human brain, but the major obstacle is not lack of machine capacity, but our inability to write programs taking full advantage of what we have.

3. Neuron Nets

How can a set of (hypothetical) neurons be arranged so as to form concepts. Considerable theoretical and experimental work has been done on this problem by Uttley, Rashevsky and his group, Farley and Clark, Pitts and McCulloch, Minsky, Rochester and Holland, and others. Partial results have been obtained but the problem needs more theoretical work.

5) Self-Improvement

Probably a truly intelligent machine will carry out activities which may best be described as self-improvement. Some schemes for doing this have been proposed and are worth further study. It seems likely that this question can be studied abstractly as well.

6) Abstractions

A number of types of "abstraction" can be distinctly defined and several others less distinctly. A direct attempt to classify these and to describe machine methods of forming abstractions from sensory and other data would seem worthwhile.



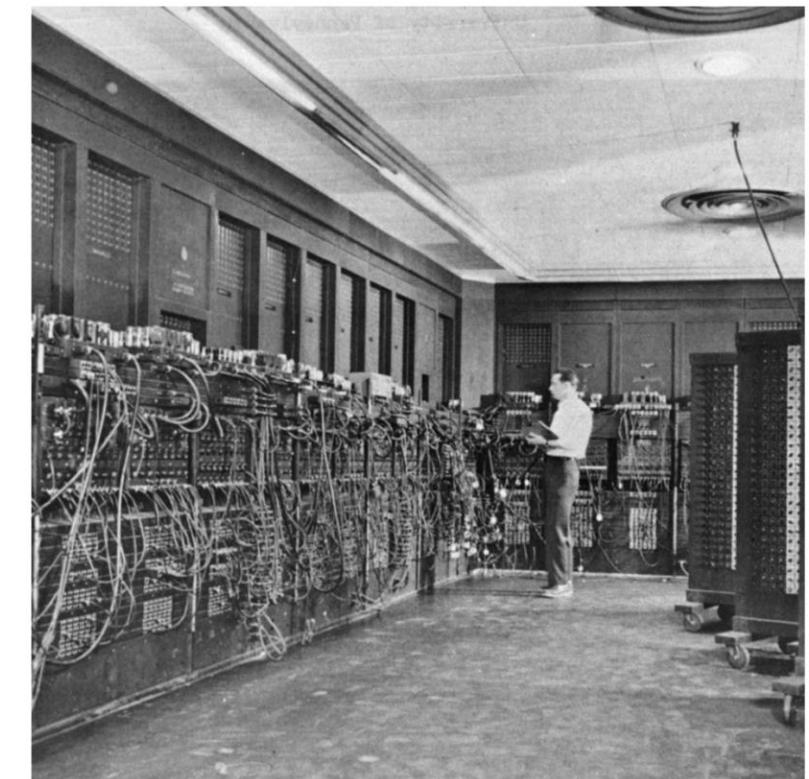
Let's go back to 1940s ...

Computers were already good at

- Doing precisely what the programmer programs them to do
- Doing arithmetic very fast

However, they would have liked them to:

- Interact with noisy data or directly with the environment
- Be massively parallel and fault tolerant
- Adapt to circumstances



Researchers were seeking a computational model beyond the Von Neumann Machine!

The Brain Computationa Model



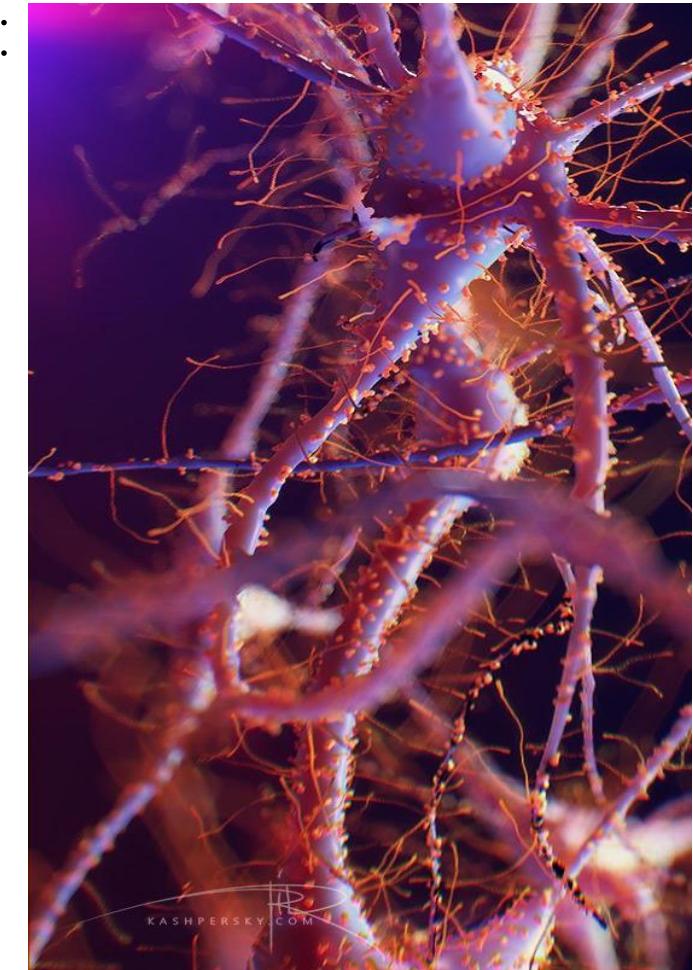
The human brain has a huge number of computing units:

- 10^{11} (one hundred billion) neurons
- 7,000 synaptic connections to other neurons
- In total from 10^{14} to 5×10^{14} (100 to 500 trillion) in adults to 10^{15} synapses (1 quadrillion) in a three year old child

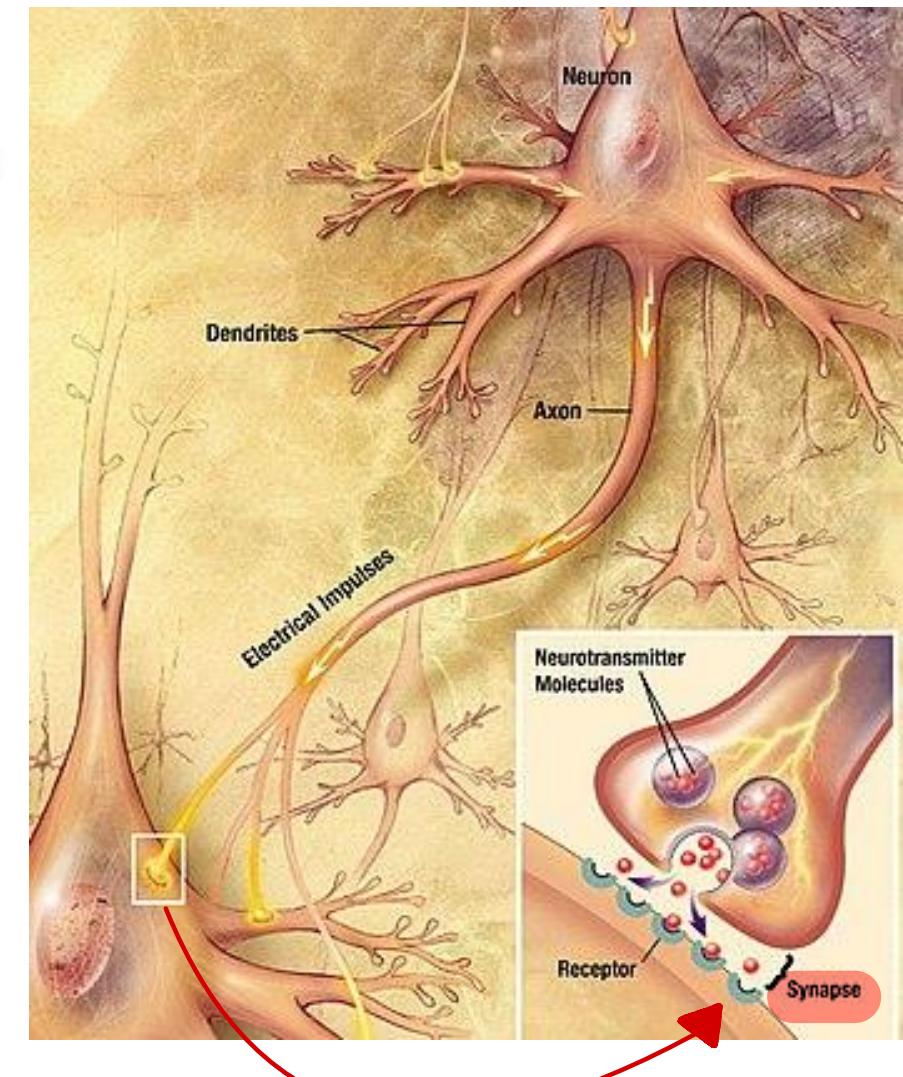
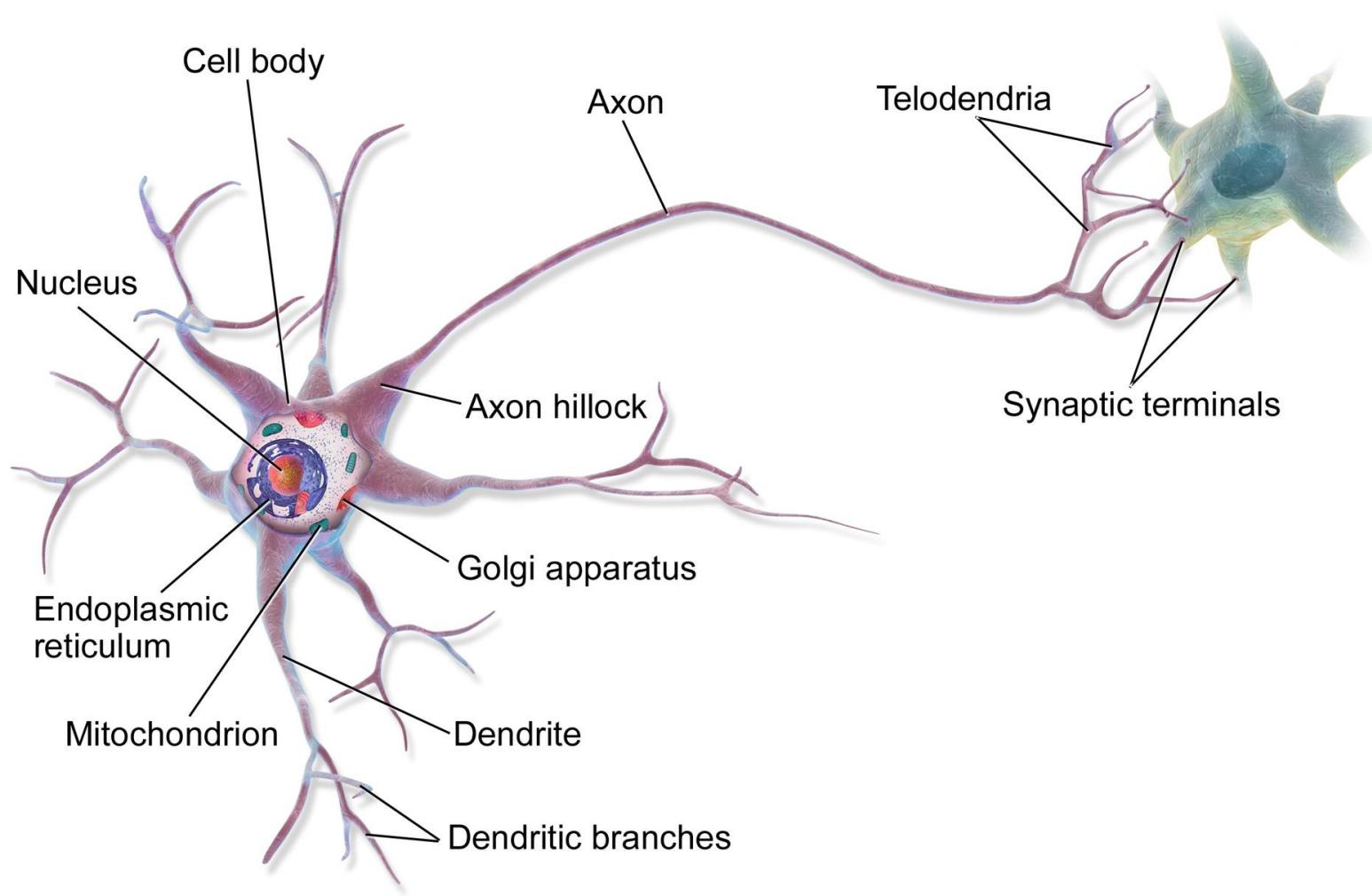
The computational model of the brain is:

- Distributed among simple non linear units
- Redundant and thus fault tolerant
- Intrinsically parallel

Perceptron: a computational model based on the brain!



Computation in Biological Neurons

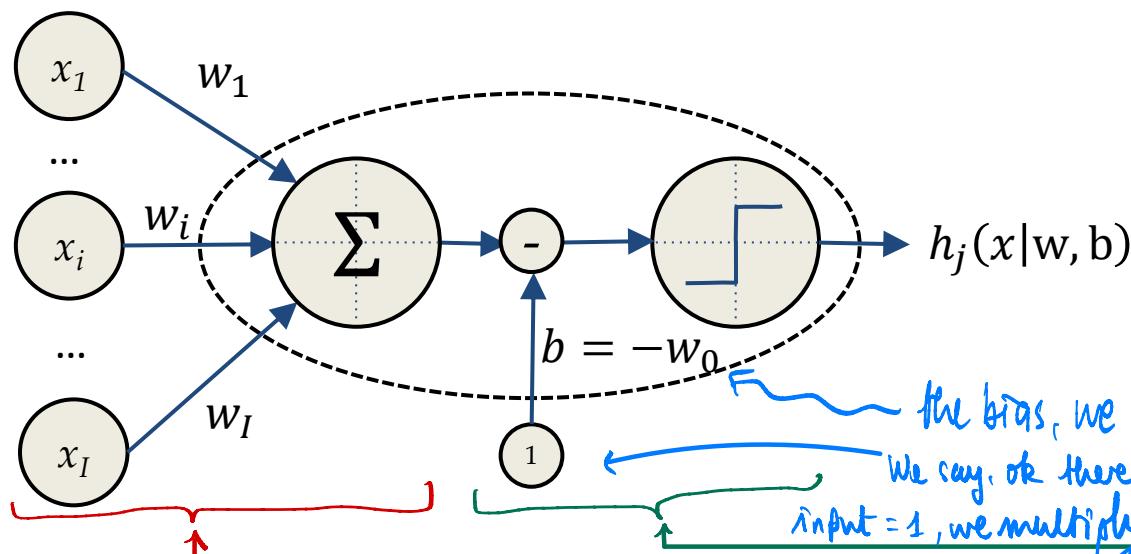


Computation in Artificial Neurons



Information is transmitted through chemical mechanisms:

- Dendrites collect charges from synapses, both Inhibitory and Excitatory
- Cumulates charge is released (neuron fires) once a Threshold is passed



We will always write it like this, but don't forget the bias!!

the bias, we call it $-w_0$.
We say, ok there is always an input = 1, we multiply it by $-w_0$: it gives us the bias.

$$h_j(x|w, b) = h_j\left(\sum_{i=1}^I w_i \cdot x_i - b\right) = h_j\left(\sum_{i=0}^I w_i \cdot x_i\right) = h_j(w^T x)$$

nonlinear function if > 0: charge is released. threshold = "bias" = $b = -w_0$

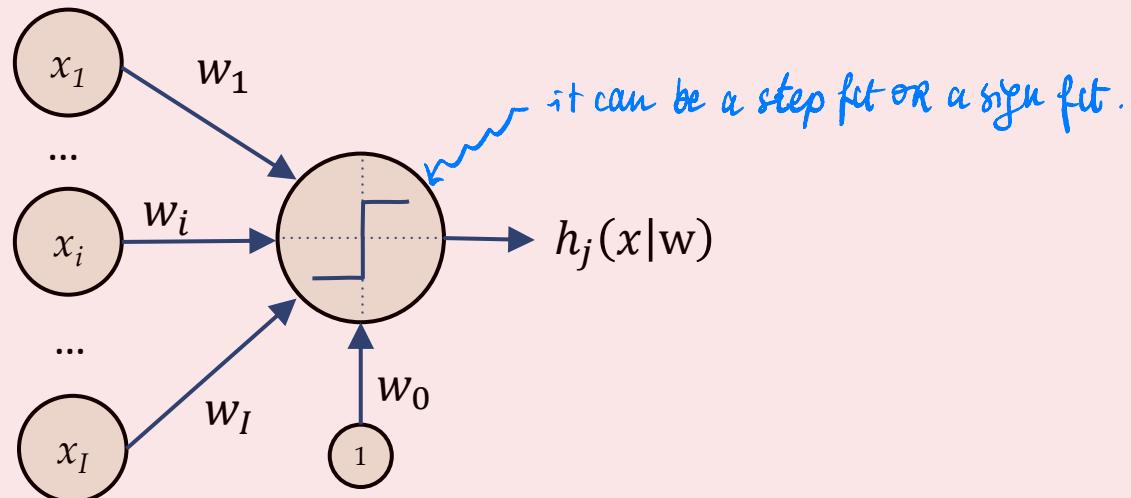
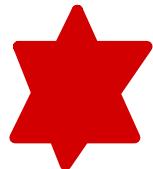
{ so using the notation $b = -w_0$, we can write it in a simpler form.



Computation in Artificial Neurons

Information is transmitted through chemical mechanisms:

- Dendrites collect charges from synapses, both Inhibitory and Excitatory
- Cumulates charge is released (neuron fires) once a Threshold is passed



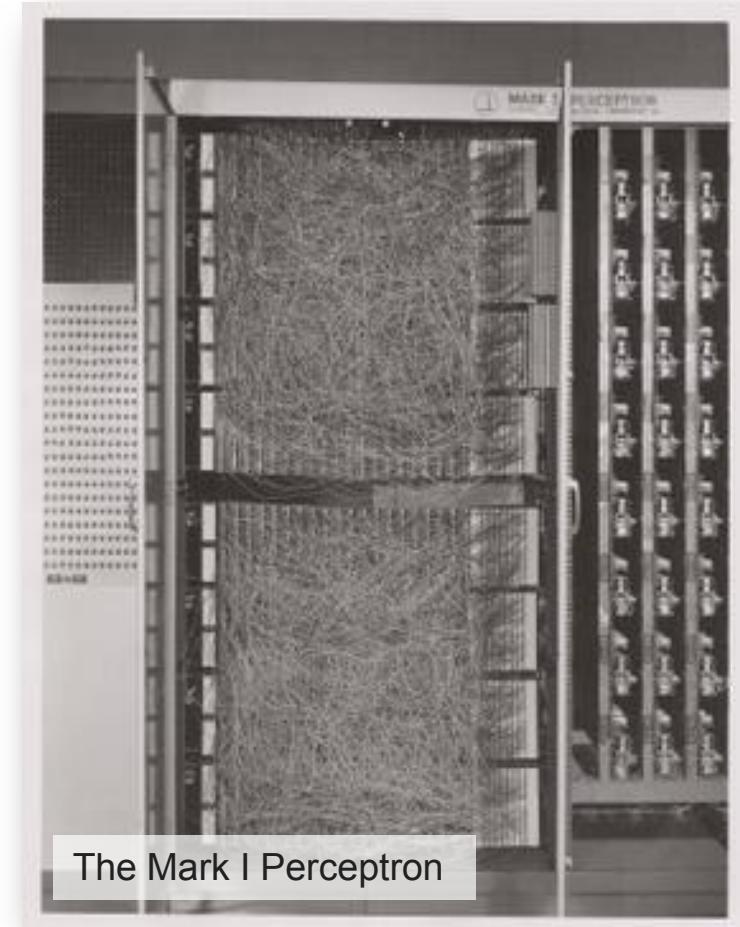
VERY
IMPORTANT
NOTATION

$$h_j(x|w, b) = h_j\left(\sum_{i=1}^I w_i \cdot x_i - b\right) = h_j\left(\sum_{i=0}^I w_i \cdot x_i\right) = h_j(w^T x)$$

Who did it first?

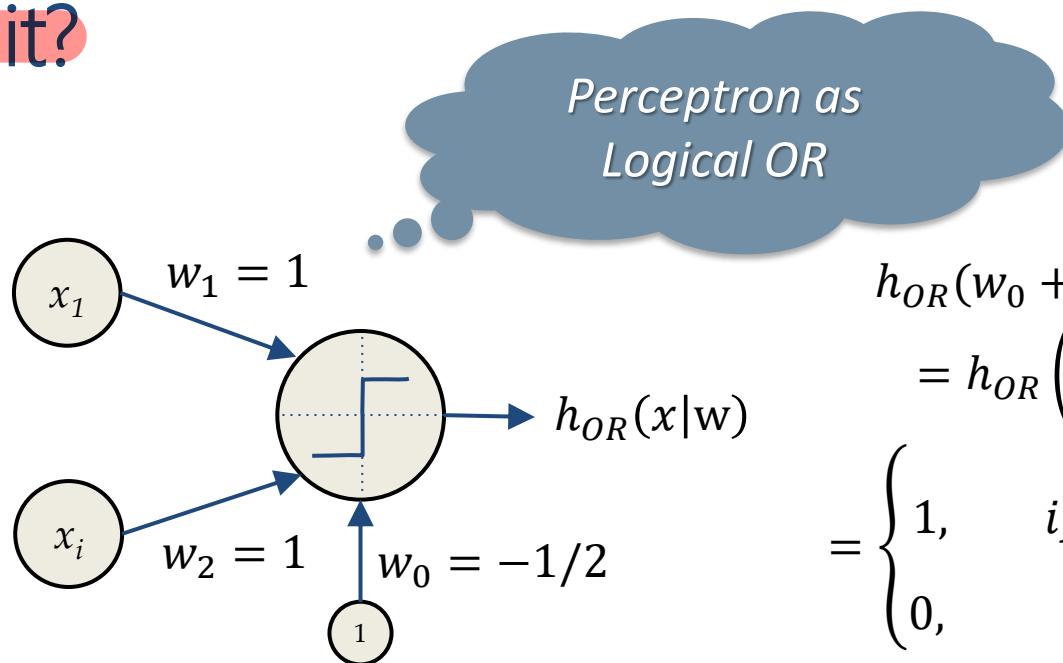
Several researchers were investigating models for the brain

- In 1943, Warren McCulloch and Walter Harry Pitts proposed the Threshold Logic Unit or Linear Unit, the activation function was a threshold unit equivalent to the Heaviside step function
- In 1957, Frank Rosenblatt developed the first Perceptron. Weights were encoded in potentiometers, and weight updates during learning were performed by electric motors
- In 1960, Bernard Widrow introduced the idea of representing the threshold value as a bias term in the ADALINE (Adaptive Linear Neuron or later Adaptive Linear Element)



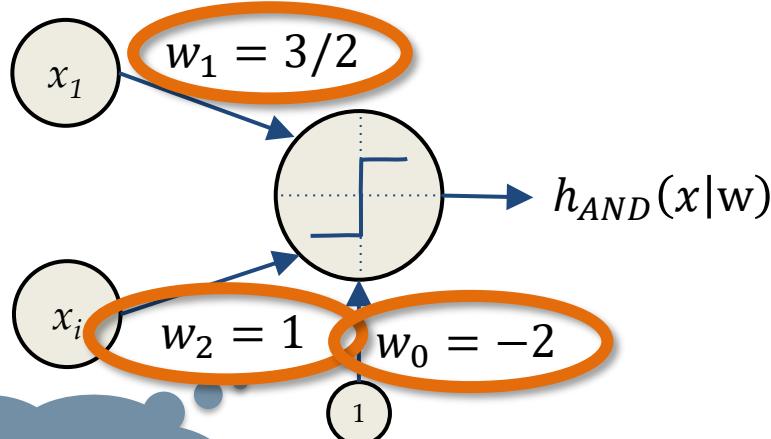
What can you do with it?

x_0	x_1	x_2	OR
1	0	0	0
1	0	1	1
1	1	0	1
1	1	1	1



$$\begin{aligned}
 h_{OR}(w_0 + w_1 \cdot x_1 + w_2 \cdot x_2) &= \\
 &= h_{OR}\left(-\frac{1}{2} + x_1 + x_2\right) = \\
 &= \begin{cases} 1, & \text{if } \left(-\frac{1}{2} + x_1 + x_2\right) > 0 \\ 0, & \text{otherwise} \end{cases}
 \end{aligned}$$

x_0	x_1	x_2	AND
1	0	0	0
1	0	1	0
1	1	0	0
1	1	1	1



$$\begin{aligned}
 h_{AND}(w_0 + w_1 \cdot x_1 + w_2 \cdot x_2) &= \\
 &= h_{AND}\left(-2 + \frac{3}{2} x_1 + x_2\right) = \\
 &= \begin{cases} 1, & \text{if } \left(-2 + \frac{3}{2} x_1 + x_2\right) > 0 \\ 0, & \text{otherwise} \end{cases}
 \end{aligned}$$

↑
input always = 1



Hebbian Learning



"The strength of a synapse increases according to the simultaneous activation of the relative input and the desired target"

(Donald Hebb, *The Organization of Behavior*, 1949)

Hebbian learning can be summarized by the following equations:

$$w_i^{k+1} = w_i^k + \Delta w_i^k$$

$$\Delta w_i^k = \eta \cdot x_i^k \cdot t^k$$

Where we have:

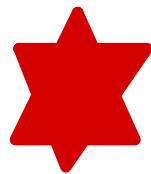
- η : learning rate
- x_i^k : the i^{th} perceptron input at time k
- t^k : the desired output at time k

Start from a random initialization

Fix the weights one sample at the time (online), and only if the sample is not correctly predicted



Perceptron Example

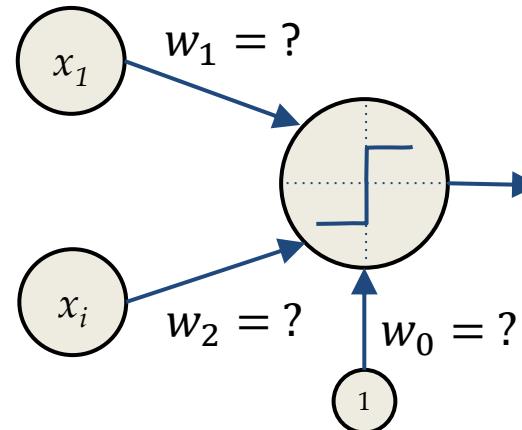


Life is easier with Sign(.)



Learn the weights to implement the OR operator

- Start from random weights, e.g.,
 $w = [0 \ 0 \ 0]$
- Chose a learning rate, e.g.,
 $\eta = 0.5$
- Cycle through the records by fixing those which are not correct
- End once all the records are correctly predicted



$$h(w^T x) = \begin{cases} 1 & \text{if } w^T x > 0 \\ 0 & \text{if } w^T x = 0 \\ -1 & \text{if } w^T x < 0 \end{cases}$$

Bias	x_0	x_1	x_2	OR
1	1	-1	-1	-1
1	1	-1	1	1
1	1	1	-1	1
1	1	1	1	1

Does the procedure always converge?

YES, if ... see slide n°13.

Does it always converge to the same sets of weights?

⚠ The weights depend on the order of records, η , & initialisation!



Let's do record #1: $\begin{cases} x_0 = 1 \\ x_1 = -1 \\ x_2 = -1 \end{cases}$. With the random weights: $w = [0, 0, 0]$. so we get:

$h(0+0+0) = h(0) = 0$. \rightarrow we wanted $-1 \rightarrow$ so we have to fix it.
 according to the def of h (of previous slide).

Fixing by hebbian learning: • $w_0 = w_0 + \Delta w_0 = w_0 + \gamma x_0 t = 0 + \frac{1}{2} \cdot 1 \cdot (-1) = -\frac{1}{2}$

• $w_1 = w_1 + \Delta w_1 = w_1 + \gamma x_1 t$
 $= 0 + \frac{1}{2} (-1)(-1) = +\frac{1}{2}$

• $w_2 = w_2 + \Delta w_2 = w_2 + \gamma x_2 t = 0 + \frac{1}{2} (1) (-1) = +\frac{1}{2}$

(*) to check that the updates we did with records 2,3,4...

did not screw up the answer

on record 1. And again, & again...

Many epochs \rightarrow here if we do it, after epoch 2 it's finished!

AND then we repeat the procedure, with the new w & the
 3 next records. At the end we get: $w = [\dots]$: 1 EPOCH again (*)

A pass through
the dataset

We have
to pass

Perceptron Math



don't forget that we
don't write the bias.

A perceptron computes a weighted sum, returns its Sign (Thresholding)

$$h_j(x|w) = h_j\left(\sum_{i=0}^I w_i \cdot x_i\right) = \text{Sign}(w_0 + w_1 \cdot x_1 + \cdots + w_I \cdot x_I)$$

It is a linear classifier for which the decision boundary is the hyperplane

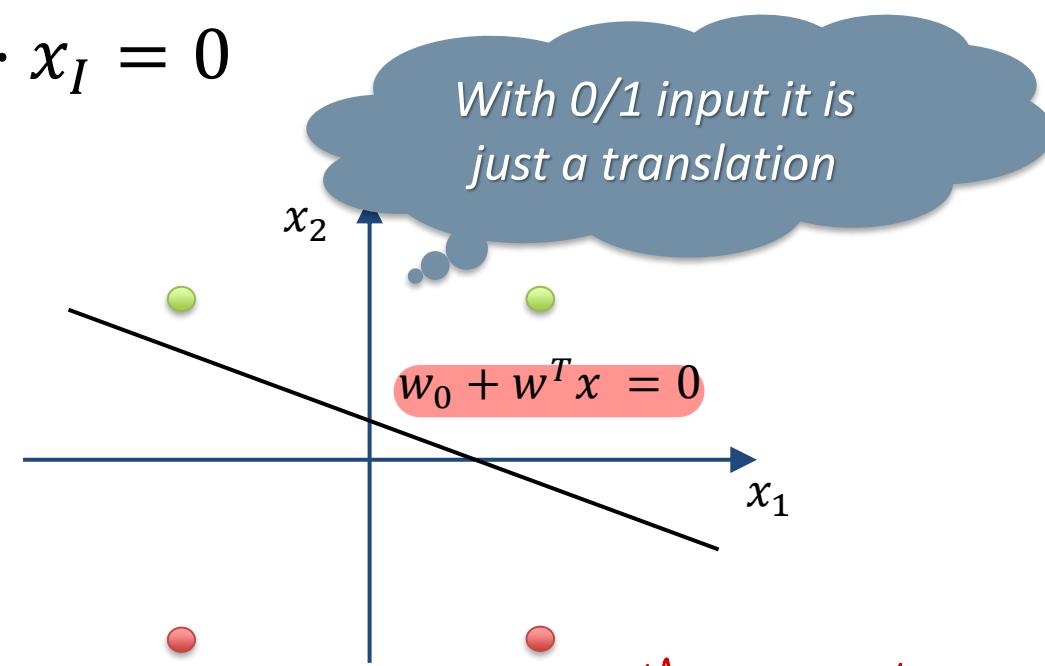
$$w_0 + w_1 \cdot x_1 + \cdots + w_I \cdot x_I = 0$$

In 2D, this turns into

$$w_0 + w_1 \cdot x_1 + w_2 \cdot x_2 = 0$$

$$w_2 \cdot x_2 = -w_0 - w_1 \cdot x_1$$

$$x_2 = -\frac{w_0}{w_2} - \frac{w_1}{w_2} \cdot x_1 : \text{it's a line.}$$



So the perceptron finds a line that keeps all the points >0 : above; & all <0 : below

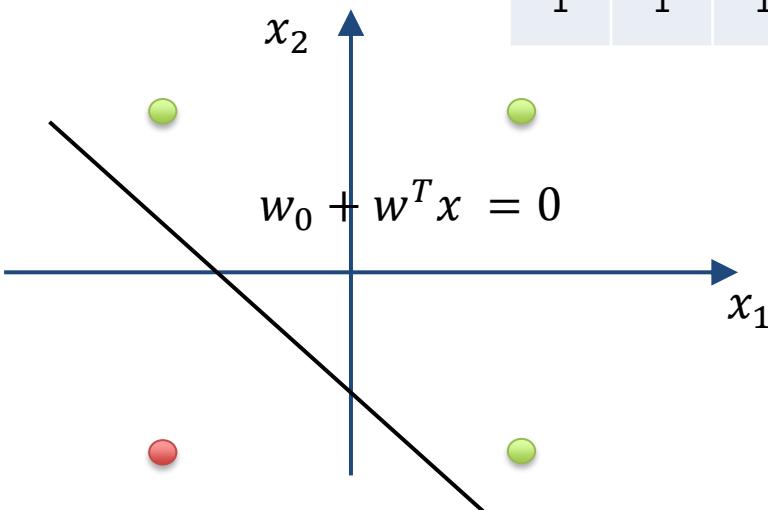


Boolean Operators are Linear Boundaries

What's about it? We had already Boolean operators

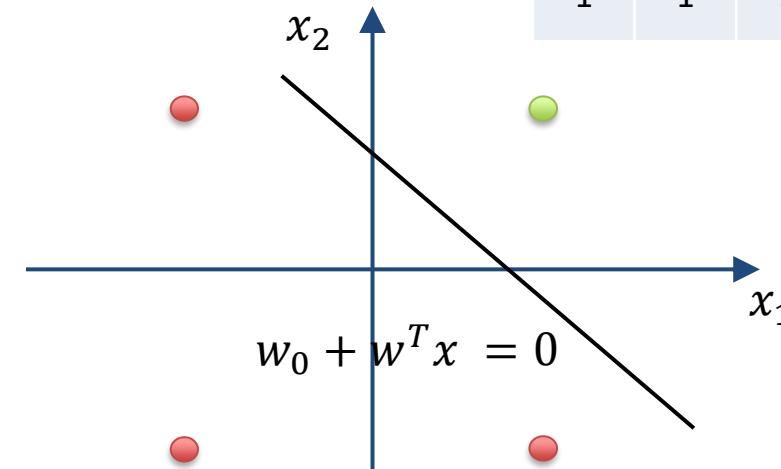
Linear boundary explains how Perceptron implements Boolean operators

x_0	x_1	x_2	OR
1	-1	-1	-1
1	-1	1	1
1	1	-1	1
1	1	1	1

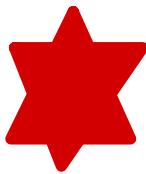


Yes! But this is a single trainable HW!

x_0	x_1	x_2	AND
1	-1	-1	-1
1	-1	1	-1
1	1	-1	-1
1	1	1	1



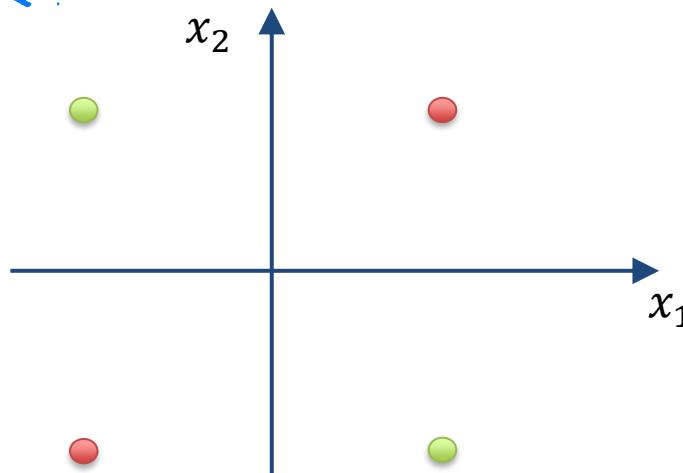
What can't you do with it?



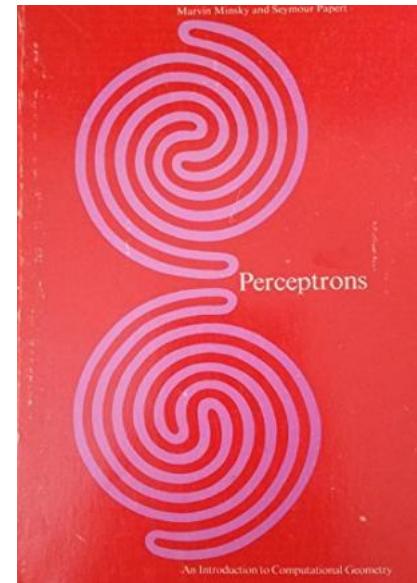
What if the dataset we want to learn does not have a linear separation boundary?

↳ Like the XOR.

x_0	x_1	x_2	XOR
1	-1	-1	-1
1	-1	1	1
1	1	-1	1
1	1	1	-1



Marvin Minsky, Seymour Papert
"Perceptrons: an introduction to computational geometry" 1969.

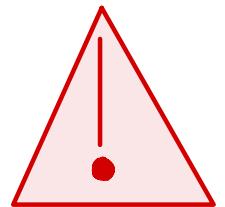


The Perceptron does not work any more and we need alternative solutions

- Non linear boundary
- Alternative input representations

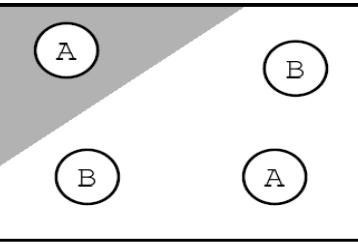
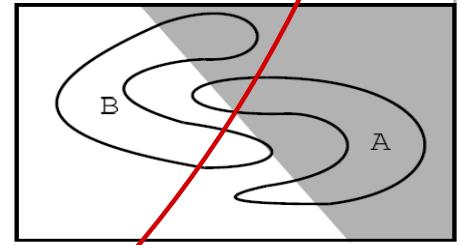
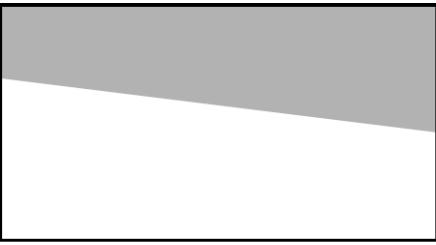
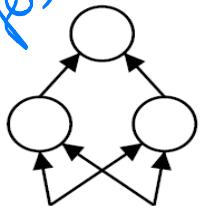
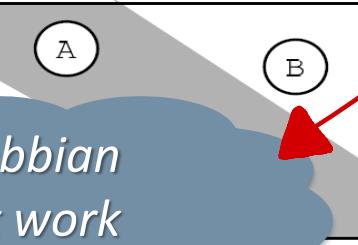
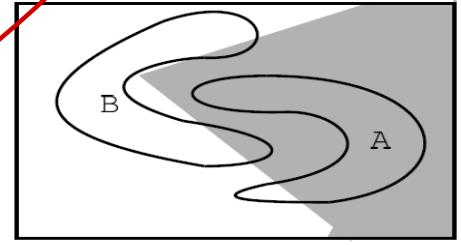
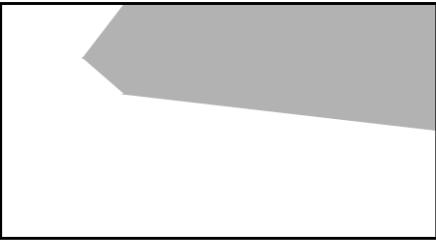
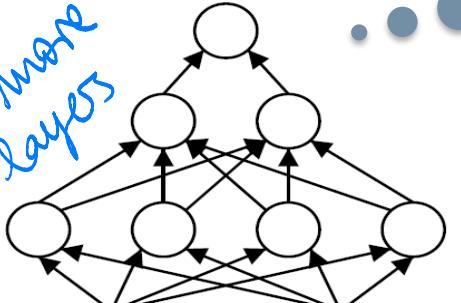
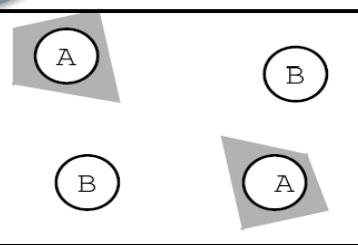
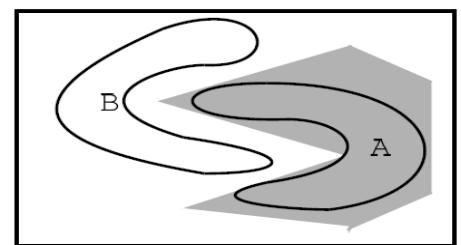
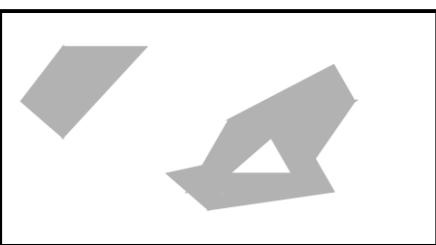


*The idea behind Multi
Layer Perceptrons*



What can't you do with it?

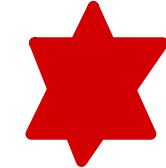
Difficulty : TRAINING

Topology	Type of Decision Region	XOR Problem	Classes with Meshed Regions	Most General Region Shapes
	Half bounded by hyperplanes			
 2 layers	Convex Open or <i>Unfortunately Hebbian learning does not work any more ...</i>			
 more layers	Arbitrary Regions (Complexity limited by the number of nodes)			

Layer Perceptrons

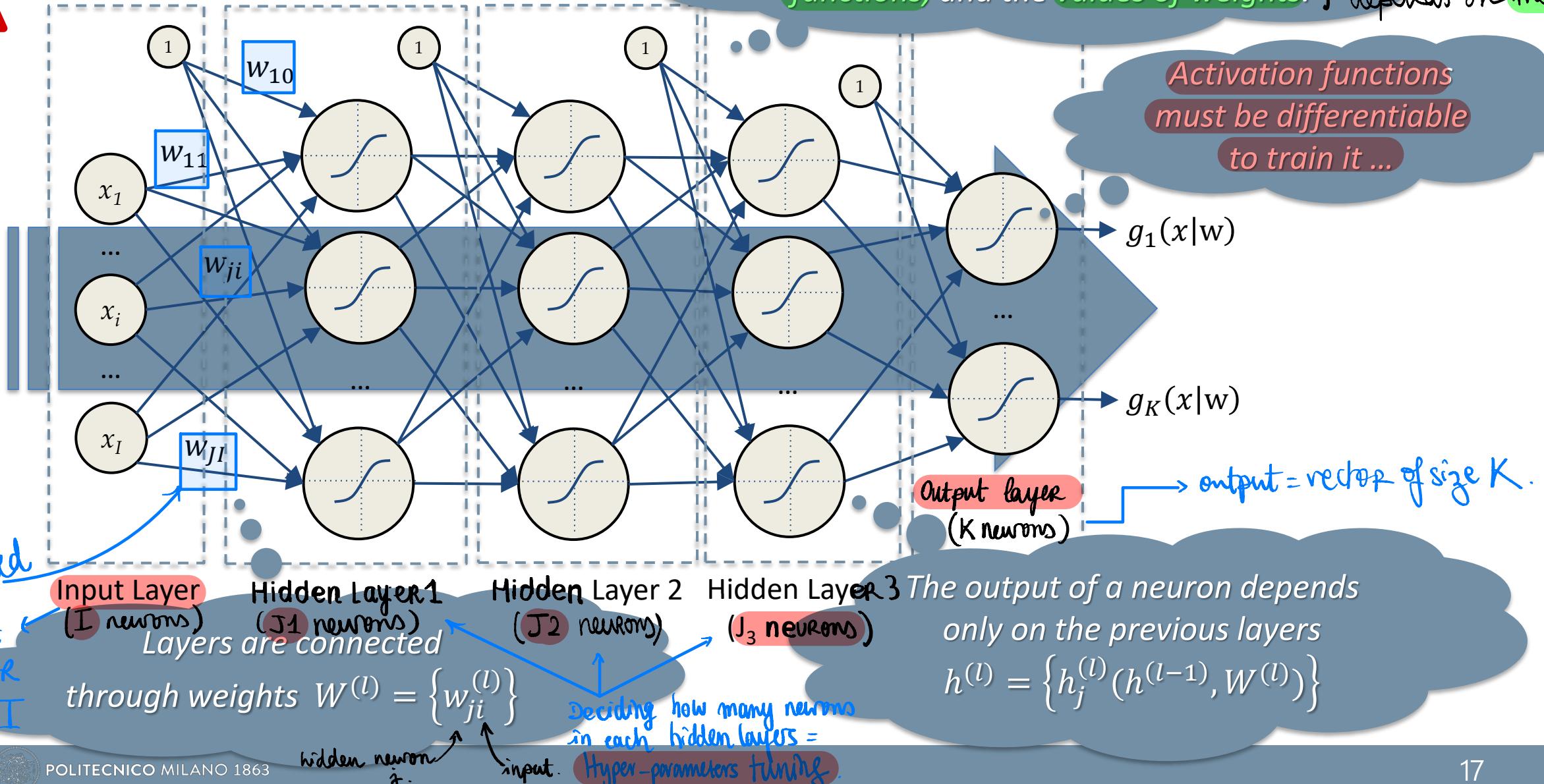


Feed Forward Neural Networks



Non-linear model characterized by the number of neurons, activation functions, and the values of weights.

} output vector is generated via a highly non-linear fn which depends on these.



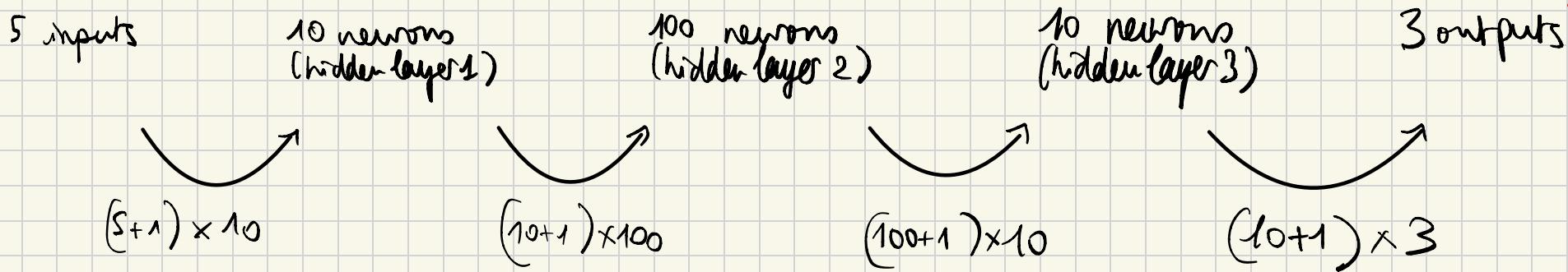


How many weights
in my network?

Between input
and H.L.1 :

$$J_x(I+1)$$

BIAS!



As you can notice (above) the nb of weights is quadratic:

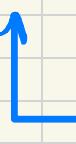
If you have H_L_i --- H_L_{i+1} } you get $(100+1) \times 100 \approx 100^2$ weights.
 $\frac{100}{100}$ neurons $\frac{100}{100}$ neurons

So the nb of parameters can be huge in those models.

STRUCTURE: EACH LAYER IS FULLY CONNECTED w/ THE PREVIOUS ONE .
 (or NEXT)



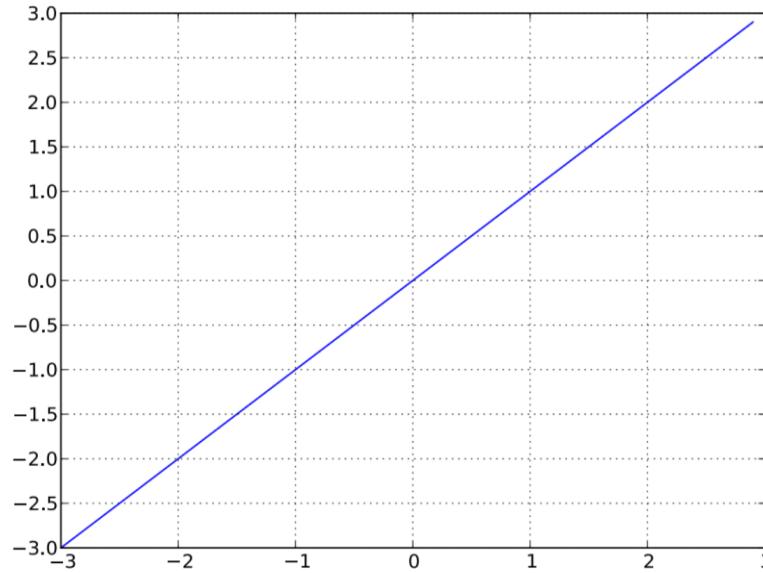
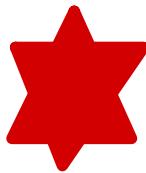
Here we work with FEED FORWARD NEURAL NETWORK.



Why do we like this? Because we have BACKPROPAGATION

Which Activation Function?

We will introduce other activation functions...

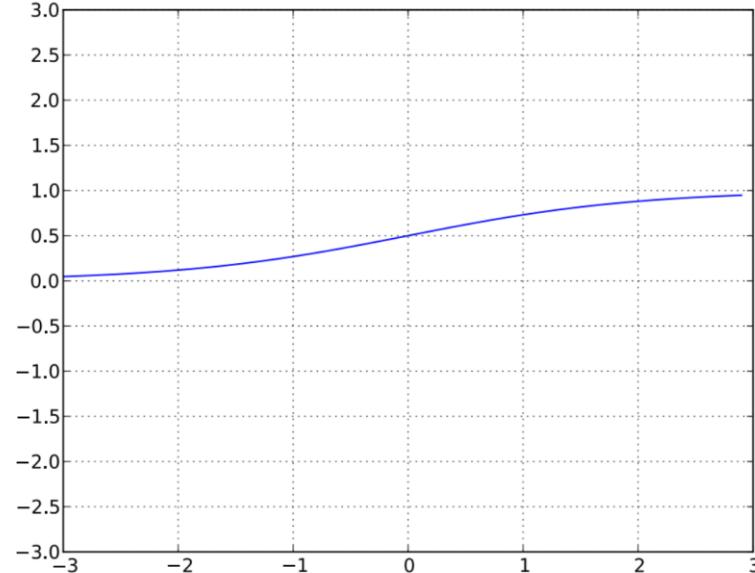


Linear activation function

Don't we
only linear act. fct!

$$g(a) = a$$
$$g'(a) = 1$$

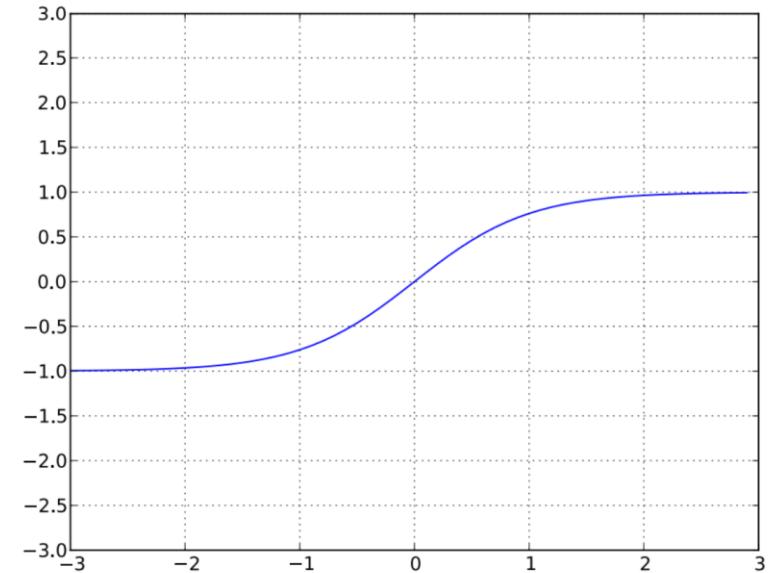
⚠ Problem: the final output will be a linear combination. That's why we need non-linear act. fct.



Sigmoid activation function

$$g(a) = \frac{1}{1 + \exp(-a)}$$
$$g'(a) = g(a)(1 - g(a))$$

One of the most used. easy to compute



Tanh activation function

$$g(a) = \frac{\exp(a) - \exp(-a)}{\exp(a) + \exp(-a)}$$
$$g'(a) = 1 - g(a)^2$$



Output Layer in Regression and Classification



Choice of activation fct for output layer in certain cases:

(1)

In **Regression** the output spans the whole \mathbb{R} domain:

- Use a Linear activation function for the output neuron

↳ tanh or sigmoid won't work for \mathbb{R} -output.

For all hidden neurons use sigmoid or tanh (see later)

(2)

In **Classification** with **two classes**, chose according to their coding:

- Two classes $\{\Omega_0 = -1, \Omega_1 = +1\}$ then use **Tanh output activation**
- Two classes $\{\Omega_0 = 0, \Omega_1 = 1\}$ then use **Sigmoid output activation**
(it can be interpreted as **class posterior probability**)

«One hot» encoding

(3)

When dealing with **multiple classes** (K) use **as many neuron as classes**

- Classes are coded as $\{\Omega_0 = [0 \ 0 \ 1], \Omega_1 = [0 \ 1 \ 0], \Omega_2 = [1 \ 0 \ 0]\}$

- Output neurons use a softmax unit

$$y_k = \frac{\exp(z_k)}{\sum_k \exp(z_k)} = \frac{\exp\left(\sum_j w_{kj} h_j (\sum_i^I w_{ji} \cdot x_i)\right)}{\sum_{k=1}^K \exp\left(\sum_j w_{kj} h_j (\sum_i^I w_{ji} \cdot x_i)\right)}$$

amplify the differences.
→ normalize.

(1) Regression (\mathbb{R}) ; (2) Classification with 2 classes ; (3) Classification with K classes.



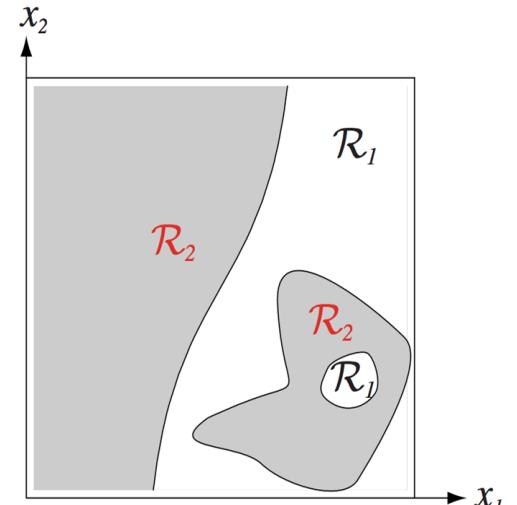
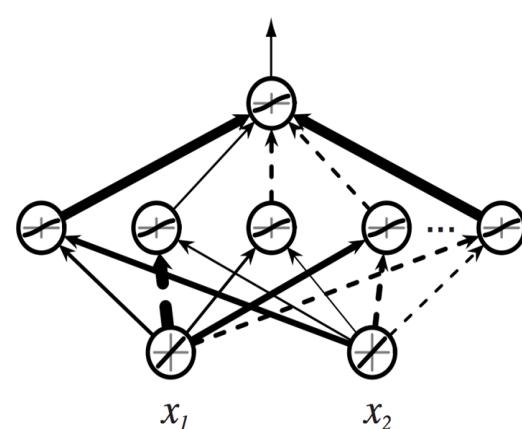
Neural Networks are Universal Approximators



"A single hidden layer feedforward neural network with S shaped activation functions can approximate any measurable function to any desired degree of accuracy on a compact set "

& it can be
hard...

Universal approximation theorem
(Kurt Hornik, 1991)



Images from Hugo Larochelle's DL Summer School Tutorial

Regardless the function we are learning, a single layer can represent it:

ie we need
to train it!

- Doesn't mean a learning algorithm can find the necessary weights!
- In the worse case, an exponential number of hidden units may be required ...
- The layer may have to be unfeasibly large and may fail to learn and generalize

Classification requires just one extra layer ... : so 2 hidden layers.

VSUAL/BASIC ARCHITECTURE.

So, firstly, we would like to learn how to find the weights. Then, we'll try to do it without overfitting!



Optimization and Learning (Supervised learning)

But @ the beginning they
didn't realize...

We will train using GRADIENT DESCENT (which is equivalent to backpropagation).

Recall learning a parametric model $y(x_n|\theta)$ in regression/classification

- Given a training set

$$D = \langle x_1, t_1 \rangle \dots \langle x_N, t_N \rangle$$

- We want to find model parameters such that for new data

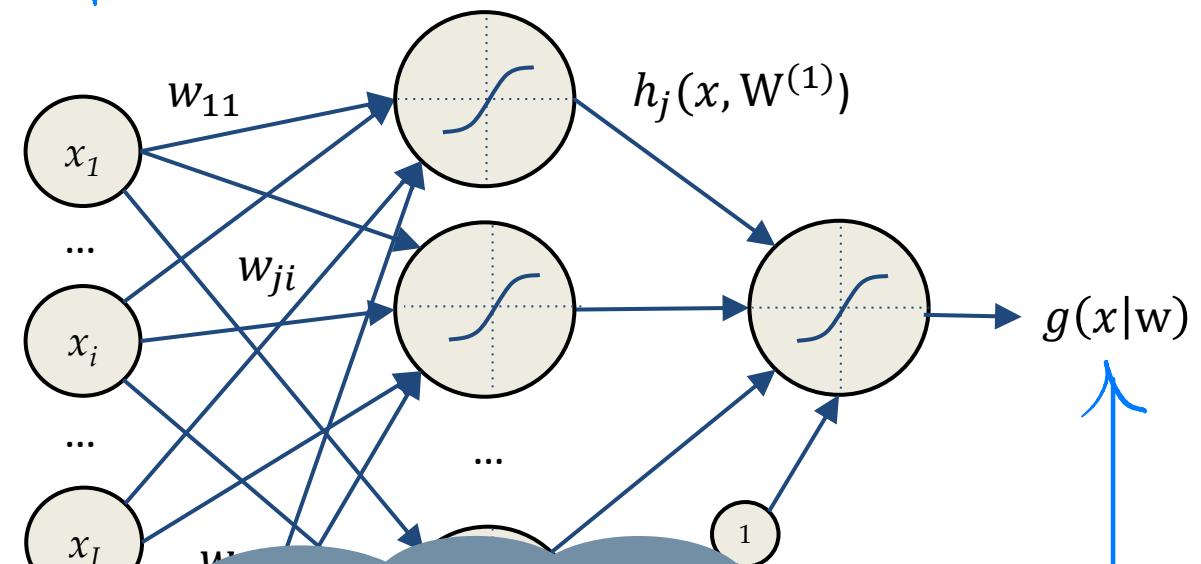
$$\underline{y(x_n|\theta) \sim t_n}$$

- In case of a Neural Network this can be rewritten as

$$\underline{g(x_n|w) \sim t_n}$$

output of our NN

↑
our input.
some set of parameters



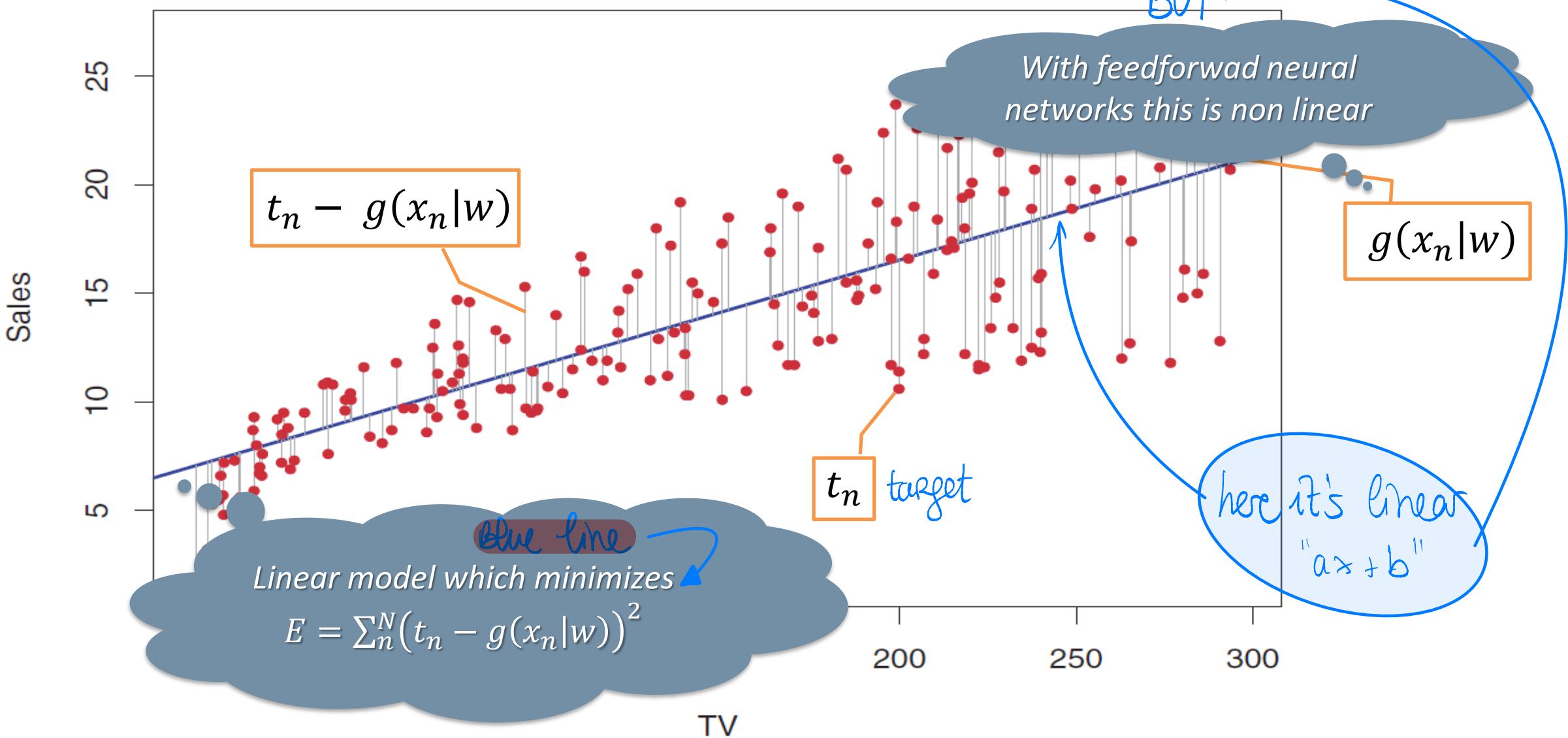
For this you can minimize

$$E = \sum_n^N (t_n - g(x_n|w))^2$$

EROR / LOSS.

} That's why we need differentiable functions.

Sum of Squared Errors



Non Linear Optimization 101



To find the minimum of a generic function, we compute the partial derivatives of the function and set them to zero

$$\rightarrow \frac{\partial J(w)}{\partial w} = 0$$

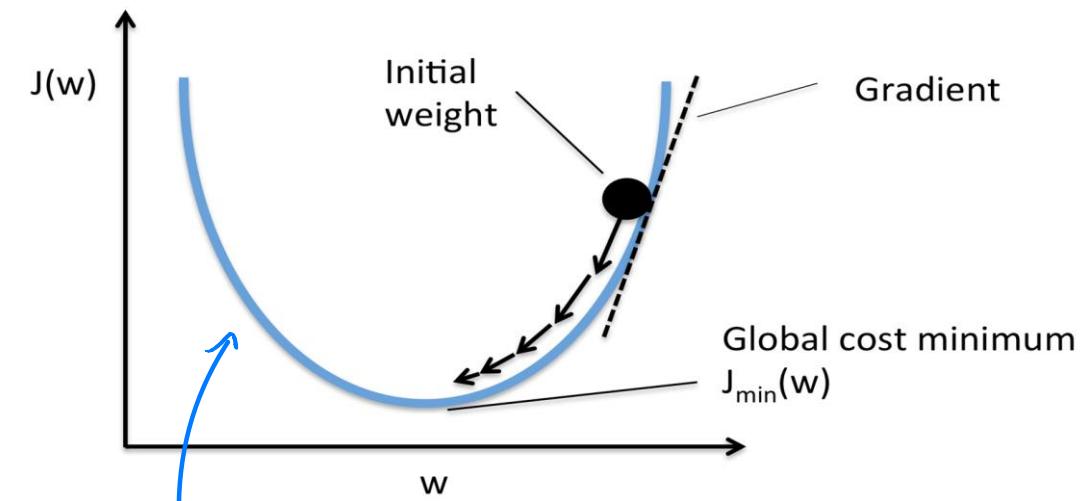
with highly non-linear fcts:
That's why we use grad. descent.

Closed-form solutions are practically never available so we can use iterative solutions (gradient descent):

- GD**
- 1) Initialize the weights to a random value
 - 2) Iterate until convergence

$$w^{k+1} = w^k - \eta \frac{\partial J(w)}{\partial w} \Big|_{w^k}$$

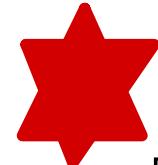
↑
Learning Rate.



that's a very "nice" shape: IRL, it's harder
(of next slide)



Gradient descent - Backpropagation



Finding the weights of a Neural Network is a

Use multiple restarts to seek for a proper global minimum.

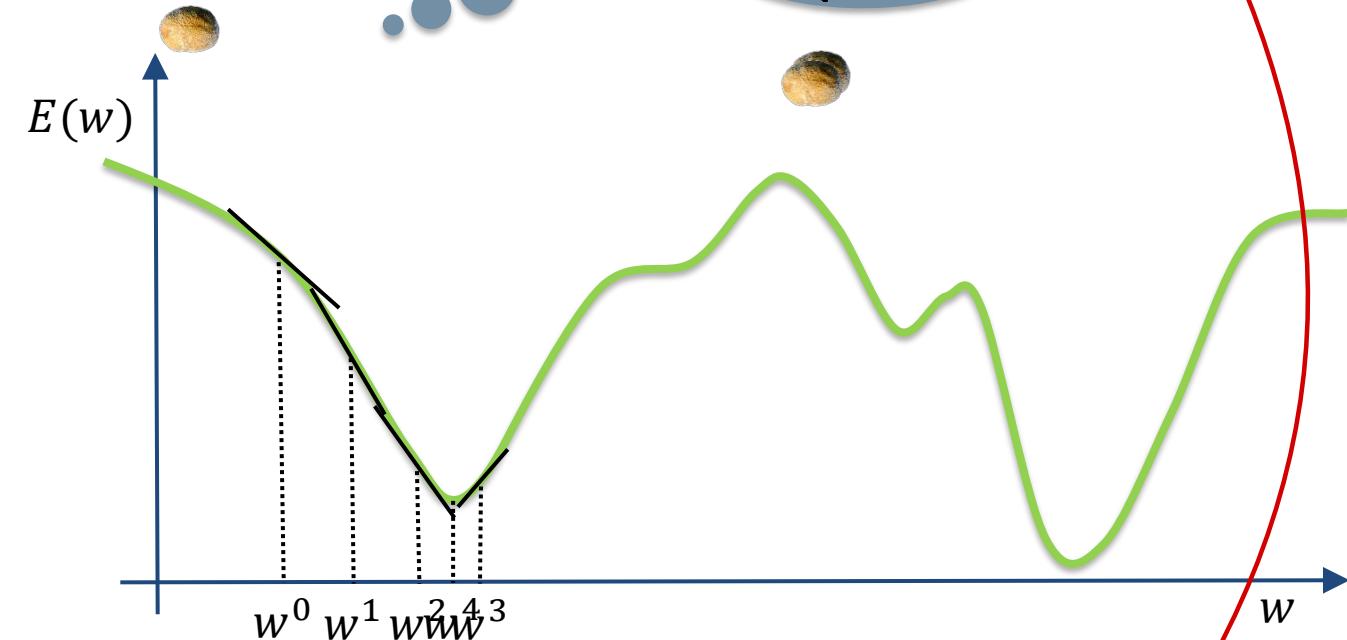
non-linear optimization

$$\operatorname{argmin}_w E(w) = \sum_{n=1}^N (t_n - g(x_n, w))^2$$

We iterate starting from an initial random configuration :

Iteration:

$$w^{k+1} = w^k - \eta \frac{\partial E(w)}{\partial w} \Big|_{w^k}$$



To avoid local minima can use momentum :

Like "ADAM":

$$w^{k+1} = w^k - \eta \frac{\partial E(w)}{\partial w} \Big|_{w^k} - \alpha \frac{\partial E(w)}{\partial w} \Big|_{w^{k-1}}$$

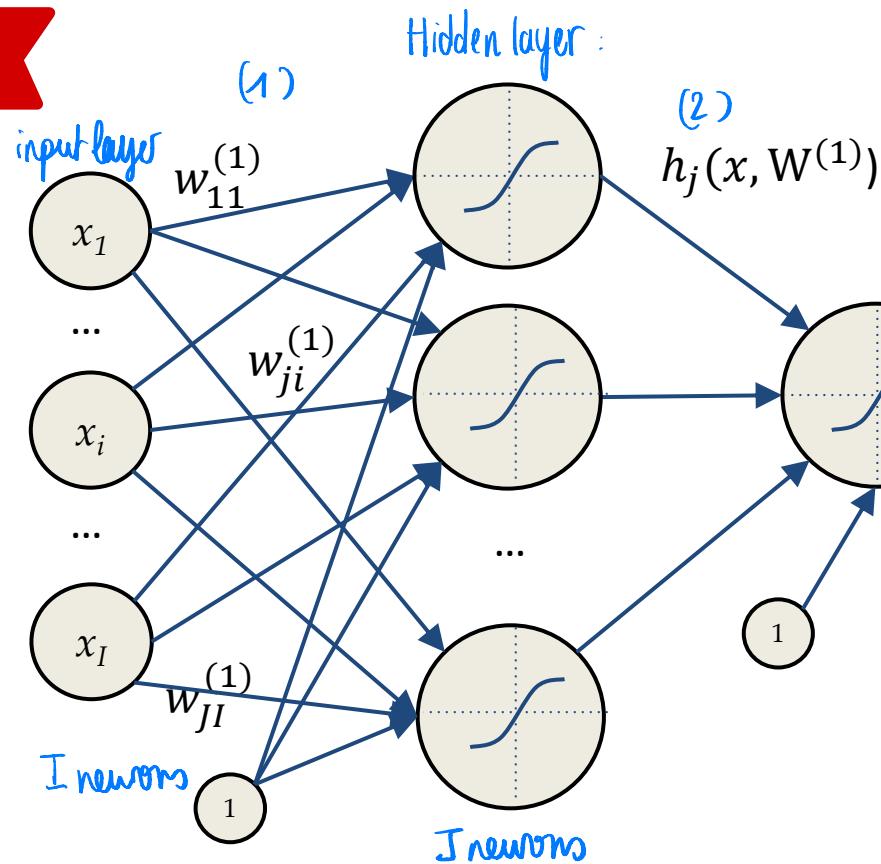
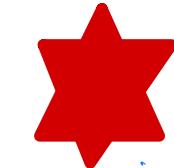
It depends on where we start from



You can use fast gradient value(s) ...



Gradient Descent Example



This is the function computed by the NN: "output function".

$$g_1(x_n|w) = g_1 \left(\sum_{j=0}^J w_{1j}^{(2)} \cdot h_j \left(\sum_{i=0}^I w_{ji}^{(1)} \cdot x_{i,n} \right) \right)$$

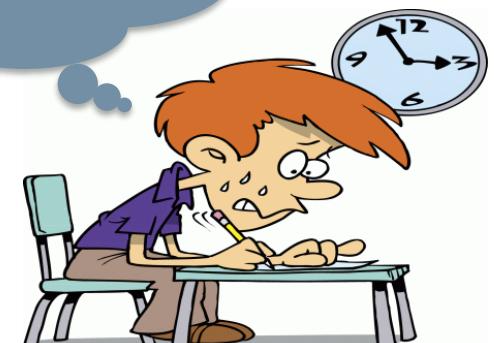
$$E(w) = \sum_{n=1}^N (t_n - g_1(x_n, w))^2$$

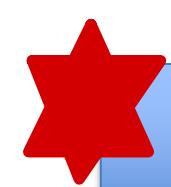
Use $j=3$ and $i=5$

$w_{3,5}$: $i=5 \rightarrow j=3$

Compute the $w_{ji}^{(1)}$ weight update formula by gradient descent

Exercise of next slide





$$g_1(x_n|w) = g_1\left(\sum_{j=0}^J w_{1j}^{(2)} \cdot h_j\left(\sum_{i=0}^I w_{ji}^{(1)} \cdot x_{i,n}\right)\right)$$

our starting points →

$$E(w) = \sum_{n=1}^N (t_n - g_1(x_n, w))^2$$

$$\frac{\partial E(w)}{\partial w_{3,5}^{(1)}} = \frac{\partial \sum_{n=1}^N (t_n - g_1(x_n, w))^2}{\partial w_{3,5}^{(1)}} = \sum_{n=1}^N \frac{\partial (t_n - g_1(x_n, w))^2}{\partial w_{3,5}^{(1)}} = -2 \sum_{n=1}^N (t_n - g_1(x_n, w)) \frac{\partial g_1(x_n, w)}{\partial w_{3,5}^{(1)}}$$

derivation formula: $(u^2)' = 2u \cdot u'$.

$$\frac{\partial g_1(x_n, w)}{\partial w_{3,5}^{(1)}} = \frac{\partial g_1\left(\sum_{j=0}^J w_{1j}^{(2)} \cdot h_j(.)\right)}{\partial w_{3,5}^{(1)}} = \underbrace{g'_1(x_n, w)}_{\text{"g' (f(x))"}} \cdot \underbrace{\frac{\partial \sum_{j=0}^J w_{1j}^{(2)} \cdot h_j(.)}{\partial w_{3,5}^{(1)}}}_{\text{"f' (x)"}} = g'_1(x_n, w) \cdot w_{1,3}^{(2)} \cdot \underbrace{\frac{\partial h_3\left(\sum_{i=0}^I w_{3i}^{(1)} \cdot x_{i,n}\right)}{\partial w_{3,5}^{(1)}}}_{\text{see previous page.}}$$

$$\frac{\partial h_3\left(\sum_{i=0}^I w_{3i}^{(1)} \cdot x_{i,n}\right)}{\partial w_{3,5}^{(1)}} = h'_3\left(\sum_{i=0}^I w_{3i}^{(1)} \cdot x_{i,n}\right) \frac{\partial \sum_{i=0}^I w_{3i}^{(1)} \cdot x_{i,n}}{\partial w_{3,5}^{(1)}} = h'_3\left(\sum_{i=0}^I w_{3i}^{(1)} \cdot x_{i,n}\right) x_{5,n}$$

only 1 term depends on $w_{3,5}^{(1)}$.

Gradient part of the backpropagation formula for $w_{3,5}^{(1)}$:

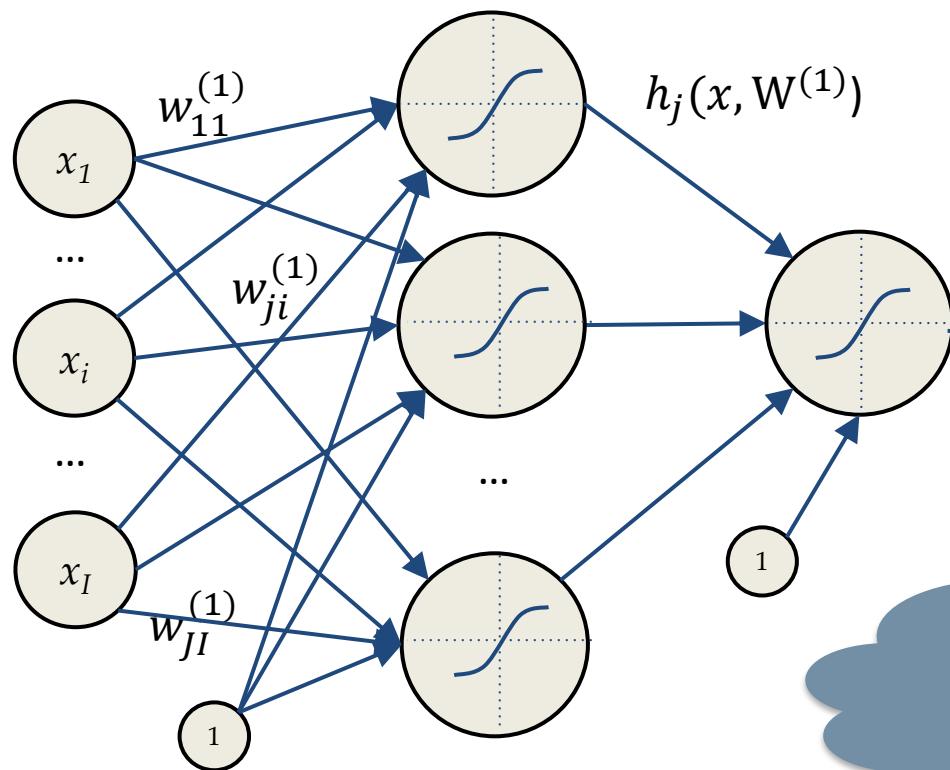
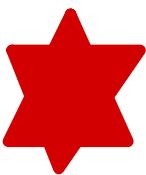
$$\frac{\partial E(w)}{\partial w_{3,5}^{(1)}} = -2 \sum_{n=1}^N (t_n - g_1(x_n, w)) g'_1(x_n, w) w_{1,3}^{(2)} h'_3\left(\sum_{i=0}^I w_{3i}^{(1)} \cdot x_{i,n}\right) x_{5,n}$$

For all weights of (1), you get the same:

you just have to change indexes (like 26 $w_{1,5}^{(2)}$)



Gradient Descent Example



$$g_1(x_n|w) = g_1\left(\sum_{j=0}^J w_{1j}^{(2)} \cdot h_j\left(\sum_{i=0}^I w_{ji}^{(1)} \cdot x_{i,n}\right)\right)$$

$$E(w) = \sum_{n=1}^N (t_n - g_1(x_n, w))^2$$

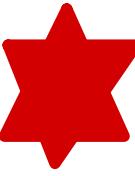
Using all the data points (**batch**) might be unpractical

sometimes all the data CANNOT FIT in the RAM of the GPU.

$$\frac{\partial E(w_{ji}^{(1)})}{\partial w_{ji}^{(1)}} = -2 \sum_n^n (t_n - g_1(x_n, w)) g'_1(x_n, w) w_{1j}^{(2)} h'_j\left(\sum_{i=0}^I w_{ji}^{(1)} \cdot x_{i,n}\right) x_i$$

cf next slide.

Gradient Descent Variations



Batch gradient descent : to solve the problem (of previous slide), we can do 1 sample @ a time.

$$\frac{\partial E(w)}{\partial w} = \frac{1}{N} \sum_n^N \frac{\partial E(x_n, w)}{\partial w}$$

Stochastic gradient descent (SGD)

$$\frac{\partial E(w)}{\partial w} \approx \frac{\partial E_{SGD}(w)}{\partial w} = \frac{\partial E(x_n, w)}{\partial w}$$

Mini-batch gradient descent

↑
Not the full
dataset: a
subset.

$$\frac{\partial E(w)}{\partial w} \approx \frac{\partial E_{MB}(w)}{\partial w} = \frac{1}{M} \sum_{n \in \text{Minibatch}}^{M < N} \frac{\partial E(x_n, w)}{\partial w}$$

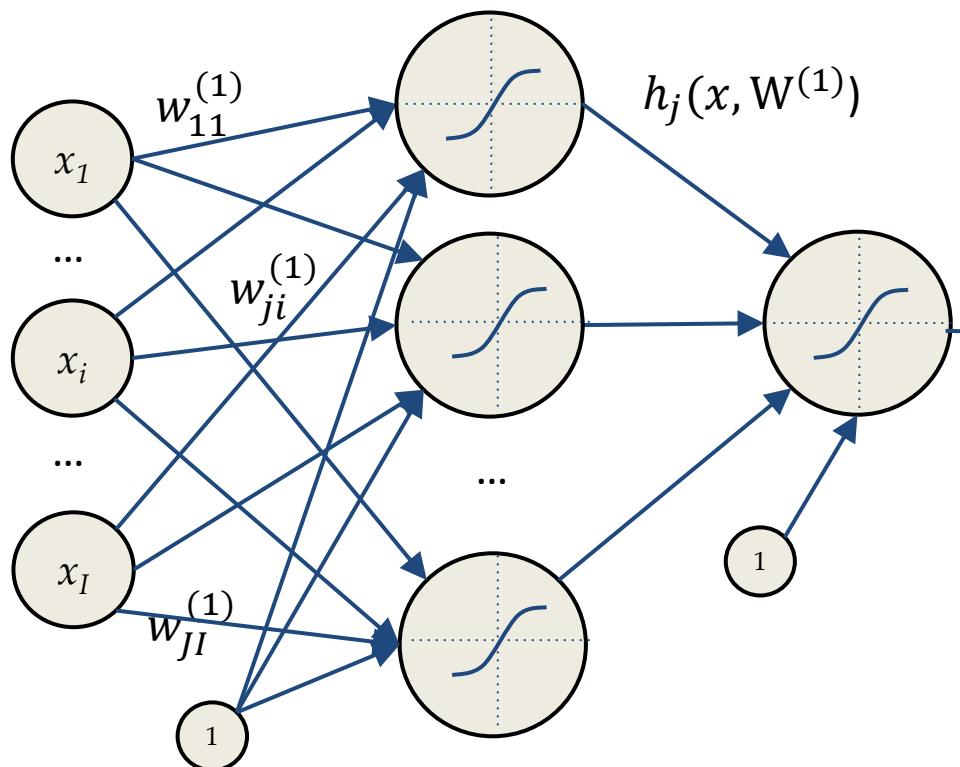
Use a single sample,
unbiased, but with
high variance

Use a subset of
samples, good trade off
variance-computation



Gradient Descent Example

Do I have to compute the gradient myself?



$$g_1(x_n|w) = g_1\left(\sum_{j=0}^J w_{1j}^{(2)} \cdot h_j\left(\sum_{i=0}^I w_{ji}^{(1)} \cdot x_{i,n}\right)\right)$$

$$E(w) = \sum_{n=1}^N (t_n - g_1(x_n, w))^2$$

Can I make it
automatic?

$$\frac{\partial E(w_{ji}^{(1)})}{\partial w_{ji}^{(1)}} = -2 \sum_n^n (t_n - g_1(x_n, w)) g'_1(x_n, w) w_{1j}^{(2)} h'_j\left(\sum_{i=0}^I w_{ji}^{(1)} \cdot x_{i,n}\right) x_i$$

Backpropagation and Chain Rule (1)



We have "weighted sum of non-linear functions of weighted sum...": NESTED Structure.



Weights update can be done in parallel, locally, and requires just 2 passes

- Let x be a real number and two functions $f: \mathbb{R} \rightarrow \mathbb{R}$ and $g: \mathbb{R} \rightarrow \mathbb{R}$
- Consider the composed function $z = f(g(x)) = f(y)$ where $y = g(x)$
- The derivative of z w.r.t. x can be computed applying the chain rule

$$\frac{dz}{dx} = \frac{dz}{dy} \frac{dy}{dx} = f'(y)g'(x) = f'(g(x))g'(x)$$

$$\frac{d f(g(x))}{dx} = g'(x) \times f'(g(x)) = \frac{dy}{dx} \times \frac{df(g(x))}{dg(x)} = \frac{dy}{dx} \times \frac{dz}{dy}$$

→ The same holds for backpropagation

$$\frac{\partial E(w_{ji}^{(1)})}{\partial w_{ji}^{(1)}} = -2 \sum_n^N (t_n - g_1(x_n, w)) \cdot g'_1(x_n, w) \cdot w_{1j}^{(2)} \cdot h'_j \left(\sum_{i=0}^I w_{ji}^{(1)} \cdot x_{i,n} \right) \cdot x_i$$

↓ ↓ ↓ ↓ ↓ ↓

$$\frac{\partial E}{\partial w_{ji}^{(1)}} \quad \frac{\partial E}{\partial g(x_n, w)} \quad \frac{\partial g(x_n, w)}{\partial w_{1j}^{(2)} h_j(.)} \quad \frac{\partial w_{1j}^{(2)} h_j(.)}{\partial h_j(.)} \quad \frac{\partial h_j(.)}{\partial w_{ji}^{(1)} x_i} \quad \frac{\partial w_{ji}^{(1)} x_i}{\partial w_{ji}^{(1)}}$$

If we want to change one activation fn: we don't need to recompute everything.

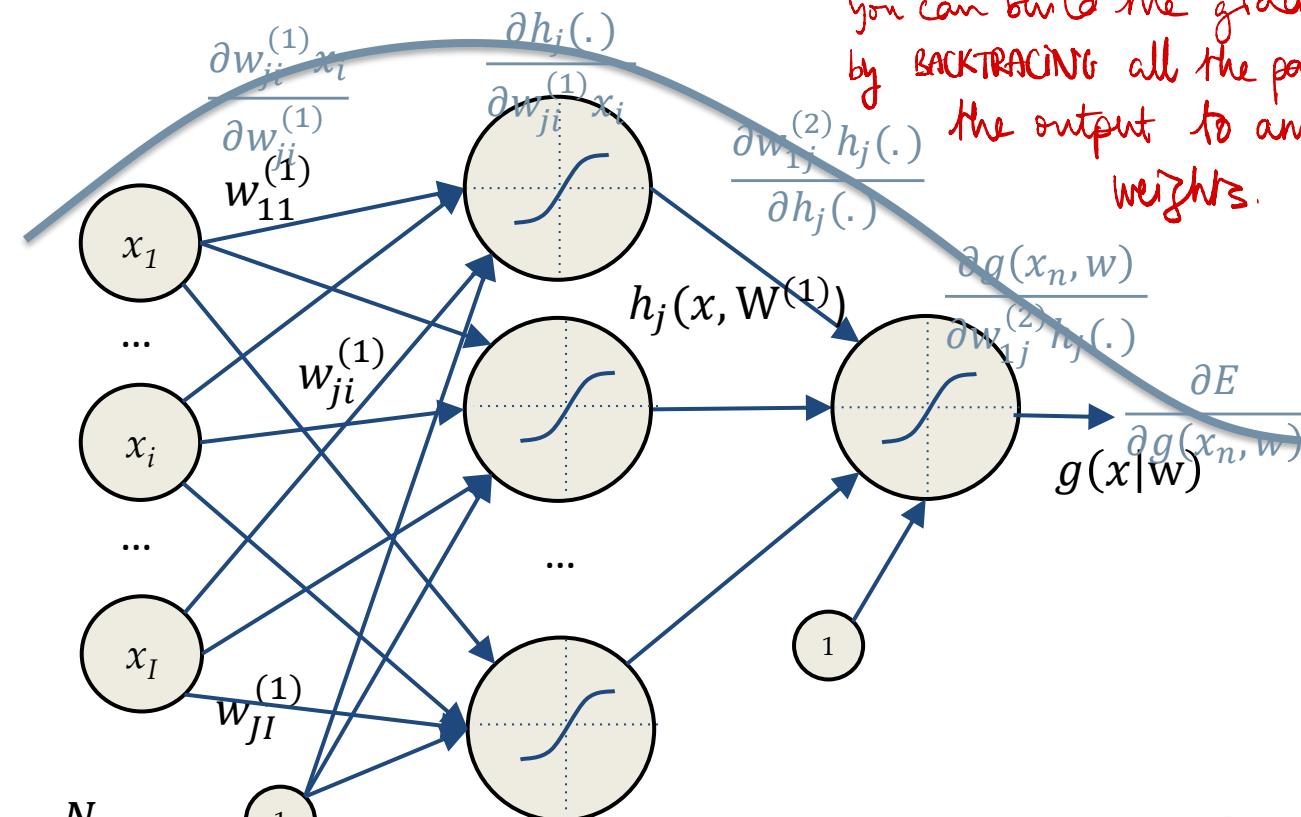


Backpropagation and Chain Rule (2)



"forward-backward algorithm"

you can build the gradient by BACKTRACING all the paths from the output to any single weights.



Forward pass :
you compute all the weighted sums \oplus
all the blocks : they are LOCAL, they don't depend on other neurons .

$$\frac{\partial E(w_{ji}^{(1)})}{\partial w_{ji}^{(1)}} = -2 \sum_n^N (t_n - g_1(x_n, w)) \cdot g'_1(x_n, w) \cdot w_{1j}^{(2)} \cdot h'_j \left(\sum_{i=0}^I w_{ji}^{(1)} \cdot x_{i,n} \right) \cdot x_i$$

$\frac{\partial E}{\partial w_{ji}^{(1)}}$ $\frac{\partial E}{\partial g(x_n, w)}$ $\frac{\partial g(x_n, w)}{\partial w_{1j}^{(2)} h_j(.)}$ $\frac{\partial w_{1j}^{(2)} h_j(.)}{\partial h_j(.)}$ $\frac{\partial h_j(.)}{\partial w_{ji}^{(1)} x_i}$ $\frac{\partial w_{ji}^{(1)} x_i}{\partial w_{ji}^{(1)}}$



Backpropagation and Chain Rule (2)



"forward-backward algorithm"

you can build the gradient by BACKTRACING all the paths from the output to any single weights.

you can evaluate every weights

Forward pass
(you compute every values)

It is very efficient on
GPUs. (perfect for // computing)

(you multiply everything)
Backward pass

you multiply
everything together
to get $\frac{\partial E(w_{ji}^{(1)})}{\partial w_{ji}^{(1)}}$

$$\frac{\partial E(w_{ji}^{(1)})}{\partial w_{ji}^{(1)}} = -2 \sum_n^N (t_n - g_1(x_n, w)) \cdot g'_1(x_n, w) \cdot w_{1j}^{(2)} \cdot h_j' \left(\sum_{i=0}^I w_{ji}^{(1)} \cdot x_{i,n} \right) \cdot x_i$$

$$\frac{\partial E}{\partial w_{ji}^{(1)}}$$

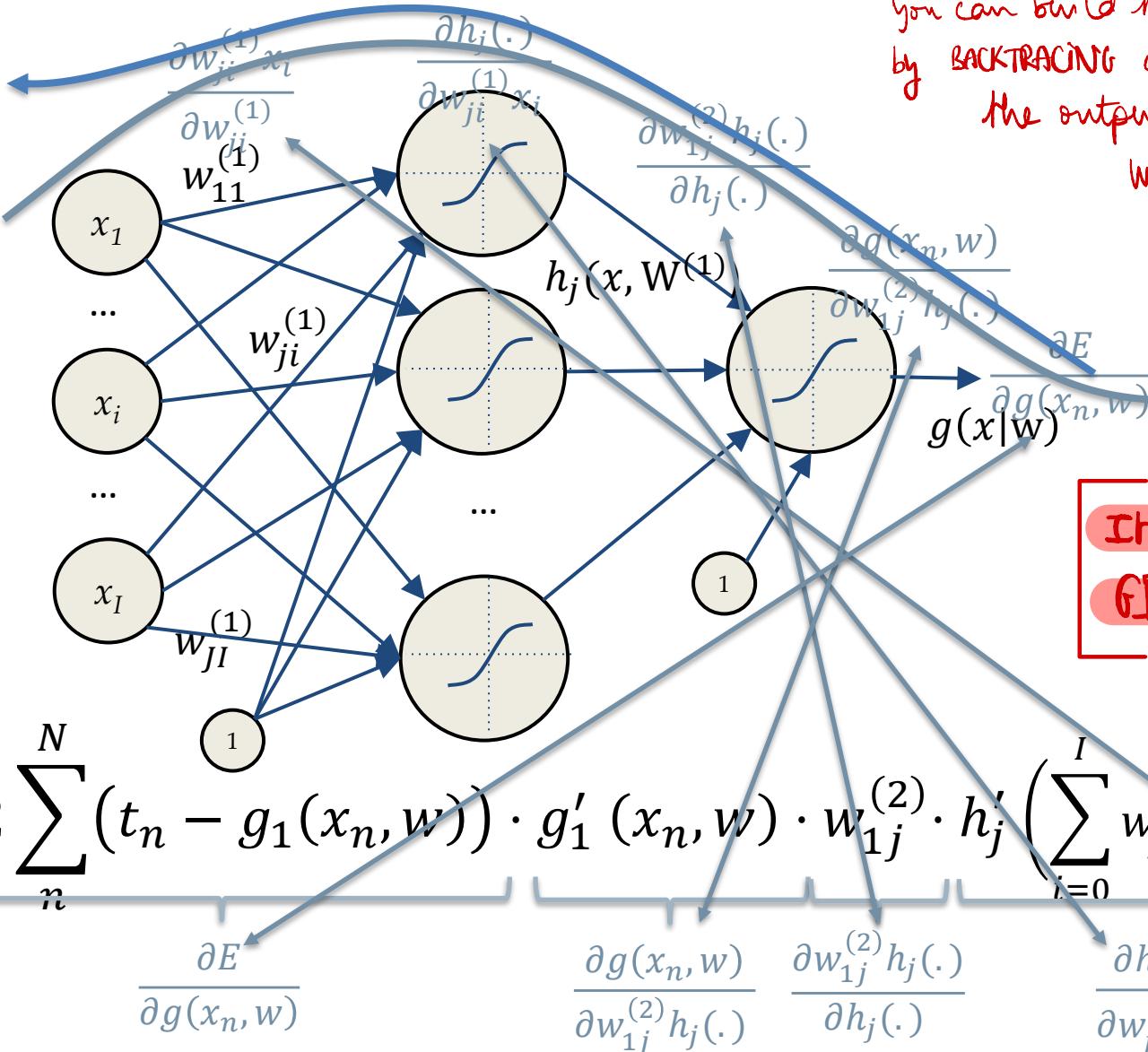
$$\frac{\partial E}{\partial g(x_n, w)}$$

$$\frac{\partial g(x_n, w)}{\partial w_{1j}^{(2)} h_j(.)}$$

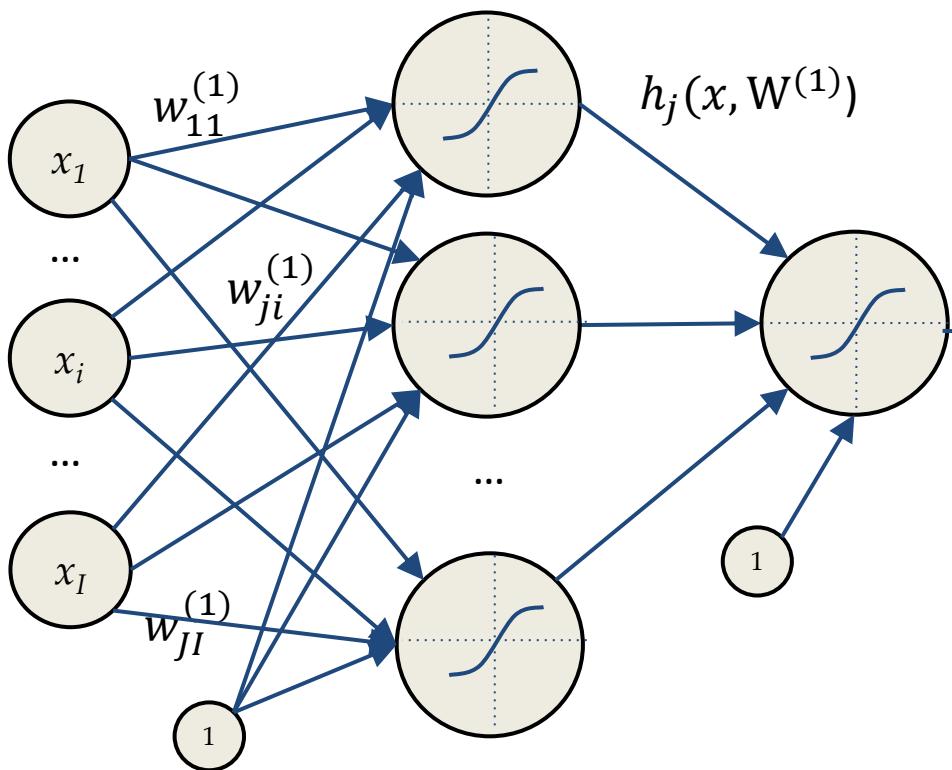
$$\frac{\partial w_{1j}^{(2)} h_j(.)}{\partial h_j(.)}$$

$$\frac{\partial h_j(.)}{\partial w_{ji}^{(1)} x_i}$$

$$\frac{\partial w_{ji}^{(1)} x_i}{\partial w_{ji}^{(1)}}$$



Gradient Descent Example



$$g_1(x_n|w) = g_1\left(\sum_{j=0}^J w_{1j}^{(2)} \cdot h_j\left(\sum_{i=0}^I w_{ji}^{(1)} \cdot x_{i,n}\right)\right)$$

$$E(w) = \sum_{n=1}^N (t_n - g_1(x_n, w))^2$$

Why is this a
good "Loss" (or "Error") function?

$$\frac{\partial E(w_{ji}^{(1)})}{\partial w_{ji}^{(1)}} = -2 \sum_n (t_n - g_1(x_n, w)) g'_1(x_n, w) w_{1j}^{(2)} h'_j\left(\sum_{i=0}^I w_{ji}^{(1)} \cdot x_{i,n}\right) x_{i,n}$$

Why should I use this?

Let's go back to the basics...

A Note on Maximum Likelihood Estimation

Let's observe i.i.d. samples from a Gaussian distribution with known σ^2

$$x_1, x_2, \dots, x_N \sim N(\mu, \sigma^2)$$

$$p(x|\mu, \sigma^2) = \frac{1}{\sqrt{2\pi}\sigma} e^{-\frac{(x-\mu)^2}{2\sigma^2}}$$

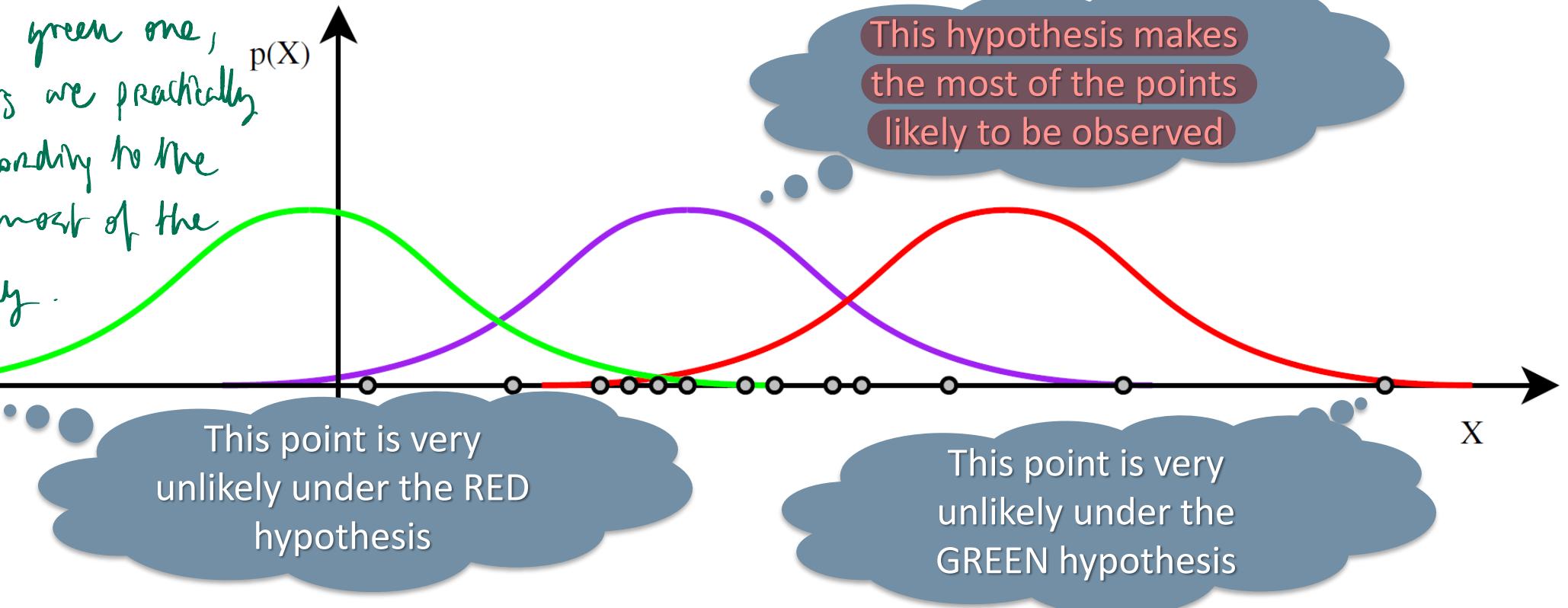
According to the green one, a lot of points are practically impossible. According to the purple one, most of the points are likely.

↓
MAXIMUM LIKELIHOOD.

This point is very unlikely under the RED hypothesis

This hypothesis makes the most of the points likely to be observed

This point is very unlikely under the GREEN hypothesis



A Note on Maximum Likelihood Estimation

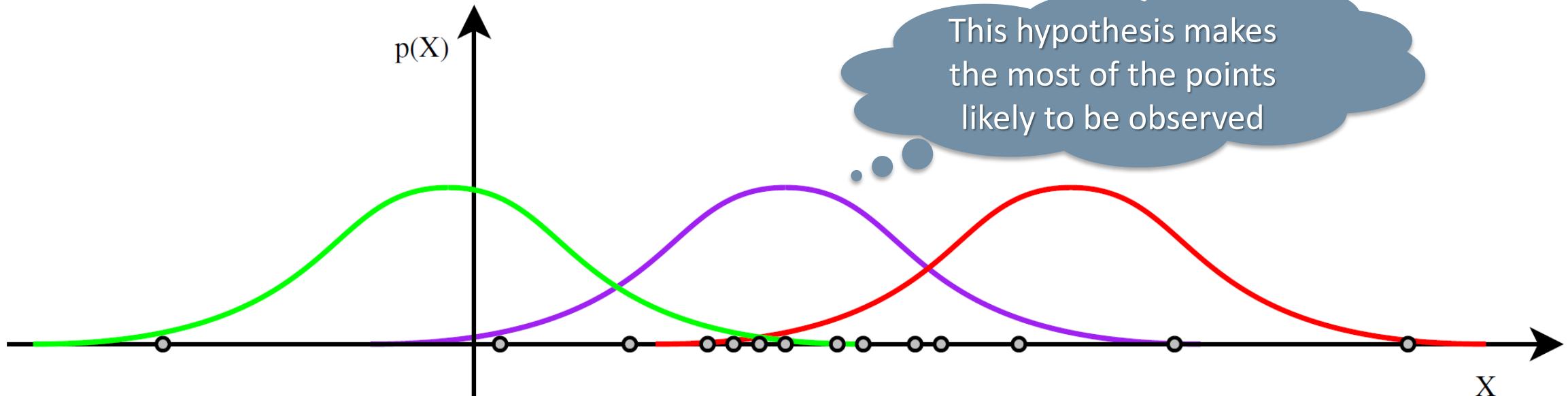


Let's observe i.i.d. samples from a Gaussian distribution with known σ^2

$$x_1, x_2, \dots, x_N \sim N(\mu, \sigma^2)$$

$$p(x|\mu, \sigma^2) = \frac{1}{\sqrt{2\pi}\sigma} e^{-\frac{(x-\mu)^2}{2\sigma^2}}$$

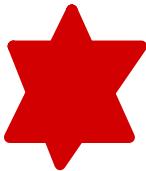
This hypothesis makes
the most of the points
likely to be observed



Maximum Likelihood: Choose parameters which maximize data probability



Maximum Likelihood Estimation: The Recipe



Let $\theta = (\theta_1, \theta_2, \dots, \theta_p)^T$ a vector of parameters, find the MLE for θ :

- Write the likelihood $L = P(Data|\theta)$ for the data
- [Take the logarithm of likelihood $l = \log P(Data|\theta)$] ...
- Work out $\frac{\partial L}{\partial \theta}$ or $\frac{\partial l}{\partial \theta}$ using high-school calculus
- Solve the set of simultaneous equations $\frac{\partial L}{\partial \theta_i} = 0$ or $\frac{\partial l}{\partial \theta_i} = 0$
- Check that θ^{MLE} is a maximum

Optional

To maximize/minimize the (log)likelihood you can use:

We know already about
gradient descent, let's try
with some analytical stuff ...

- Analytical Techniques (i.e., solve the equations)
- Optimization Techniques (e.g., Lagrange multipliers)
- Numerical Techniques (e.g., gradient descend)

The 3 main
techniques!

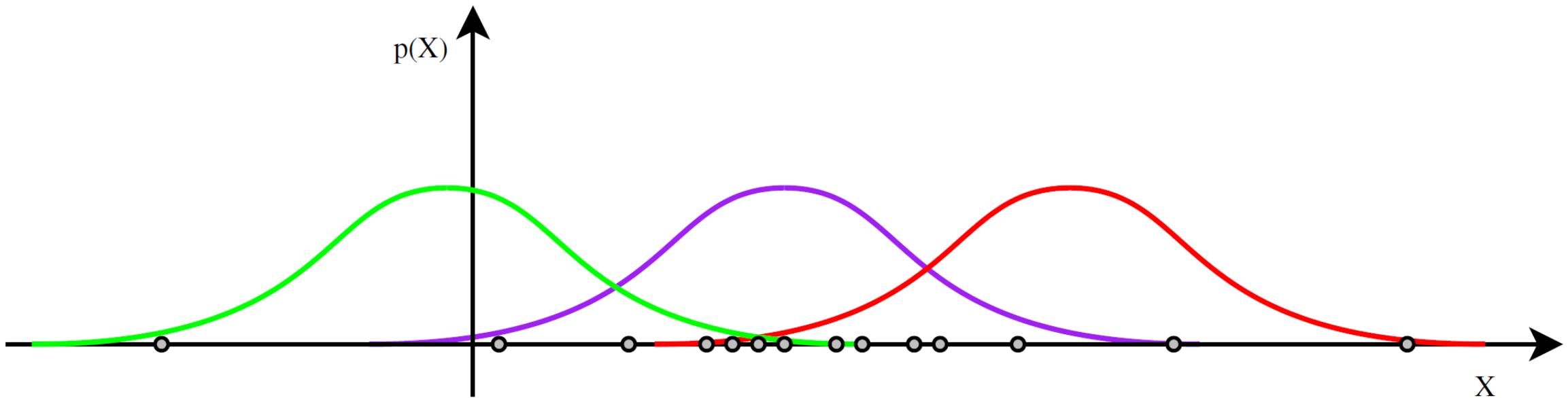


Maximum Likelihood Estimation Example

Let's observe i.i.d. samples coming from a Gaussian with known σ^2

$$x_1, x_2, \dots, x_N \sim N(\mu, \sigma^2)$$

$$p(x|\mu, \sigma^2) = \frac{1}{\sqrt{2\pi}\sigma} e^{-\frac{(x-\mu)^2}{2\sigma^2}}$$



Find the Maximum Likelihood Estimator for μ

Maximum Likelihood Estimation Example

Let's observe i.i.d. samples coming from a Gaussian with known σ^2

$$x_1, x_2, \dots, x_N \sim N(\mu, \sigma^2)$$

$$p(x|\mu, \sigma^2) = \frac{1}{\sqrt{2\pi}\sigma} e^{-\frac{(x-\mu)^2}{2\sigma^2}}$$

1. Write the likelihood $L = P(Data|\theta)$ for the data

$$\begin{aligned} L(\mu) &= p(x_1, x_2, \dots, x_N | \mu, \sigma^2) = \prod_{n=1}^N p(x_n | \mu, \sigma^2) = \\ &= \prod_{n=1}^N \frac{1}{\sqrt{2\pi}\sigma} e^{-\frac{(x_n-\mu)^2}{2\sigma^2}} \end{aligned}$$

INDEPENDENT RV



Maximum Likelihood Estimation Example

Let's observe i.i.d. samples coming from a Gaussian with known σ^2

$$x_1, x_2, \dots, x_N \sim N(\mu, \sigma^2) \quad p(x|\mu, \sigma^2) = \frac{1}{\sqrt{2\pi}\sigma} e^{-\frac{(x-\mu)^2}{2\sigma^2}}$$

- 2 • Take the logarithm $l = \log P(Data|\theta)$ of the likelihood

$$\begin{aligned} l(\mu) &= \log \left(\prod_{n=1}^N \frac{1}{\sqrt{2\pi}\sigma} e^{-\frac{(x_n-\mu)^2}{2\sigma^2}} \right) = \sum_{n=1}^N \log \frac{1}{\sqrt{2\pi}\sigma} e^{-\frac{(x_n-\mu)^2}{2\sigma^2}} = \\ &= N \cdot \log \frac{1}{\sqrt{2\pi}\sigma} - \frac{1}{2\sigma^2} \sum_{n=1}^N (x_n - \mu)^2 \end{aligned}$$



Maximum Likelihood Estimation Example

Let's observe i.i.d. samples coming from a Gaussian with known σ^2

$$x_1, x_2, \dots, x_N \sim N(\mu, \sigma^2) \quad p(x|\mu, \sigma^2) = \frac{1}{\sqrt{2\pi}\sigma} e^{-\frac{(x-\mu)^2}{2\sigma^2}}$$

3. Work out $\partial l/\partial \theta$ using high-school calculus

$$\begin{aligned} \frac{\partial l(\mu)}{\partial \mu} &= \frac{\partial}{\partial \mu} \left(N \cdot \log \frac{1}{\sqrt{2\pi}\sigma} - \frac{1}{2\sigma^2} \sum_n^N (x_n - \mu)^2 \right) = \\ &= -\frac{1}{2\sigma^2} \frac{\partial}{\partial \mu} \sum_n^N (x_n - \mu)^2 = \frac{1}{2\sigma^2} \sum_n^N 2(x_n - \mu) \end{aligned}$$



Maximum Likelihood Estimation Example

Let's observe i.i.d. samples coming from a Gaussian with known σ^2

$$x_1, x_2, \dots, x_N \sim N(\mu, \sigma^2)$$

$$p(x|\mu, \sigma^2) = \frac{1}{\sqrt{2\pi}\sigma} e^{-\frac{(x-\mu)^2}{2\sigma^2}}$$

- 4 • Solve the set of simultaneous equations $\frac{\partial l}{\partial \theta_i} = 0$

$$\frac{1}{2\sigma^2} \sum_n^N 2(x_n - \mu) = 0$$

$$\sum_n^N (x_n - \mu) = 0$$

$$\sum_n^N x_n = \sum_n^N \mu$$



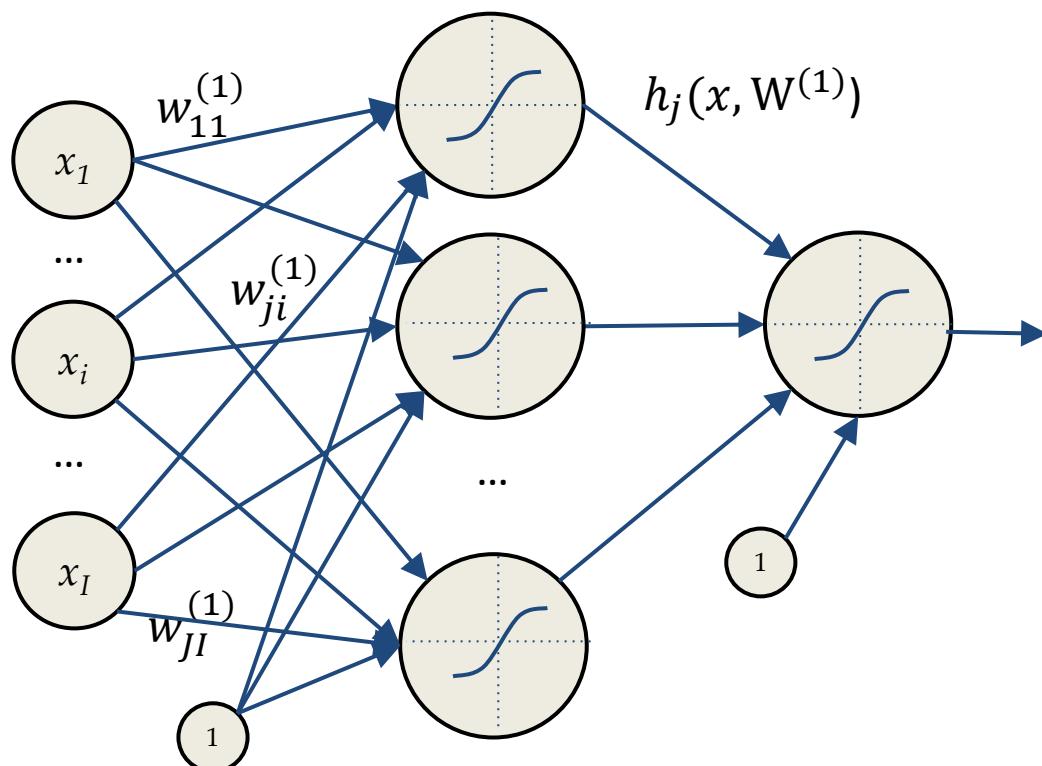
$$\mu^{MLE} = \frac{1}{N} \sum_n^N x_n$$

Let's apply this all to
Neural Networks!

VERY NATURAL !



Neural Networks for Regression



$$g(x_n|w) = g\left(\sum_{j=0}^J w_{1j}^{(2)} \cdot h_j\left(\sum_{i=0}^I w_{ji}^{(1)} \cdot x_{i,n}\right)\right)$$

NN is a universal approximator so let's assume
that \exists a set of weights for which we can say
that the target comes from the NN.

Goal: approximate a target function t having N observations

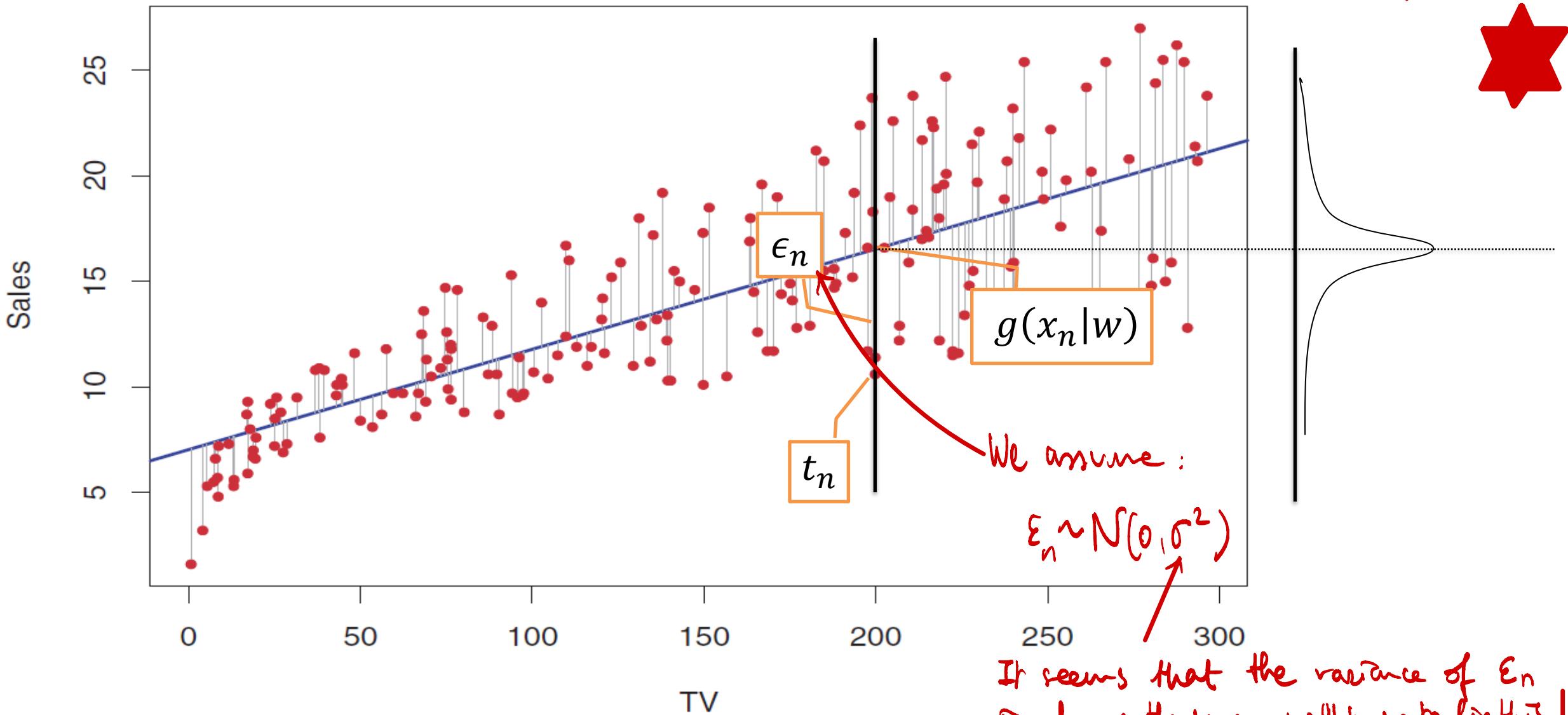
$$t_n = g(x_n|w) + \epsilon_n, \quad \epsilon_n \sim N(0, \sigma^2)$$

However, we
assume that these observations
come from the NN but have
been corrupted by some noise

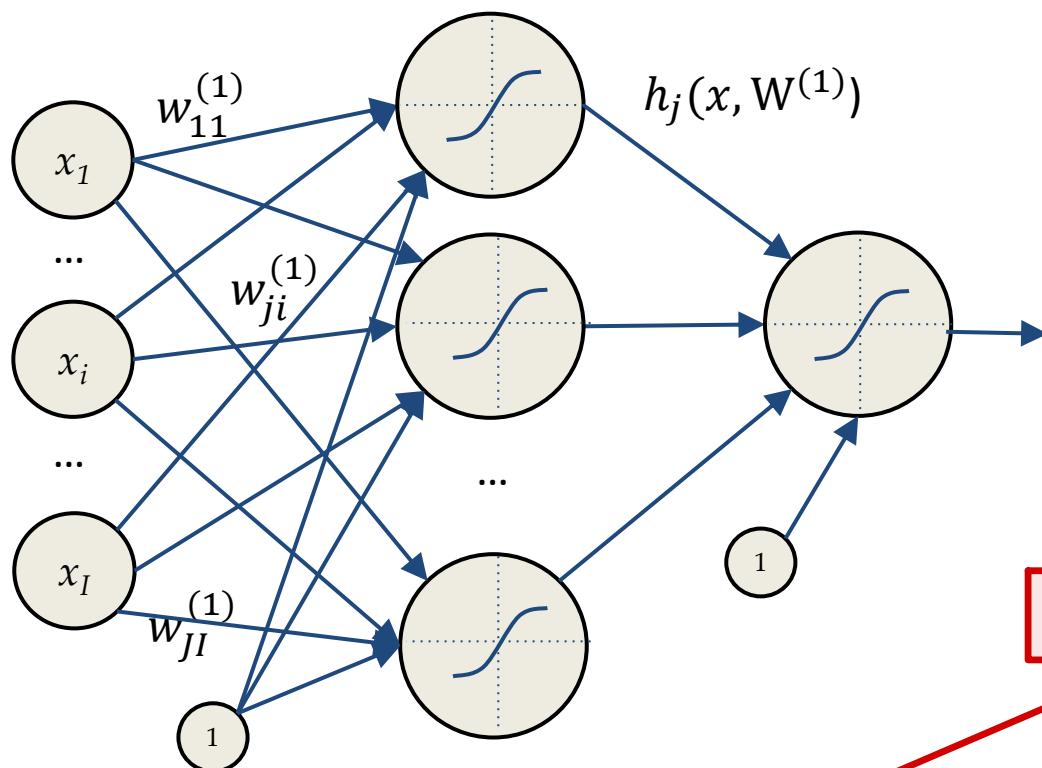


Statistical Learning Framework

Here : observations = \bullet : we have noise!
 target = --- : fct I want to approximate.



Neural Networks for Regression



$$g(x_n|w) = g\left(\sum_{j=0}^J w_{1j}^{(2)} \cdot h_j\left(\sum_{i=0}^I w_{ji}^{(1)} \cdot x_{i,n}\right)\right)$$

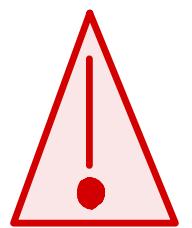
DETERMINISTIC

→ So:

Goal: approximate a target function t having N observations

$$t_n = g(x_n|w) + \epsilon_n, \quad \epsilon_n \sim N(0, \sigma^2) \quad \rightarrow \quad t_n \sim N(g(x_n|w), \sigma^2)$$

so let's find the weights of my NN
such that the likelihood of my
data is MAXIMIZED!



Maximum Likelihood Estimation for Regression



AS WE DID IN A PREVIOUS EXAMPLE (cf above) :

We have i.i.d. samples coming from a Gaussian with known σ^2

$$t_n \sim N(g(x_n|w), \sigma^2) \quad p(t|g(x|w), \sigma^2) = \frac{1}{\sqrt{2\pi}\sigma} e^{-\frac{(t-g(x|w))^2}{2\sigma^2}}$$

Write the likelihood $L = P(Data|\theta)$ for the data

$$\begin{aligned} L(w) &= p(t_1, t_2, \dots, t_N | g(x|w), \sigma^2) = \prod_{n=1}^N p(t_n | g(x_n|w), \sigma^2) = \\ &= \prod_{n=1}^N \frac{1}{\sqrt{2\pi}\sigma} e^{-\frac{(t_n-g(x_n|w))^2}{2\sigma^2}} \end{aligned}$$



Maximum Likelihood Estimation for Regression



We have i.i.d. samples coming from a Gaussian with known σ^2

$$t_n \sim N(g(x_n|w), \sigma^2) \quad p(t|g(x|w), \sigma^2) = \frac{1}{\sqrt{2\pi}\sigma} e^{-\frac{(t-g(x|w))^2}{2\sigma^2}}$$

Write the loglikelihood $l = \log P(Data|\theta)$ for the data

$$\begin{aligned} l(w) &= \log \prod_{n=1}^N \frac{1}{\sqrt{2\pi}\sigma} e^{-\frac{(t_n-g(x_n|w))^2}{2\sigma^2}} = \sum_n \log \left(\frac{1}{\sqrt{2\pi}\sigma} e^{-\frac{(t_n-g(x_n|w))^2}{2\sigma^2}} \right) \\ &= \sum_n \log \frac{1}{\sqrt{2\pi}\sigma} - \frac{1}{2\sigma^2} (t_n - g(x_n|w))^2 \end{aligned}$$



Maximum Likelihood Estimation for Regression



We have i.i.d. samples coming from a Gaussian with known σ^2

$$t_n \sim N(g(x_n|w), \sigma^2) \quad p(t|g(x|w), \sigma^2) = \frac{1}{\sqrt{2\pi}\sigma} e^{-\frac{(t-g(x|w))^2}{2\sigma^2}}$$

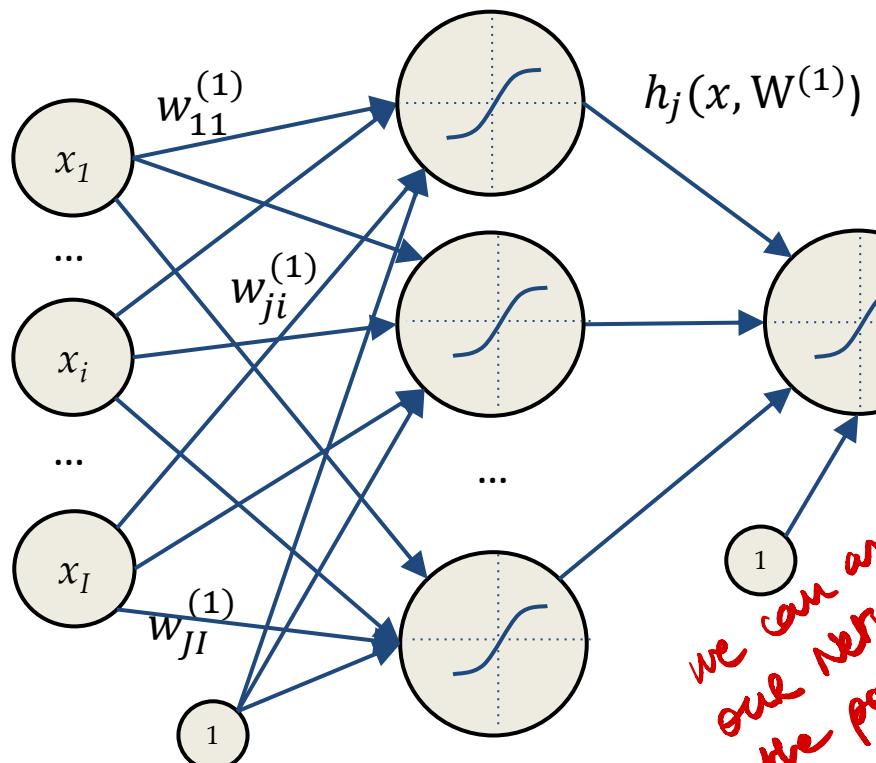
Look for the weights which maximixe the loglikelihood

$$\begin{aligned} \operatorname{argmax}_w l(w) &= \operatorname{argmax}_w \sum_n^N \log \frac{1}{\sqrt{2\pi}\sigma} - \frac{1}{2\sigma^2} (t_n - g(x_n|w))^2 = \\ &= \operatorname{argmin}_w \sum_n^N (t_n - g(x_n|w))^2 \end{aligned}$$

finding the weights of a NN under
assumption
of additive Gaussian noise w/ 0 mean and
cst variance \Leftrightarrow minimizing the
sum of squares .



Neural Networks for Classification



CLASSIFICATION this time!

$$g(x_n|w) = g\left(\sum_{j=0}^J w_{1j}^{(2)} \cdot h_j\left(\sum_{i=0}^I w_{ji}^{(1)} \cdot x_{i,n}\right)\right)$$

*we can assume
our network produces
the probabilities t_n binary data.*

**We are trying to learn the
"a posteriori probability" over the classes**

Goal: approximate a posterior probability t having N observations

$$g(x_n|w) = p(t_n|x_n), t_n \in \{0, 1\}$$

$$t_n \sim Be(g(x_n|w))$$

*Let's assume that my NN is producing the "a posteriori probability" over the classes
↳ the output = the a posteriori probability.*



Maximum Likelihood Estimation for Classification

We estimate the weights once again with MLE

In classification, noise = 0 instead of 1 (& vice versa).

We have some i.i.d. samples coming from a Bernulli distribution

$$t_n \sim Be(g(x_n|w))$$

prob of $t=1$: $g(x|w)$
 $t=0$: $1-g(x|w)$

$$p(t|g(x|w)) = g(x|w)^t \cdot (1 - g(x|w))^{1-t}$$

we can write a $B(\cdot)$ like this (try with $t=0, t=1$). t = output

Write the likelihood $L = P(Data|\theta)$ for the data

$$\begin{aligned} L(w) &= p(t_1, t_2, \dots, t_N | g(x|w)) = \prod_{n=1}^N p(t_n | g(x_n|w)) = \\ &= \prod_{n=1}^N g(x_n|w)^{t_n} \cdot (1 - g(x_n|w))^{1-t_n} \end{aligned}$$



Maximum Likelihood Estimation for Classification



We have some i.i.d. samples coming from a Bernulli distribution

$$t_n \sim Be(g(x_n|w)) \quad p(t|g(x|w)) = g(x|w)^t \cdot (1 - g(x|w))^{1-t}$$

Compute the log likelihood $l = \log P(Data|\theta)$ for the data

$$\begin{aligned} l(w) &= \log \prod_{n=1}^N g(x_n|w)^{t_n} \cdot (1 - g(x_n|w))^{1-t_n} \\ &= \sum_{n=1}^N t_n \log g(x_n|w) + (1 - t_n) \log(1 - g(x_n|w)) \end{aligned}$$



Maximum Likelihood Estimation for Classification



We have some i.i.d. samples coming from a Bernulli distribution

$$t_n \sim Be(g(x_n|w)) \quad p(t|g(x|w)) = g(x|w)^t \cdot (1 - g(x|w))^{1-t}$$

Look for the weights which maximize the loglikelihood

$$\operatorname{argmax}_w l(w) = \operatorname{argmax}_w \sum_n^N t_n \log g(x_n|w) + (1 - t_n) \log(1 - g(x_n|w))$$

This is the MLE for the parameters of the NN when the output can be described as a Be distribution with a parameter which is the a-posteriori probability predicted by the network.

can be reduced
Crossentropy
 $-\sum_n^N t_n \log g(x_n|w)$

↑
Loss function that we SHOULD use in classification! ≠ from \sum squares error!

What about perceptron



How to Choose the Error Function?



We have observed different error functions so far

$$\text{Regression: } E(w) = \sum_{n=1}^N (t_n - g_1(x_n, w))^2$$

Sum of Squared Errors

$$\text{binary classification: } E(w) = - \sum_n t_n \log g(x_n|w) + (1 - t_n) \log(1 - g(x_n|w))$$

Binary Crossentropy

Error functions define the task to be solved, but how to design them?

- Use all your knowledge/assumptions about the data distribution
- Exploit background knowledge on the task and the model
- Use your creativity!

in our regression example, the data violate
our assumption of constant std.

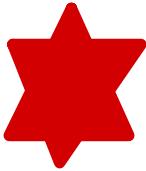
This requires lots of trial and errors ...

→ ex: SVM: linearly separate data

As for the Perceptron ...



Hyperplanes Linear Algebra



Let consider the hyperplane (affine set) $L \in \mathbb{R}^2$

$$L: w_0 + w^T x = 0$$

Any two points x_1 and x_2 on $L \in \mathbb{R}^2$ have

$$w^T(x_1 - x_2) = 0$$

The versor normal to $L \in \mathbb{R}^2$ is then

$$\overset{\text{normal}}{\uparrow} \quad w^* = w / \|w\|$$

$$w^* = w / \|w\|$$

It gives the direction $\perp L$.

For any point x_0 in $L \in \mathbb{R}^2$ we have

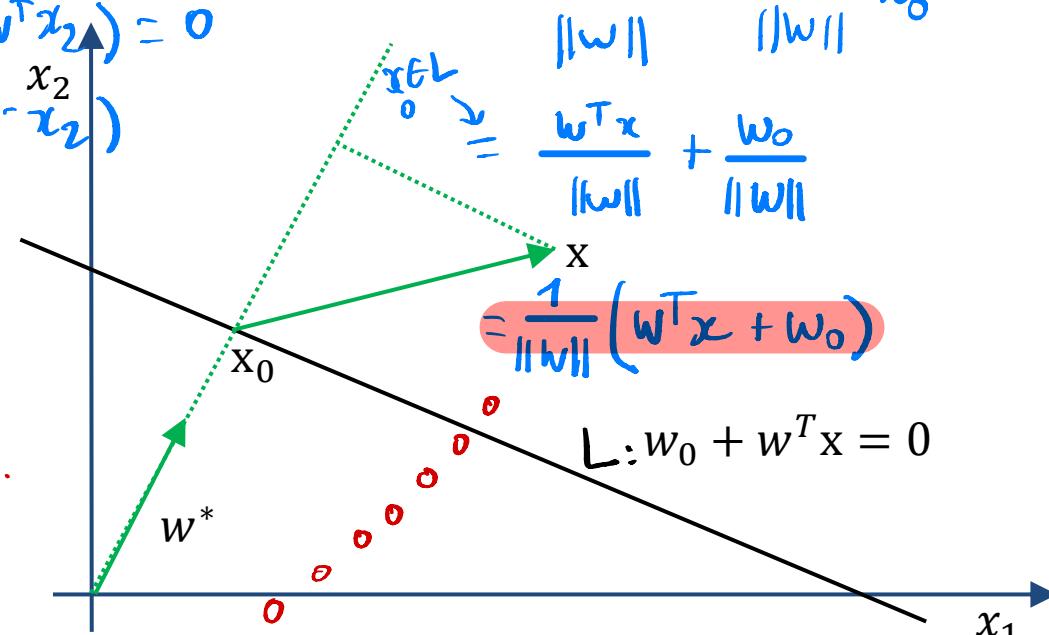
$$w^T x_0 = -w_0 \quad (\Leftrightarrow w_0 + w^T x_0 = 0)$$

The signed distance of any point x from $L \in \mathbb{R}^2$ is

$$w^{*T}(x - x_0) = \frac{1}{\|w\|} (w^T x + w_0)$$

any point on the line .

$$\begin{aligned} w^{*T}(x - x_0) &= \frac{w^T}{\|w\|} (x - x_0) \\ &= \frac{w^T x}{\|w\|} - \frac{w^T x_0}{\|w\|} \\ &\stackrel{x \in L}{=} \frac{w^T x}{\|w\|} + \frac{w_0}{\|w\|} \\ &= \frac{1}{\|w\|} (w^T x + w_0) \end{aligned}$$



$(w^T x + w_0)$ is proportional to the distance of x from the plane defined by $(w^T x + w_0) = 0$

if x above L : $w^T x + w_0 > 0$
below L : $w^T x + w_0 < 0$



Perceptron Learning Algorithm (1/2)



It can be shown, the error function the Hebbian rule is minimizing is the distance of misclassified points from the decision boundary.

Let's code the perceptron output as +1/-1

- If an output which should be +1 is misclassified then $w^T x + w_0 < 0$
- For an output with -1 we have the opposite

The goal becomes minimizing

$$D(w, w_0) = - \sum_{i \in M} t_i (w^T x_i + w_0)$$

(with \ominus : becomes positive)

$t_i (w^T x_i + w_0)$ Set of points misclassified

$w_0 \rightarrow \text{error}$

A diagram illustrating the error function. A horizontal line represents the decision boundary $w^T x + w_0 = 0$. Above the line, a set of points is labeled "Set of points misclassified". Below the line, a point is labeled $t_i (w^T x_i + w_0)$. A red arrow points from this label to the point, and another red arrow points from the label to the text "with \ominus : becomes positive". A red bracket underlines the term w_0 and is labeled "error".

This is non negative and proportional to the distance of the misclassified points from $w^T x + w_0 = 0$

so by minimizing
 $D \rightarrow$ we reach
the right
solution.

Perceptron Learning Algorithm (2/2)



Let's minimize by stochastic gradient descend the error function

$$D(w, w_0) = - \sum_{i \in M} t_i (w^T x_i + w_0)$$

The gradients with respect to the model parameters are

$$\frac{\partial D(w, w_0)}{\partial w} = - \sum_{i \in M} t_i \cdot x_i \quad \frac{\partial D(w, w_0)}{\partial w_0} = - \sum_{i \in M} t_i$$

Stochastic gradient descent applies for each misclassified point

$$\begin{pmatrix} w^{k+1} \\ w_0^{k+1} \end{pmatrix} = \begin{pmatrix} w^k \\ w_0^k \end{pmatrix} + \eta \begin{pmatrix} t_i \cdot x_i \\ t_i \end{pmatrix} = \begin{pmatrix} w^k \\ w_0^k \end{pmatrix} + \eta \begin{pmatrix} t_i \cdot x_i \\ t_i \cdot x_0 \end{pmatrix}$$

Hebbian learning
implements Stochastic
Gradient Descent

