



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

AIRLAB
ARTIFICIAL INTELLIGENCE AND ROBOTICS LAB

«Deep Learning is not AI, nor Machine Learning»

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

ARTIFICIAL INTELLIGENCE

Early artificial intelligence stirs excitement.



Neural Networks are as old as Artificial Intelligence

MACHINE LEARNING

Machine learning begins to flourish.



DEEP LEARNING

Deep learning breakthroughs drive AI boom.



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

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.

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:

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.



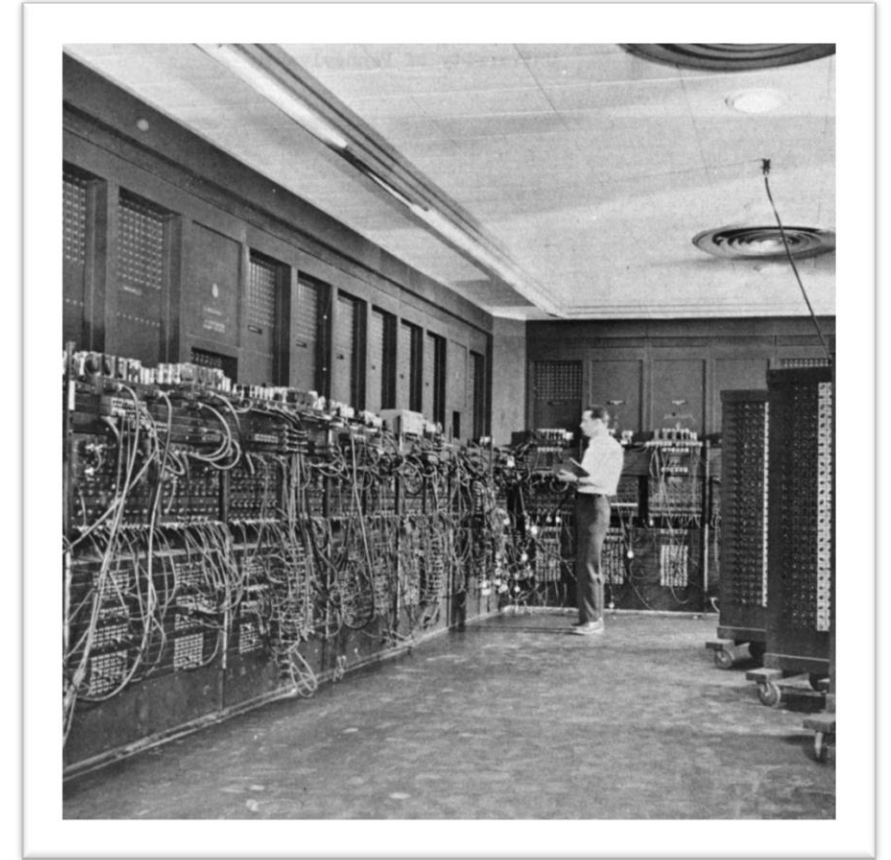
Let's go back to 1940s ...

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 Computational Model

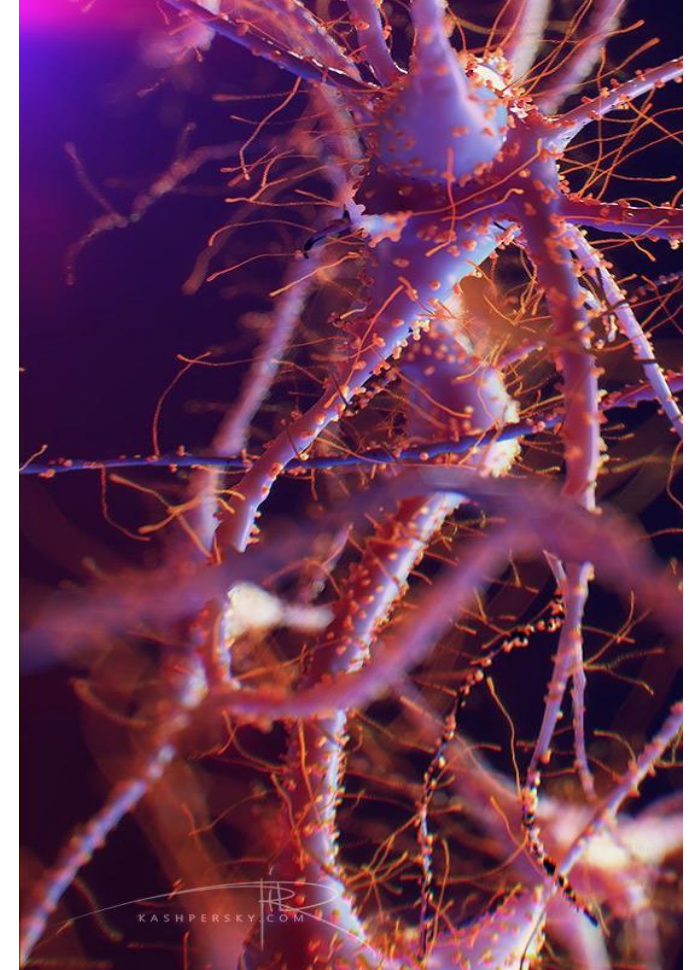
The human brain has a huge number of computing units:

- 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

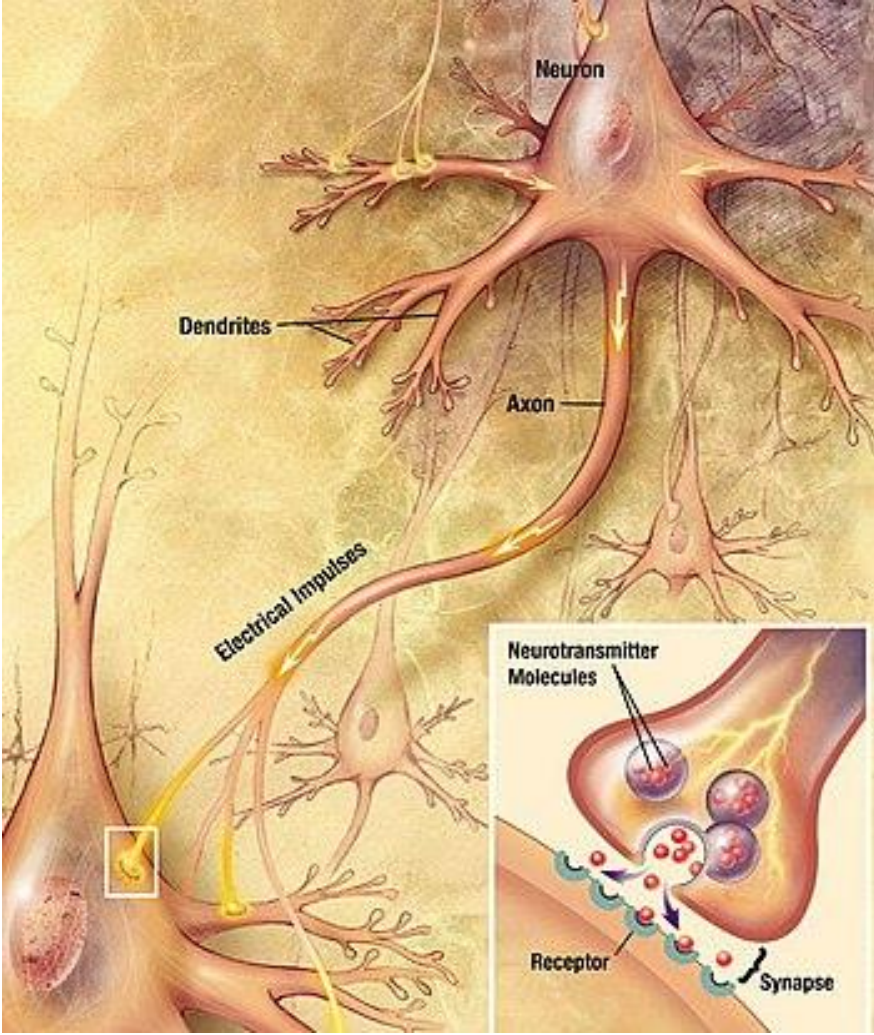
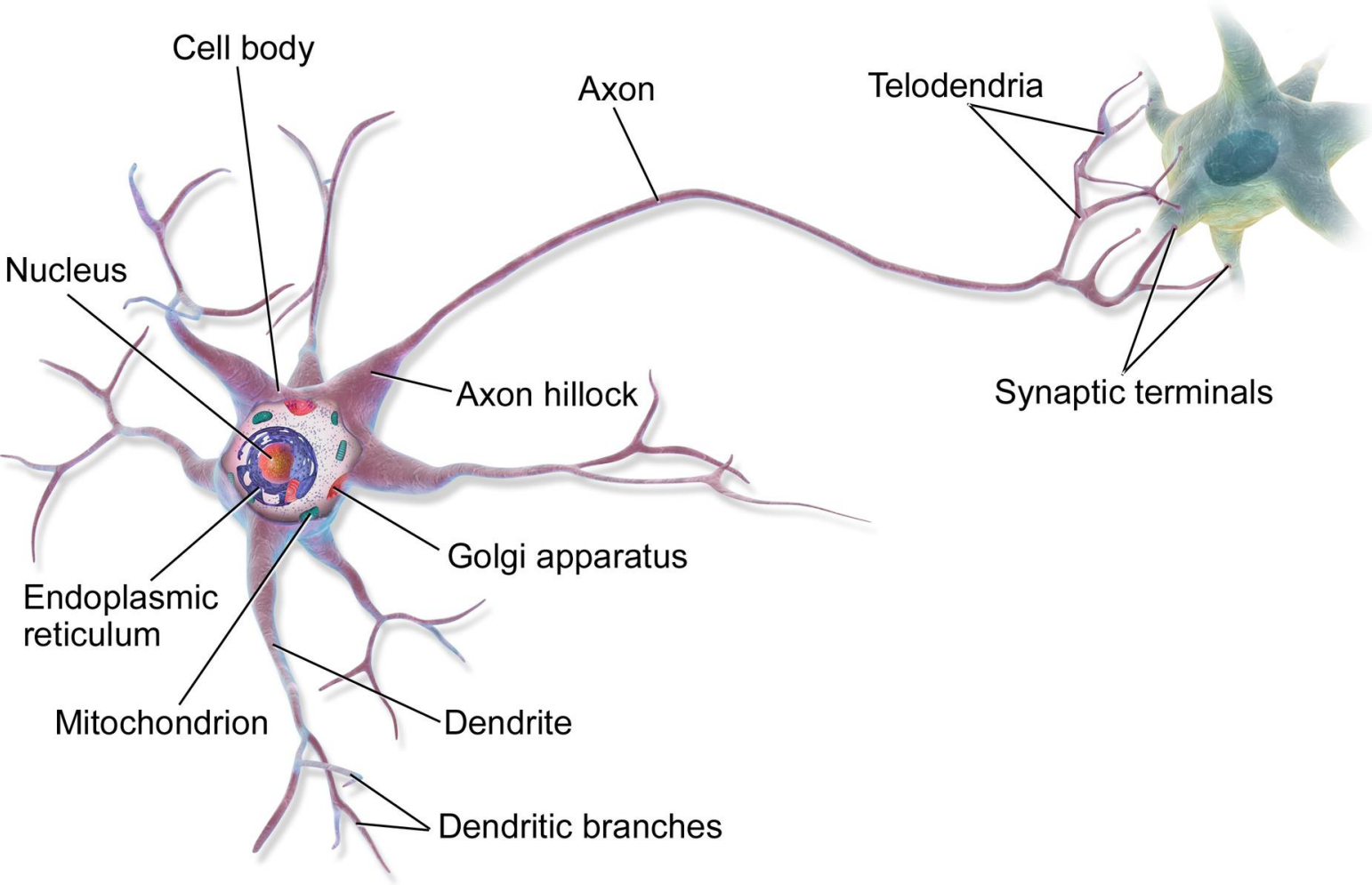
The computational model of the brain is:

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

Perceptron: a computational model based on the brain!



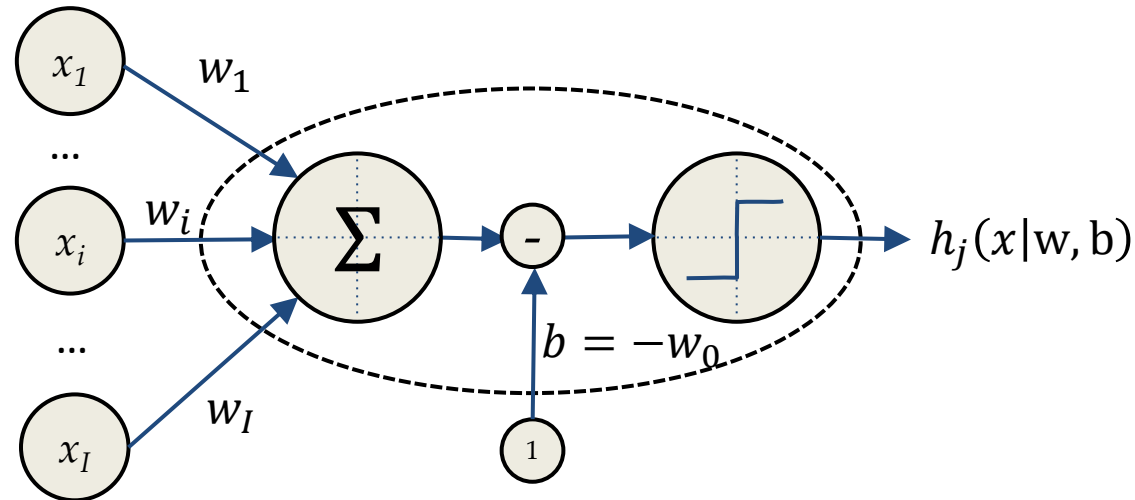
Computation in Biological Neurons



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

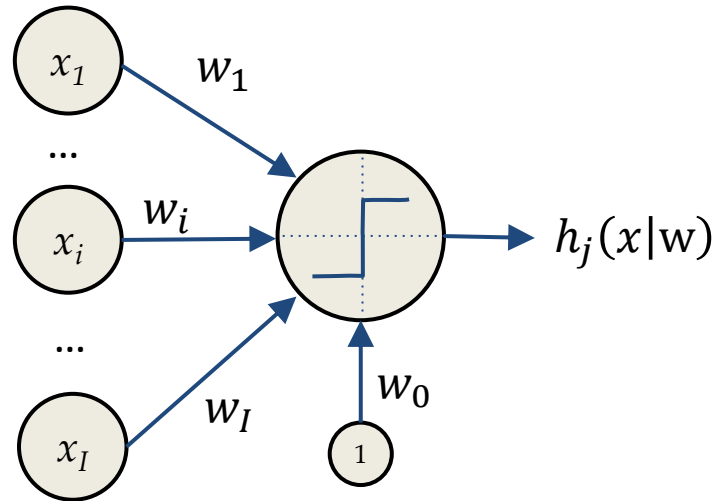


$$h_j(x|w, b) = h_j\left(\sum_{i=1}^I w_i \cdot x_i - b\right) = h_j\left(\sum_{i=0}^I w_i \cdot x_i\right) = h_j(w^T x)$$

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

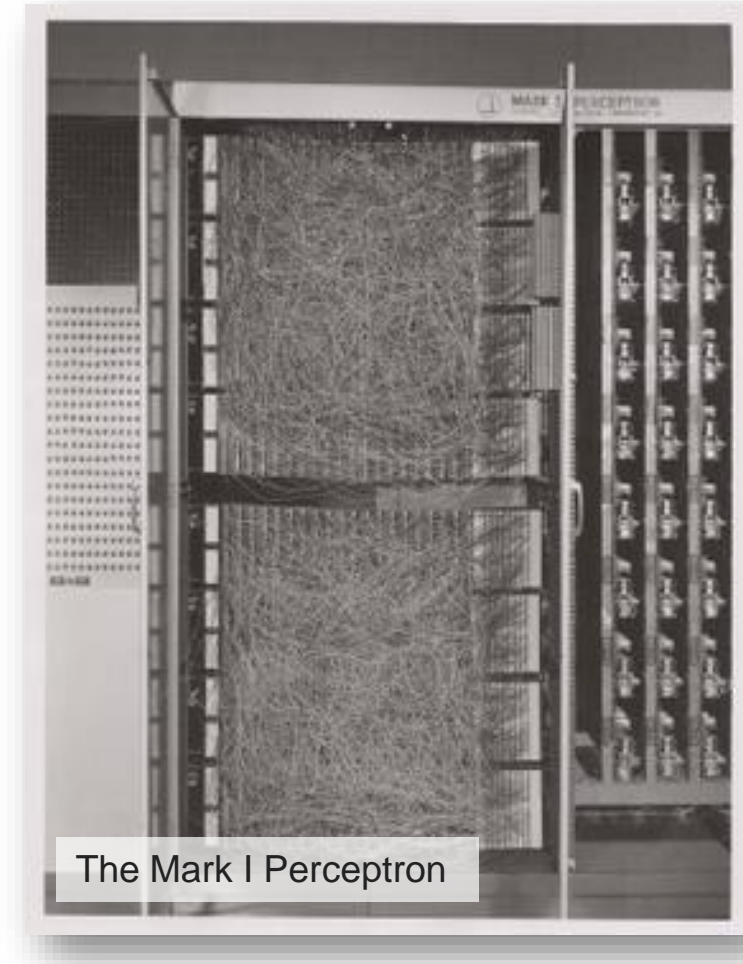


$$h_j(x|w, b) = h_j\left(\sum_{i=1}^I w_i \cdot x_i - b\right) = h_j\left(\sum_{i=0}^I w_i \cdot x_i\right) = h_j(w^T x)$$

Who did it first?

Several researchers were investigating models for the brain

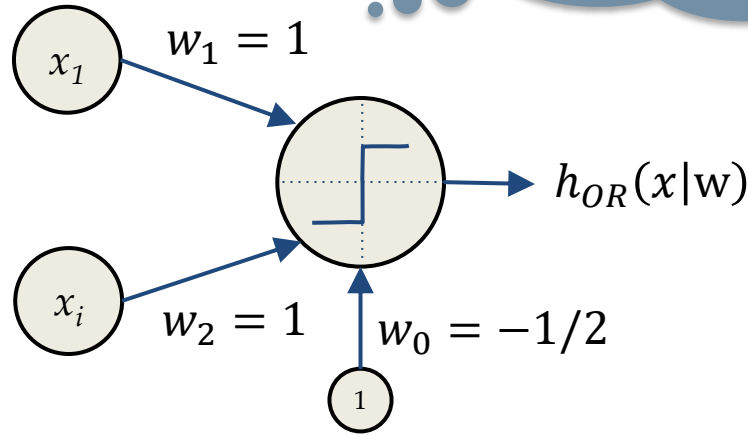
- In 1943, Warren McCulloch and Walter Harry Pitts proposed the Threshold Logic Unit or Linear Unit, the activation function was a threshold unit equivalent to the Heaviside step function
- In 1957, Frank Rosenblatt developed the first Perceptron. Weights were encoded in potentiometers, and weight updates 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)



The Mark I Perceptron

What can you do with it?

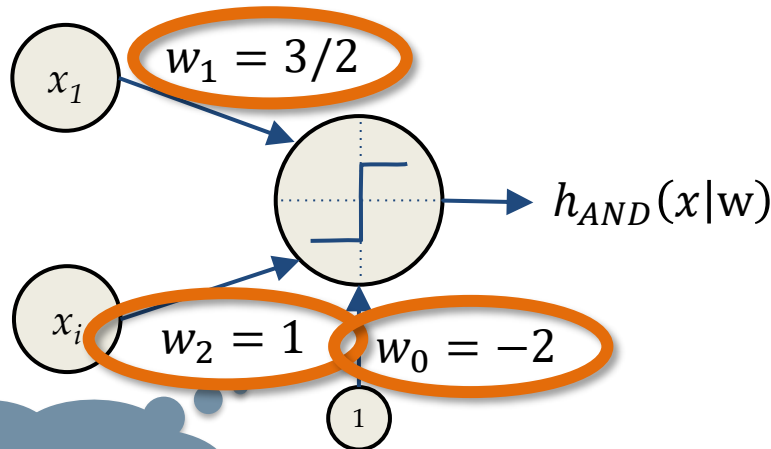
x_0	x_1	x_2	OR
1	0	0	0
1	0	1	1
1	1	0	1
1	1	1	1



Perceptron as Logical OR

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

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



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

Perceptron as Logical AND

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

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

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

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

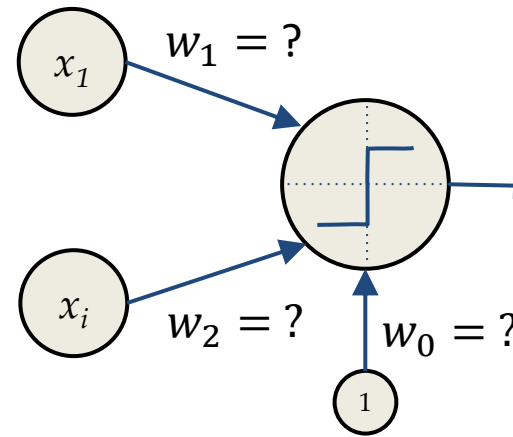
Fix the weights one sample at the time (online), and only if the sample is not correctly predicted

Perceptron Example



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



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

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

Does the procedure always converge?

Does it always converge to the same sets of weights?

Perceptron Math

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

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

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

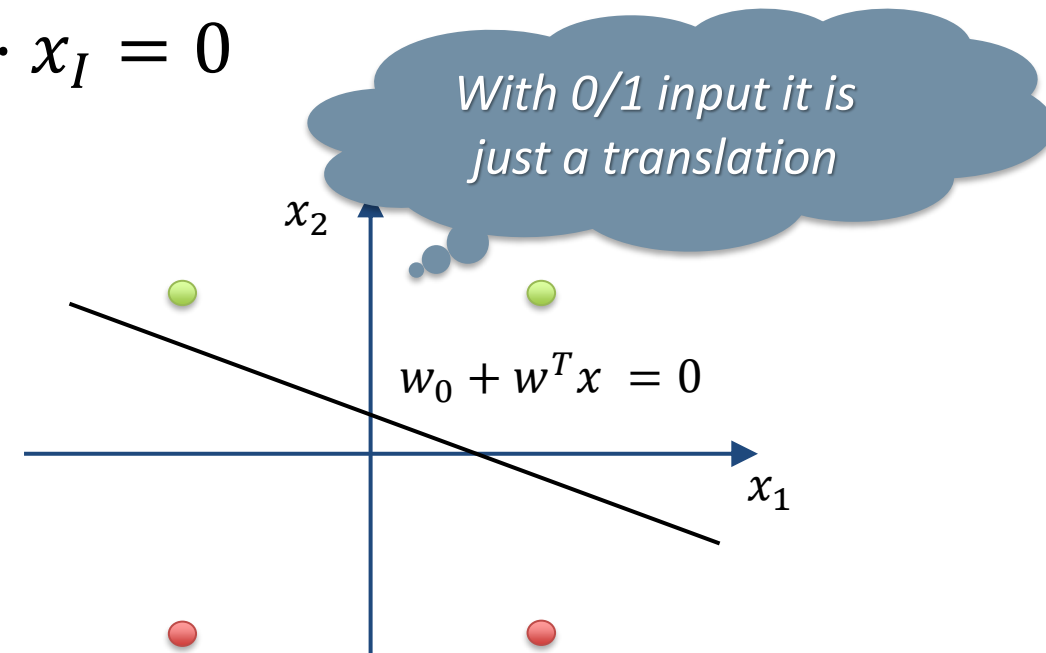
$$w_0 + w_1 \cdot x_1 + \dots + w_I \cdot x_I = 0$$

In 2D, this turns into

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

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

$$x_2 = -\frac{w_0}{w_2} - \frac{w_1}{w_2} \cdot x_1$$



Boolean Operators are Linear Boundaries

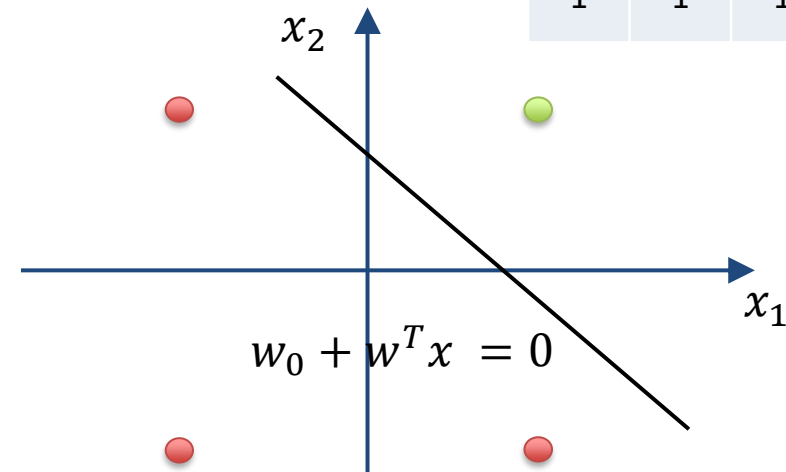
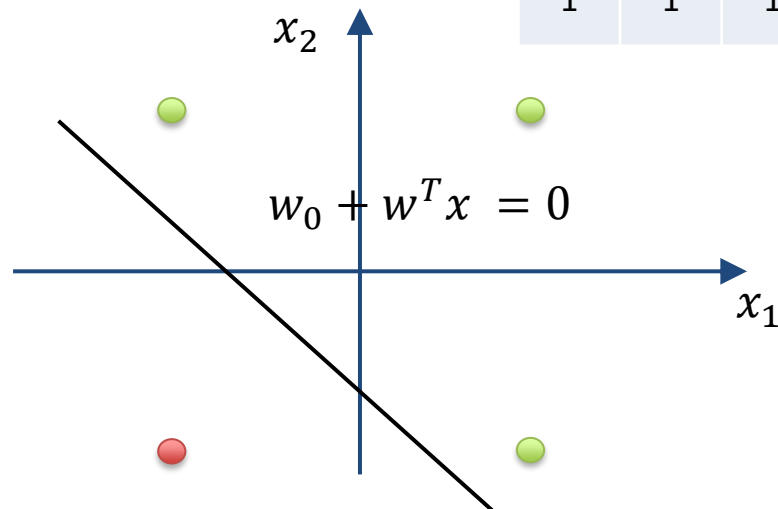
What's about it? We had already Boolean operators

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!

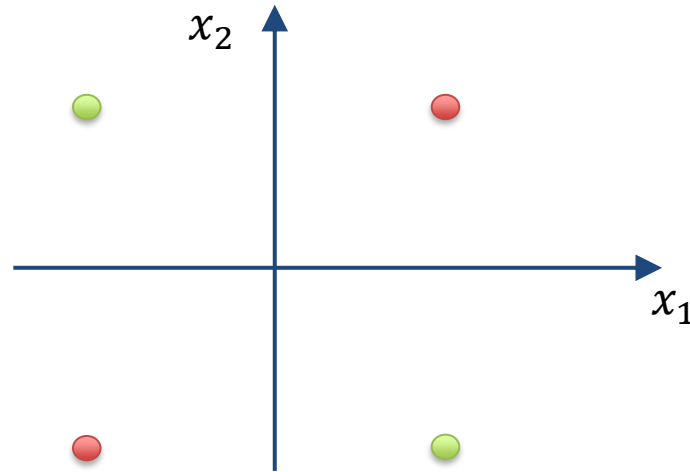
x_0	x_1	x_2	AND
1	-1	-1	-1
1	-1	1	-1
1	1	-1	-1
1	1	1	1



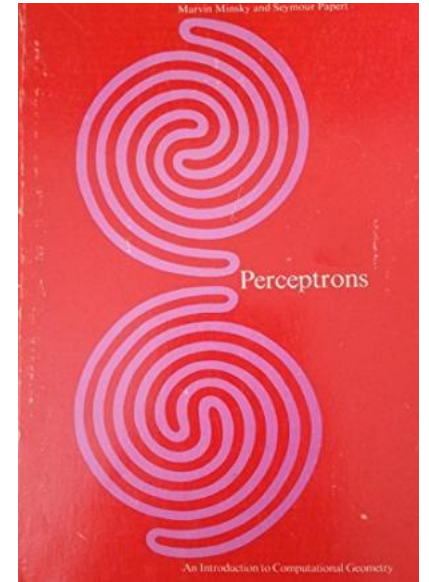
What can't you do with it?

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

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.


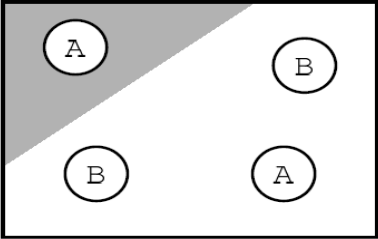
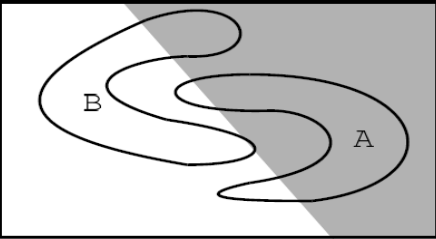
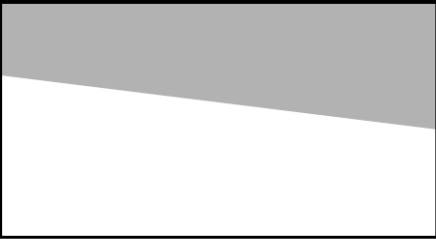
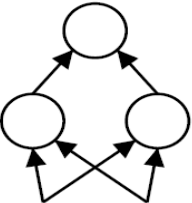
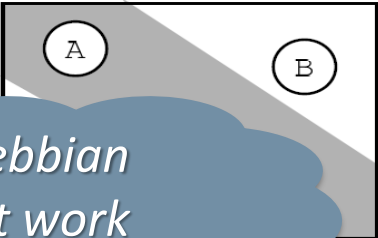
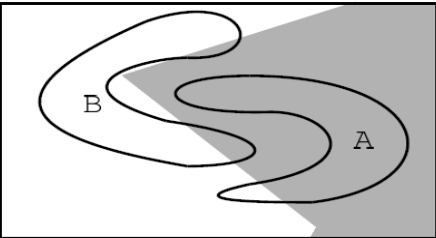
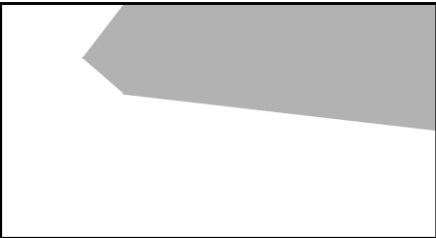
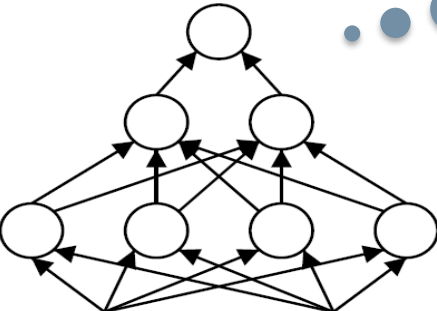
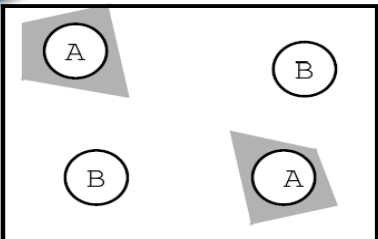
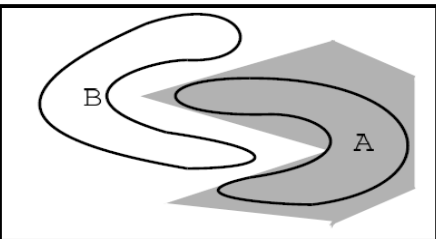
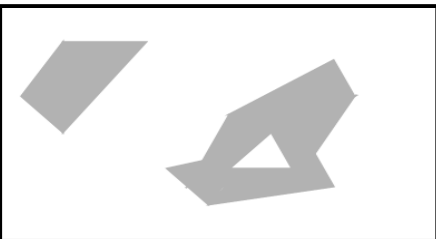


The Perceptron does not work any more and we need alternative solutions

- Non linear boundary
- Alternative input representations

*The idea behind Multi
Layer Perceptrons*

What can't you do with it?

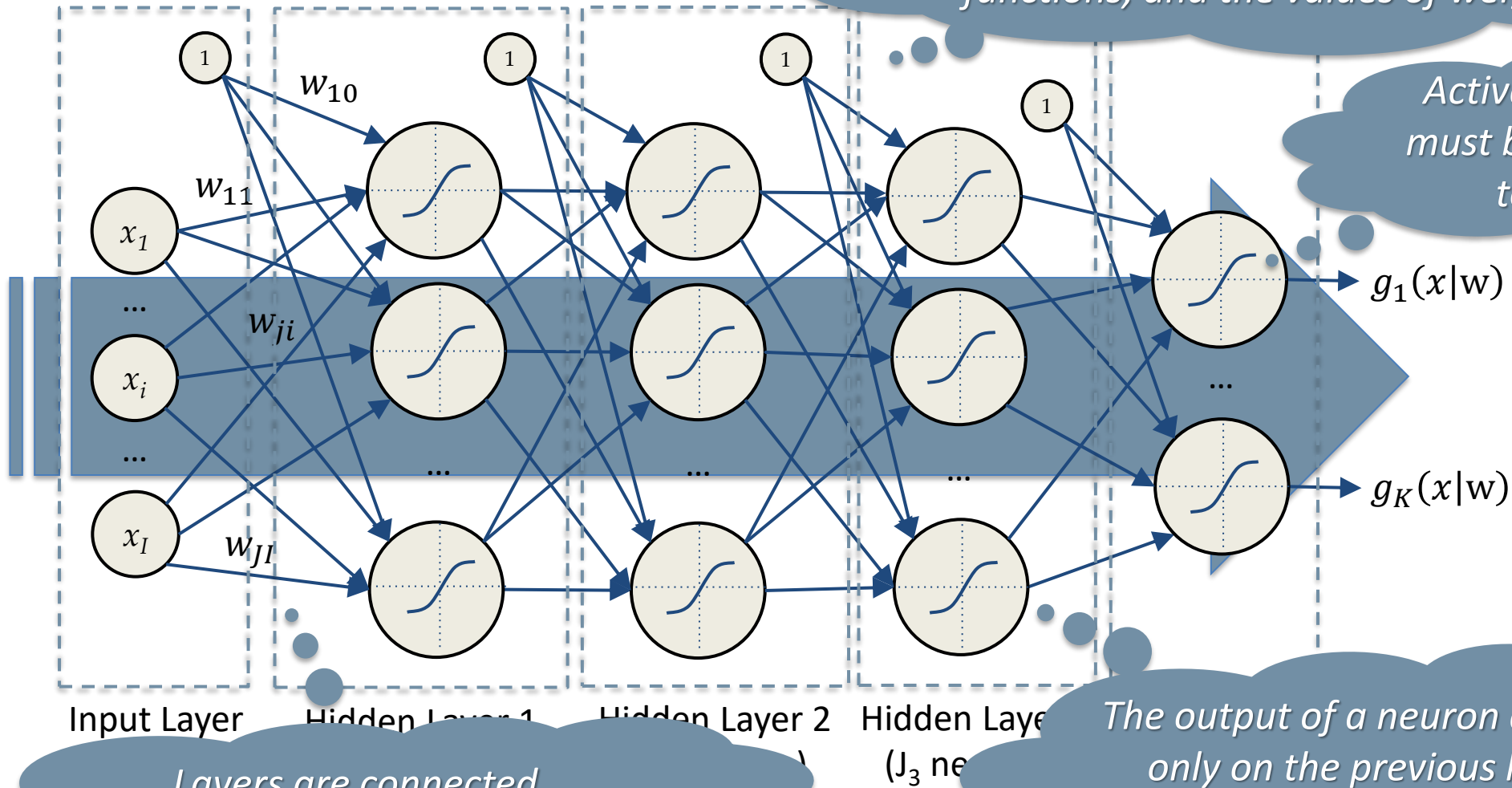
Topology	Type of Decision Region	XOR Problem	Classes with Meshed Regions	Most General Region Shapes
	<p>Half bounded by hyperplanes</p>			
	<p>Convex Open or</p>			
	<p>Arbitrary Regions (Complexity limited by the number of nodes)</p>			

Unfortunately Hebbian learning does not work any more ...

Layer Perceptrons

Feed Forward Neural Networks

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

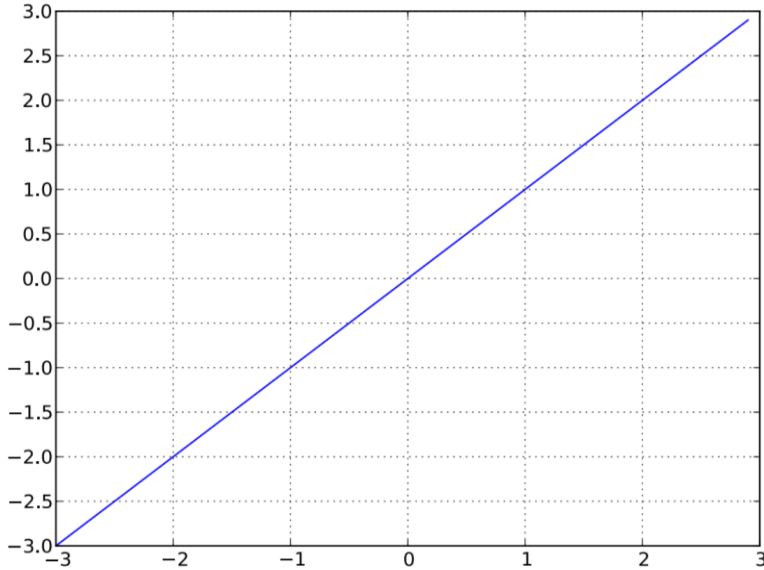


Activation functions must be differentiable to train it ...

Layers are connected through weights $W^{(l)} = \{w_{ji}^{(l)}\}$

The output of a neuron depends only on the previous layers
 $h^{(l)} = \{h_j^{(l)}(h^{(l-1)}, W^{(l)})\}$

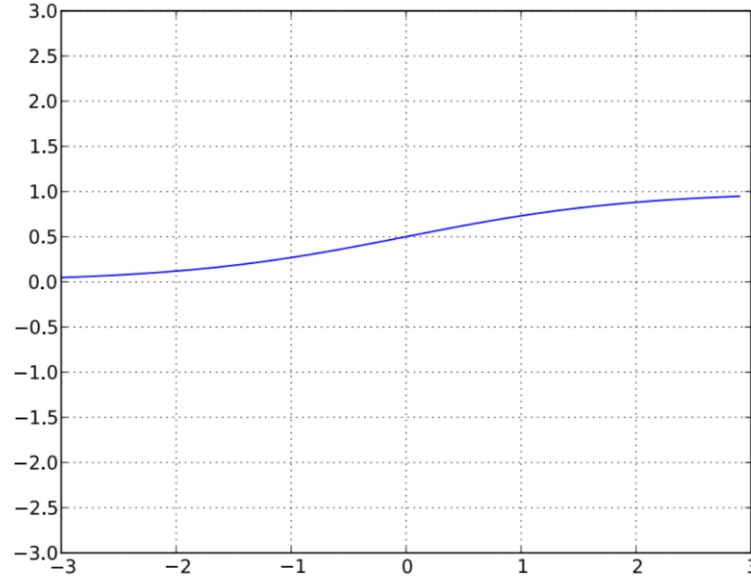
Which Activation Function?



Linear activation function

$$g(a) = a$$

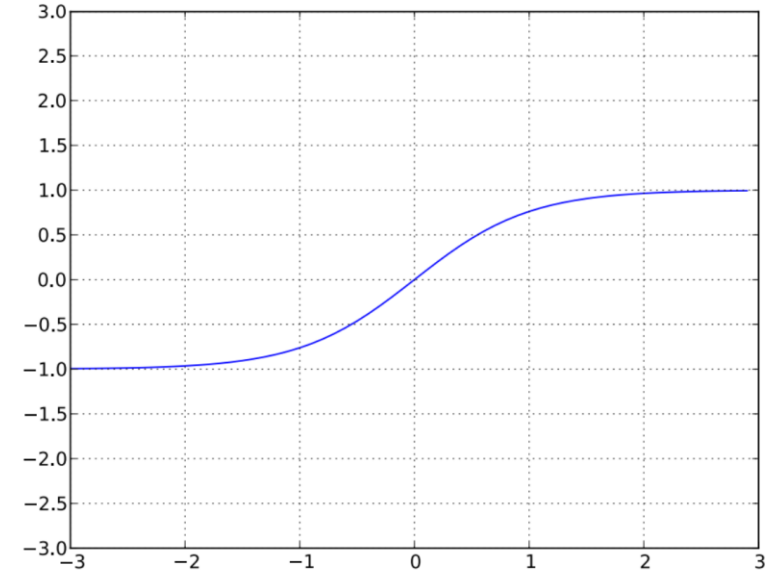
$$g'(a) = 1$$



Sigmoid activation function

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

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



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

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

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

- Use a Linear activation function for the output neuron

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

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

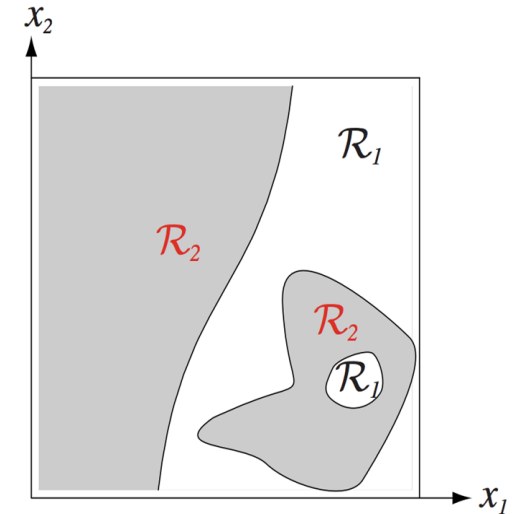
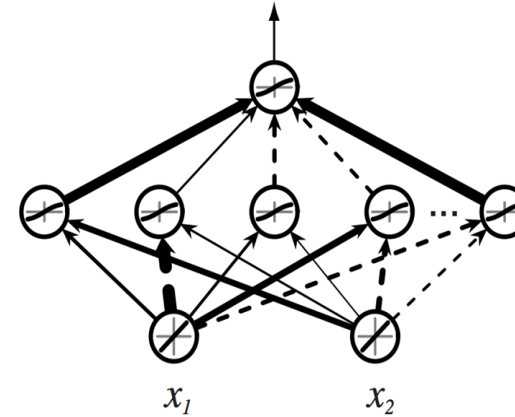
- Classes are coded as $\{\Omega_0 = [0 \ 0 \ 1], \Omega_1 = [0 \ 1 \ 0], \Omega_2 = [1 \ 0 \ 0]\}$

- Output neurons use a softmax unit
$$y_k = \frac{\exp(z_k)}{\sum_k \exp(z_k)} = \frac{\exp(\sum_j w_{kj} h_j(\sum_i^I w_{ji} \cdot x_i))}{\sum_{k=1}^K \exp(\sum_j w_{kj} h_j(\sum_i^I w_{ji} \cdot x_i))}$$

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”

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:

- 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 ...

Optimization and Learning (Supervised learning)

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

- Given a training set

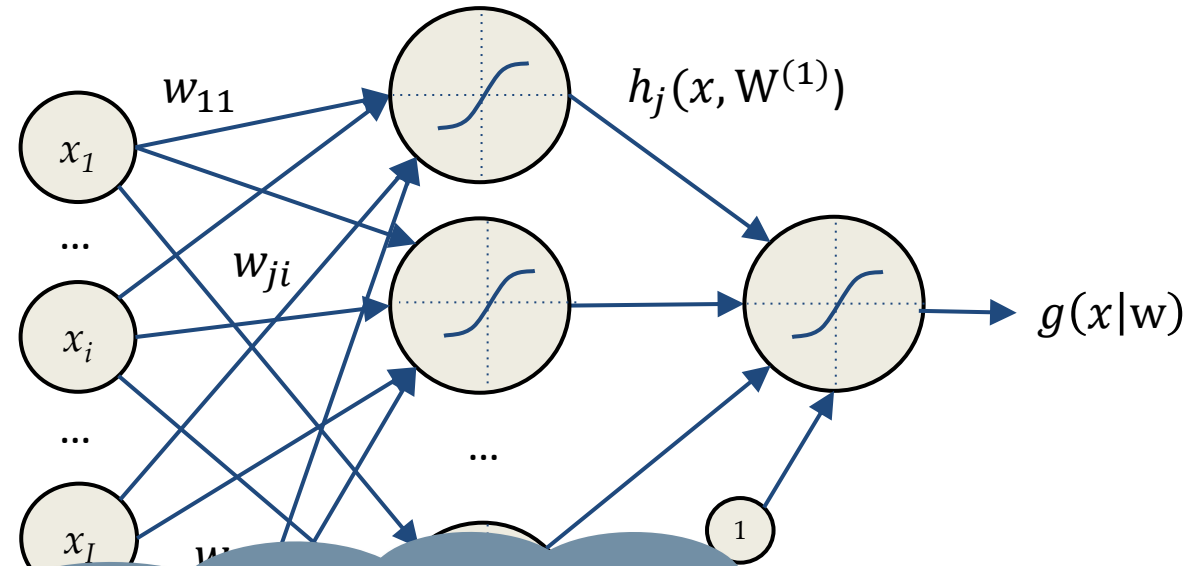
$$D = \langle x_1, t_1 \rangle \cdots \langle x_N, t_N \rangle$$

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

$$y(x_n|\theta) \sim t_n$$

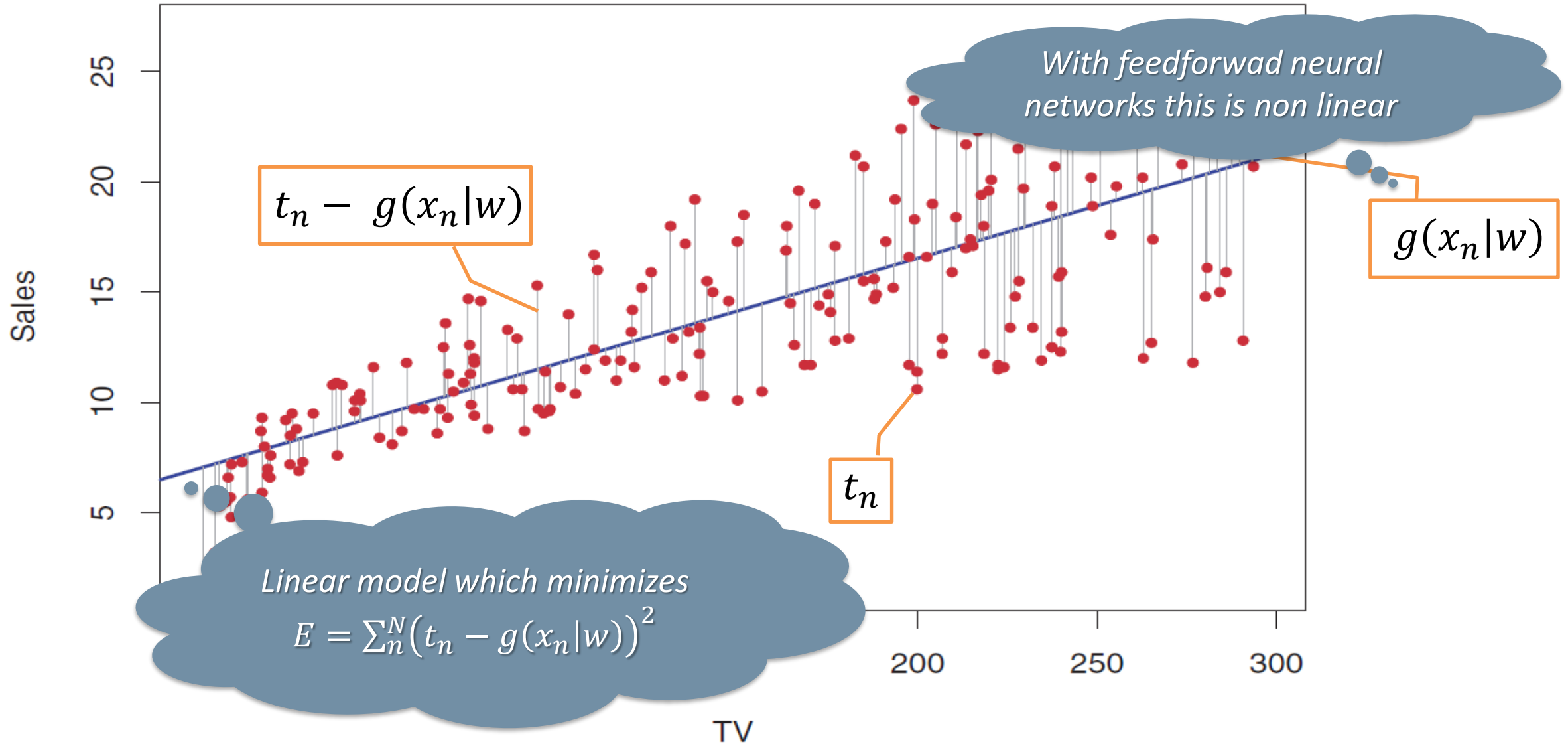
- In case of a Neural Network this can be rewritten as

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



For this you can minimize
 $E = \sum_n^N (t_n - g(x_n|w))^2$

Sum of Squared Errors



Non Linear Optimization 101

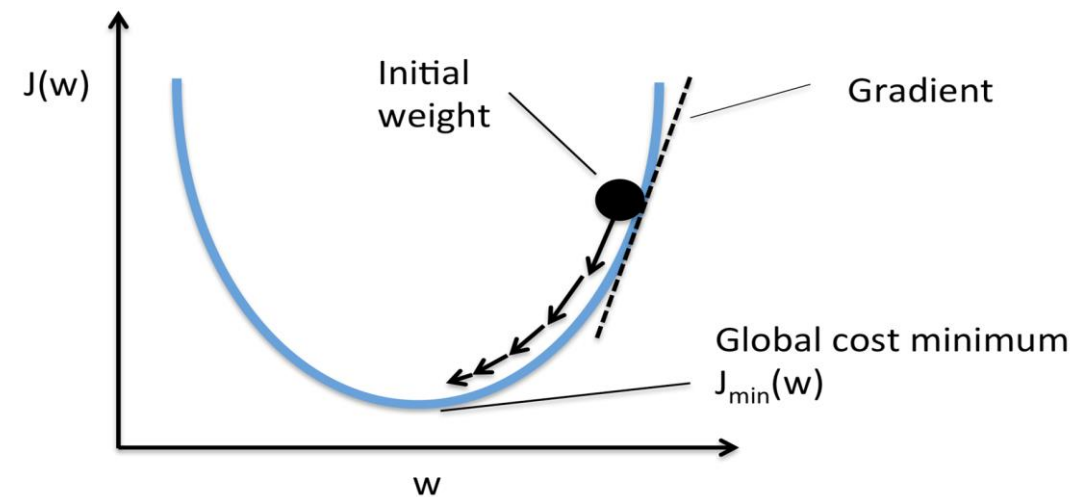
To find the minimum of a generic function, we compute the partial derivatives of the function and set them to zero

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

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

- Initialize the weights to a random value
- Iterate until convergence

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



Gradient descent - Backpropagation

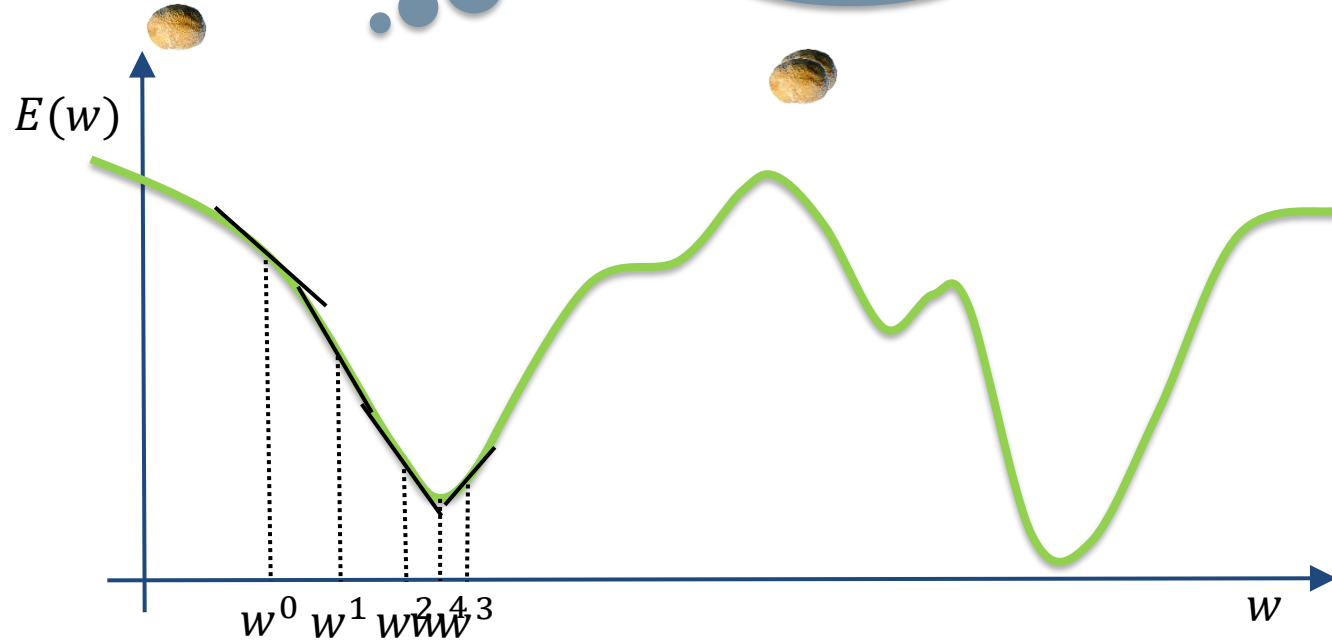
Use multiple restarts to seek for a proper global minimum.

Finding the weights of a Neural Network is a non-convex optimization problem

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

We iterate starting from an initial random configuration

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

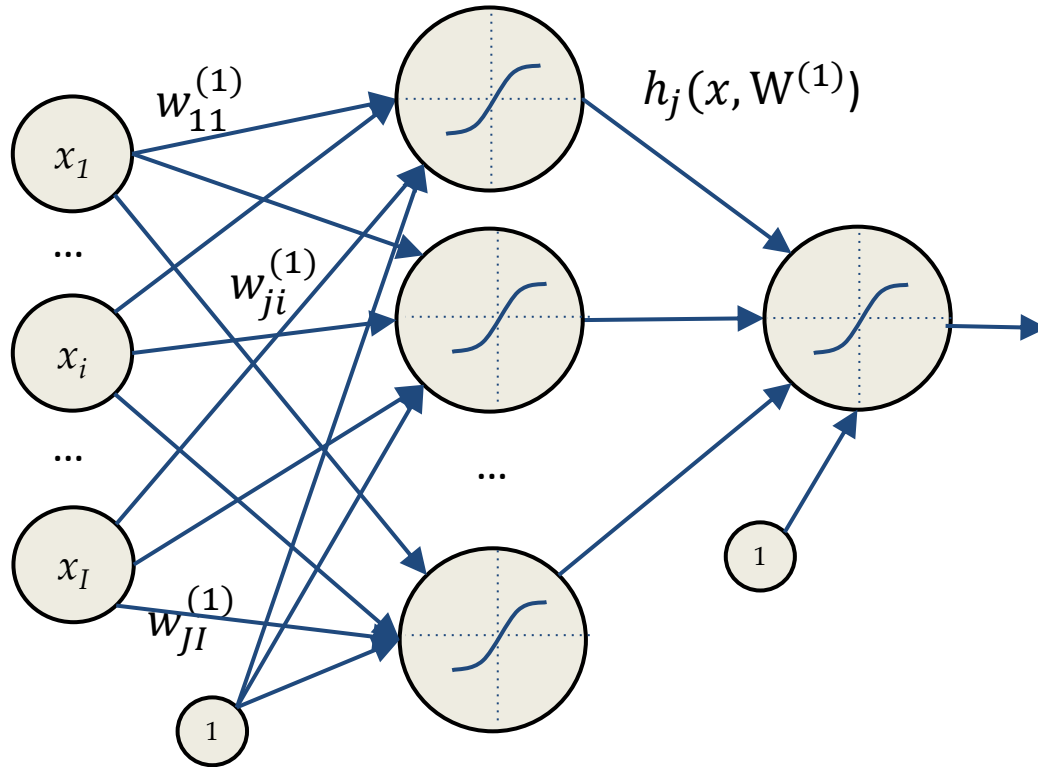


To avoid local minima can use momentum

It depends on where we start from

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

Gradient Descent Example



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

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

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

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



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

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

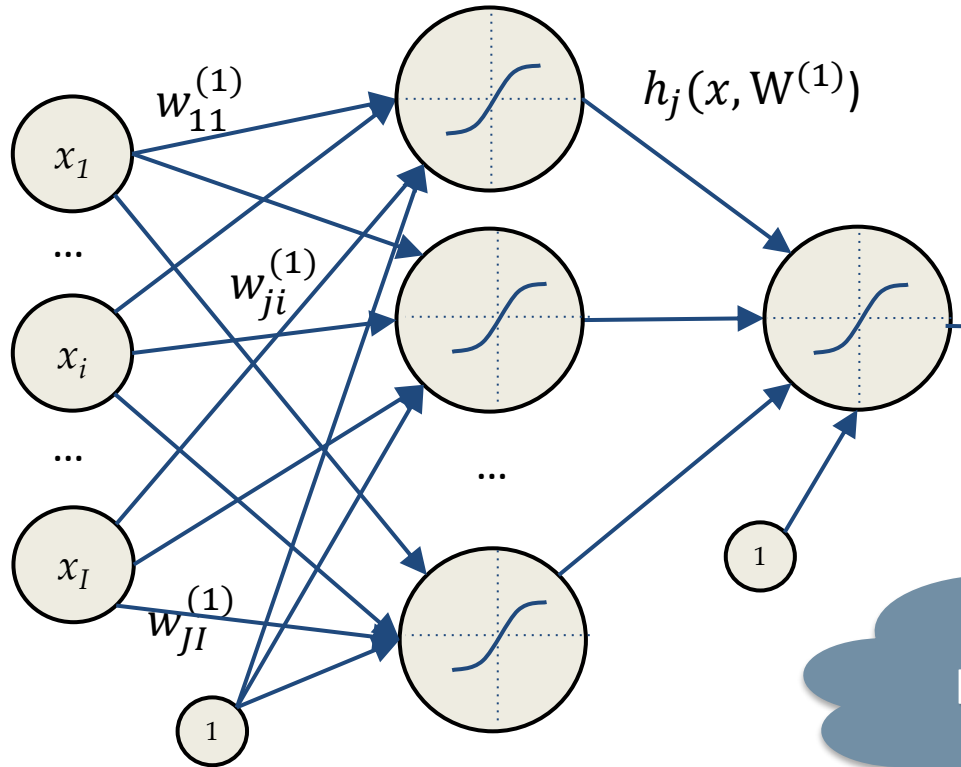
$$\frac{\partial E(w)}{\partial w_{3,5}^{(1)}} = \frac{\partial \sum_{n=1}^N (t_n - g_1(x_n, w))^2}{\partial w_{3,5}^{(1)}} = \sum_{n=1}^N \frac{\partial (t_n - g_1(x_n, w))^2}{\partial w_{3,5}^{(1)}} = -2 \sum_n (t_n - g_1(x_n, w)) \frac{\partial g_1(x_n, w)}{\partial w_{3,5}^{(1)}}$$

$$\frac{\partial g_1(x_n, w)}{\partial w_{3,5}^{(1)}} = \frac{\partial g_1 \left(\sum_{j=0}^J w_{1j}^{(2)} \cdot h_j(\cdot) \right)}{\partial w_{3,5}^{(1)}} = g_1'(x_n, w) \cdot \frac{\partial \sum_{j=0}^J w_{1j}^{(2)} \cdot h_j(\cdot)}{\partial w_{3,5}^{(1)}} = g_1'(x_n, w) \cdot w_{1,3}^{(2)} \cdot \frac{\partial h_3 \left(\sum_{i=0}^I w_{3i}^{(1)} \cdot x_{i,n} \right)}{\partial w_{3,5}^{(1)}}$$

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

$$\frac{\partial E(w)}{\partial w_{3,5}^{(1)}} = -2 \sum_n (t_n - g_1(x_n, w)) g_1'(x_n, w) w_{1,3}^{(2)} h_3' \left(\sum_{i=0}^I w_{3,i}^{(1)} \cdot x_{i,n} \right) x_{5,n}$$

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

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

Gradient Descent Variations

Batch gradient descent

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

Use a single sample, unbiased, but with high variance

Stochastic gradient descent (SGD)

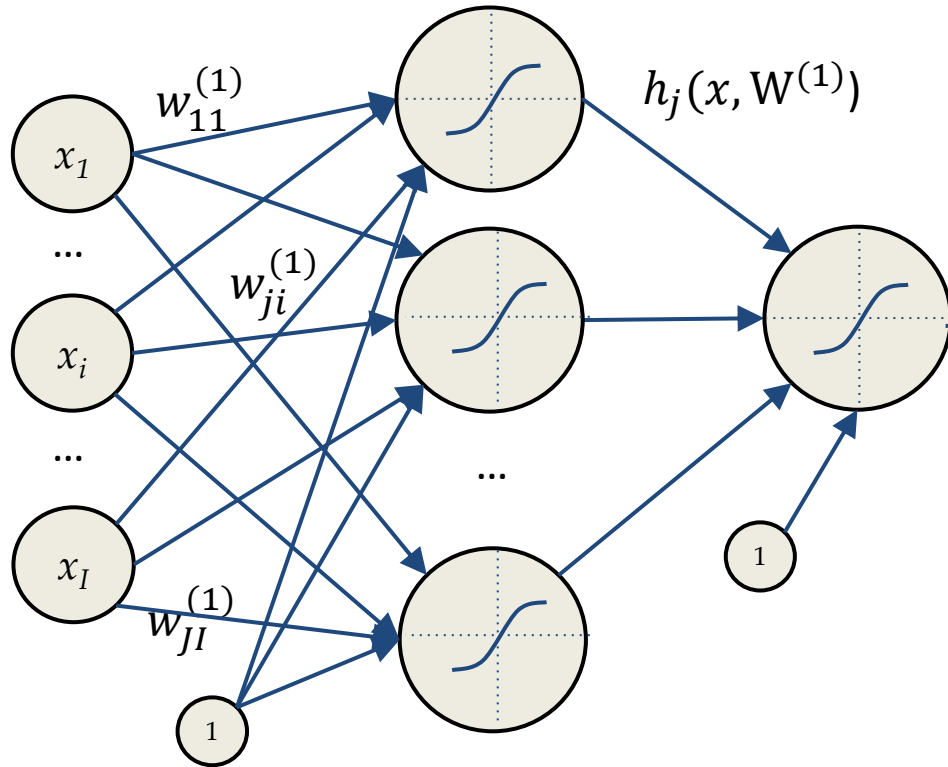
$$\frac{\partial E(w)}{\partial w} \approx \frac{\partial E_{SGD}(w)}{\partial w} = \frac{\partial E(x_n, w)}{\partial w}$$

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

Mini-batch gradient descent

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

Gradient Descent Example



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

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

Can I make it automatic?

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

Backpropagation and Chain Rule (1)

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

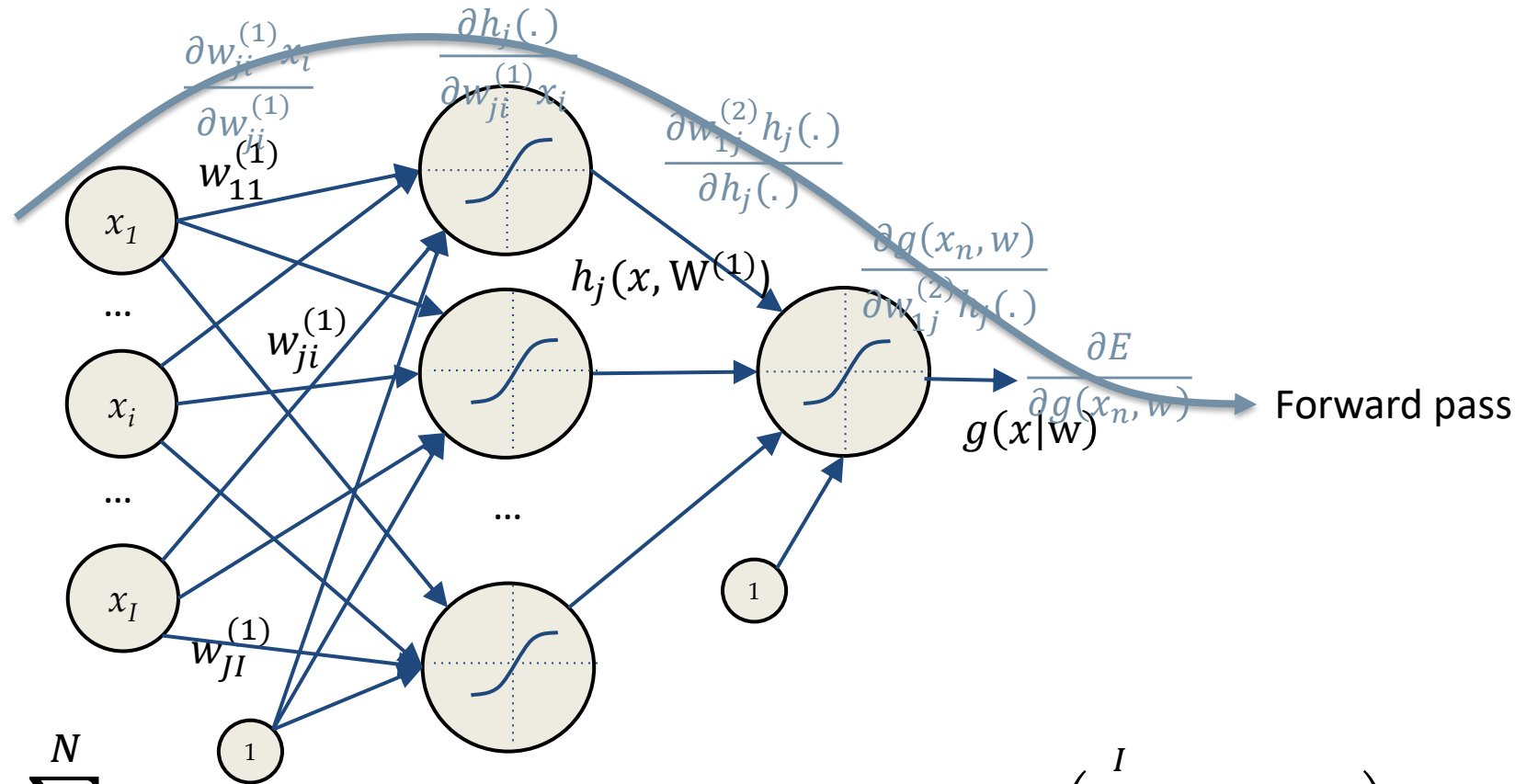
- Let x be a real number and two functions $f: \mathfrak{R} \rightarrow \mathfrak{R}$ and $g: \mathfrak{R} \rightarrow \mathfrak{R}$
- Consider the composed function $z = f(g(x)) = f(y)$ where $y = g(x)$
- The derivative of z w.r.t. x can be computed applying the chain rule

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

The same holds for backpropagation

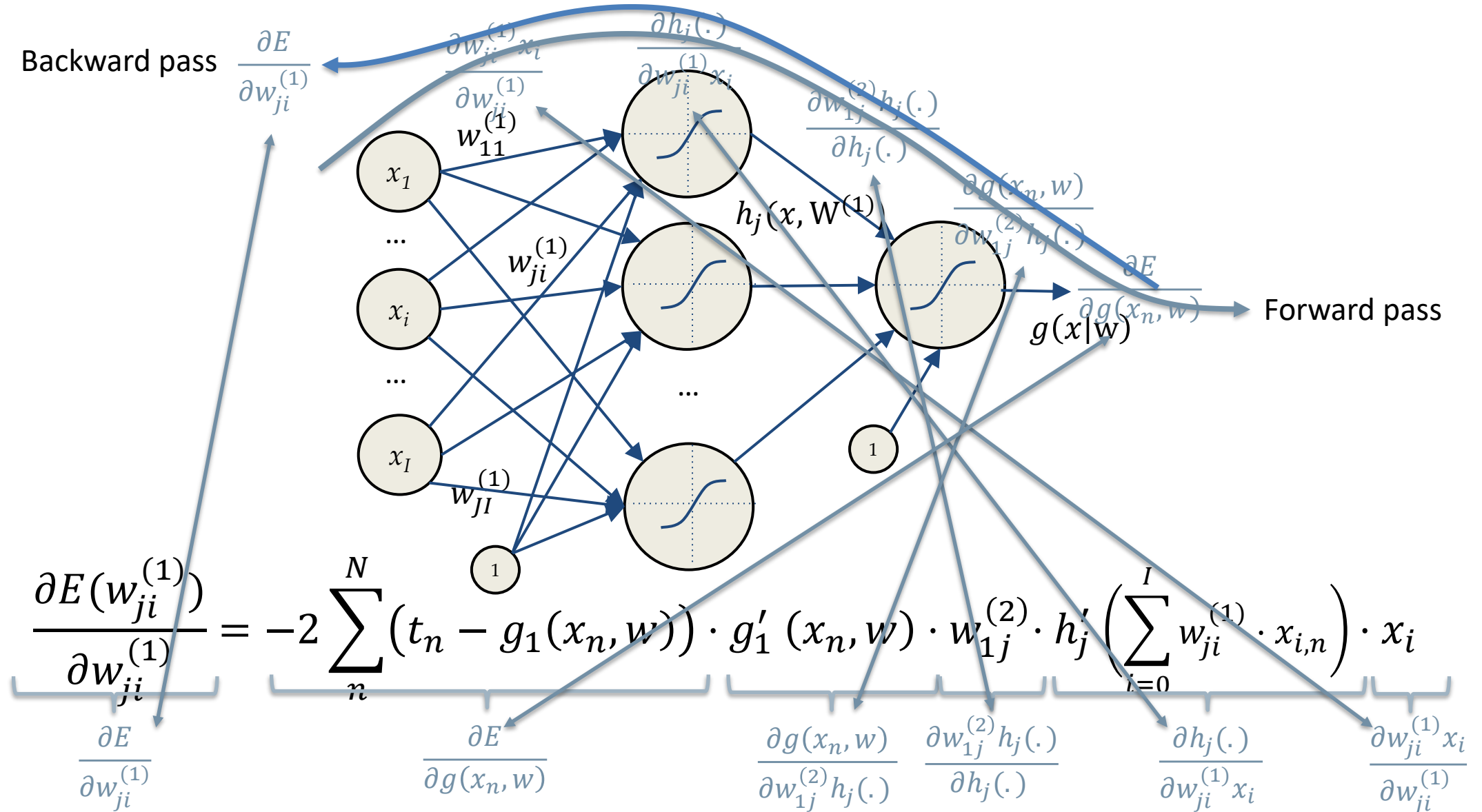
$$\underbrace{\frac{\partial E(w_{ji}^{(1)})}{\partial w_{ji}^{(1)}}}_{\frac{\partial E}{\partial w_{ji}^{(1)}}} = -2 \underbrace{\sum_n}_{\frac{\partial E}{\partial g(x_n, w)}} (t_n - g_1(x_n, w)) \cdot \underbrace{g'_1(x_n, w)}_{\frac{\partial g(x_n, w)}{\partial w_{1j}^{(2)} h_j(\cdot)}} \cdot \underbrace{w_{1j}^{(2)}}_{\frac{\partial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

Backpropagation and Chain Rule (2)

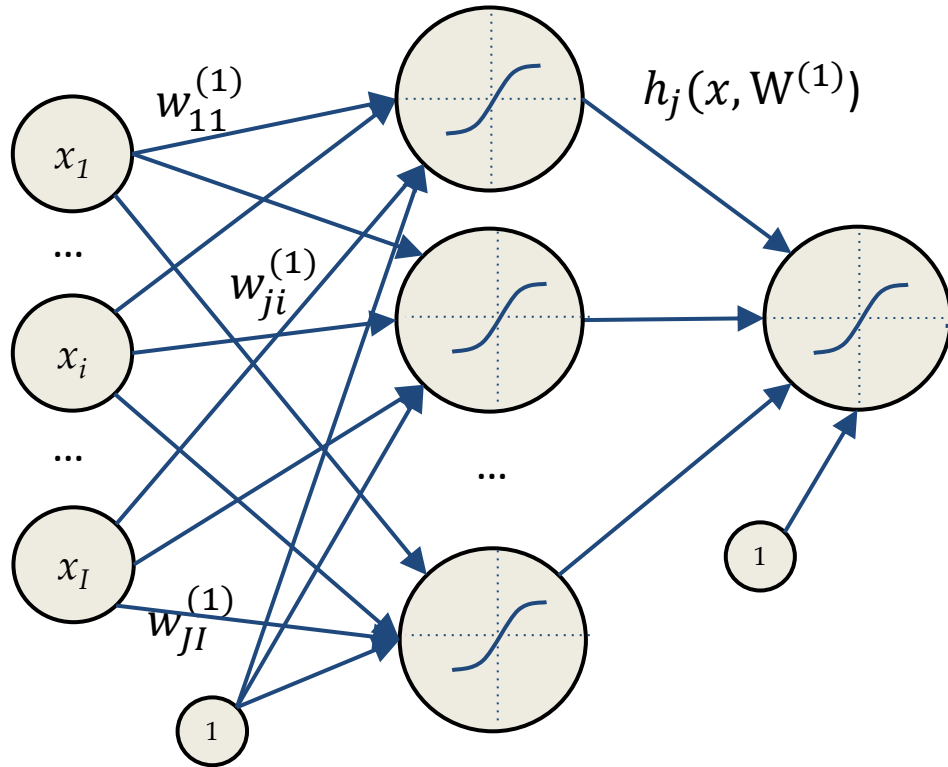


$$\underbrace{\frac{\partial E(w_{ji}^{(1)})}{\partial w_{ji}^{(1)}}}_{\frac{\partial E}{\partial w_{ji}^{(1)}}} = -2 \sum_n \underbrace{(t_n - g_1(x_n, w))}_{\frac{\partial E}{\partial g(x_n, w)}} \cdot \underbrace{g'_1(x_n, w)}_{\frac{\partial g(x_n, w)}{\partial w_{1j}^{(2)} h_j(.)}} \cdot \underbrace{w_{1j}^{(2)}}_{\frac{\partial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

Backpropagation and Chain Rule (2)



Gradient Descent Example



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

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

Why should I use this?

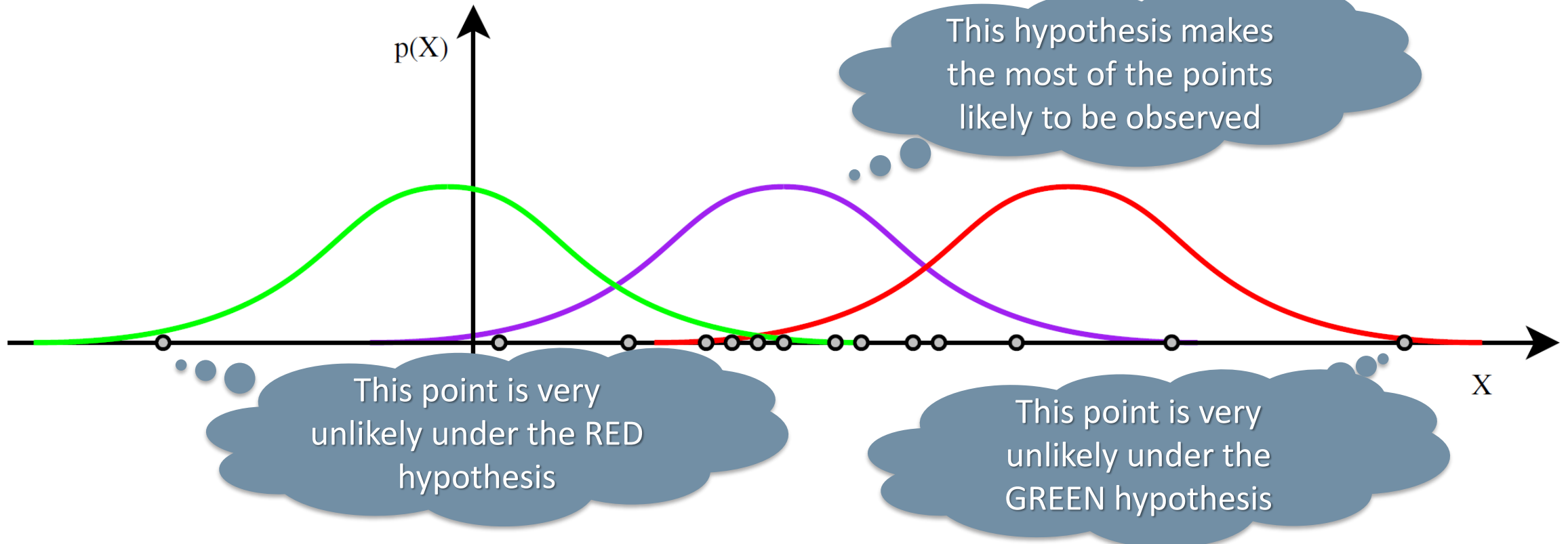
$$\frac{\partial E(w_{ji}^{(1)})}{\partial w_{ji}^{(1)}} = -2 \sum_n (t_n - g_1(x_n, w)) g_1'(x_n, w) w_{1j}^{(2)} h_j' \left(\sum_{i=0}^I w_{ji}^{(1)} \cdot x_{i,n} \right) x_{i,n}$$

A Note on Maximum Likelihood Estimation

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

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

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



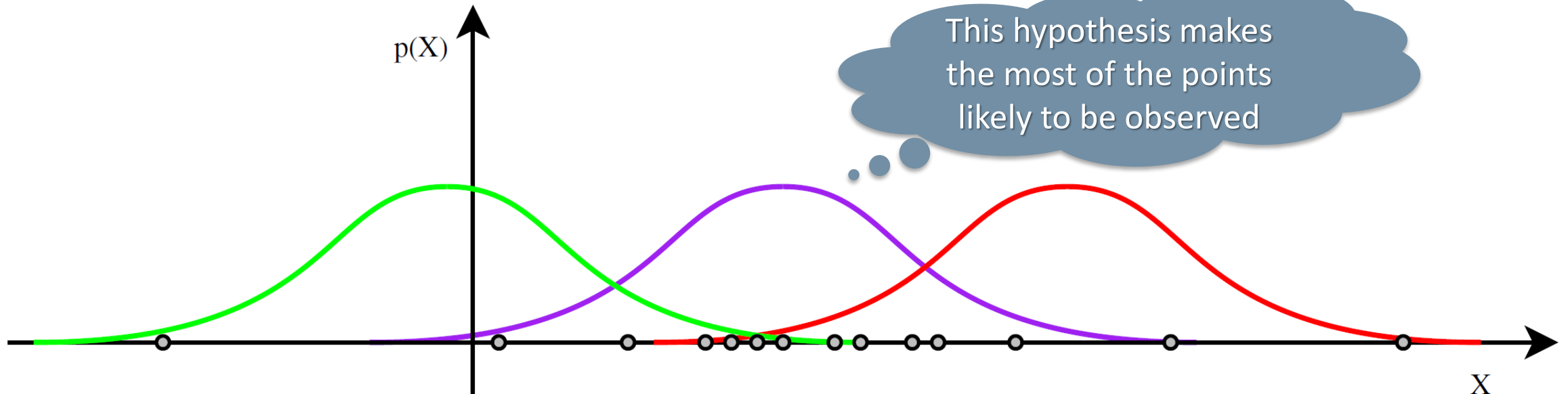
A Note on Maximum Likelihood Estimation

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: Chose 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 $L = P(\text{Data}|\theta)$ for the data
- [Take the logarithm of likelihood $l = \log P(\text{Data}|\theta)$]
- 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

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

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

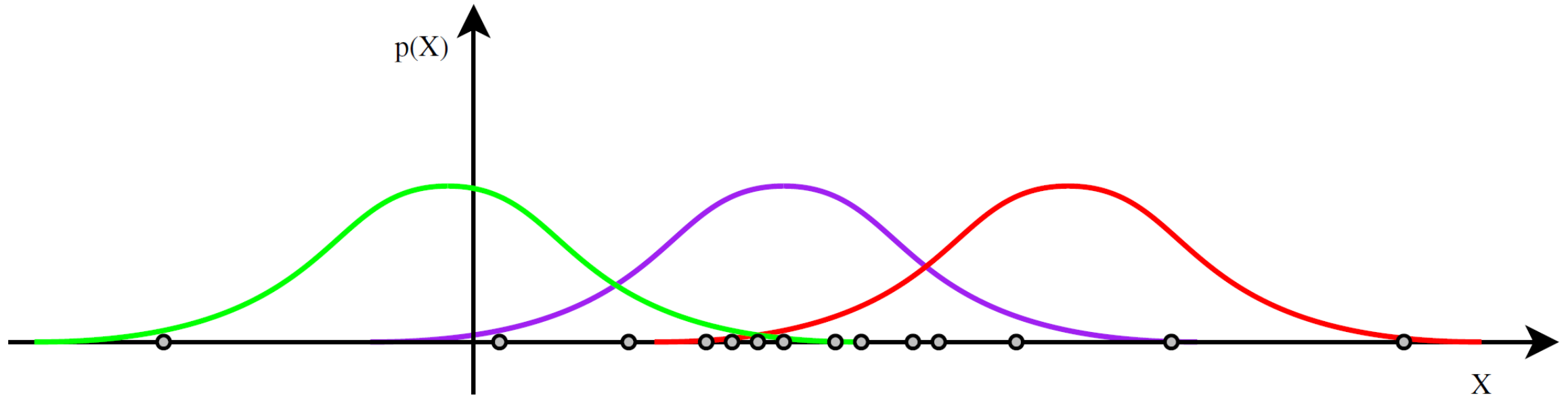
- Analytical Techniques (i.e., solve the equations)
- Optimization Techniques (e.g., Lagrange multipliers)
- Numerical Techniques (e.g., gradient descend)

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

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

- Write the likelihood $L = P(\text{Data}|\theta)$ for the data

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

Maximum Likelihood Estimation Example

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

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

- Take the logarithm $l = \log P(\text{Data}|\theta)$ of the likelihood

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

Maximum Likelihood Estimation Example

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

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

- 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

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

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

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

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

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

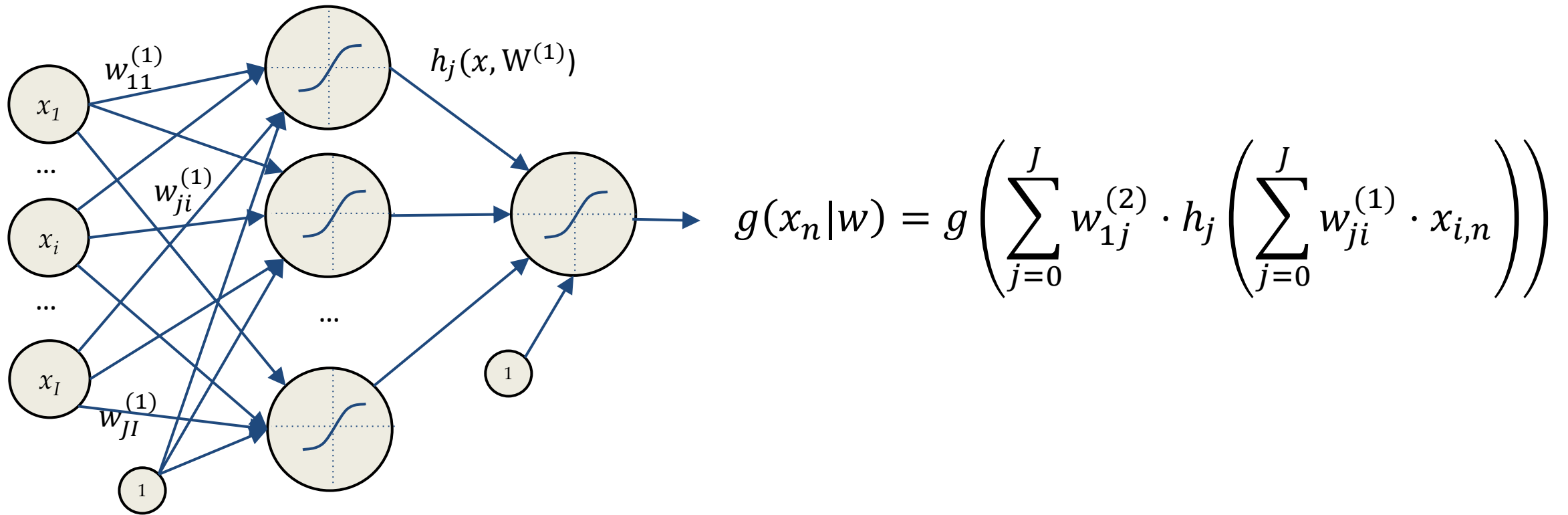
$$\sum_n^N x_n = \sum_n^N \mu$$



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

Let's apply this all to
Neural Networks!

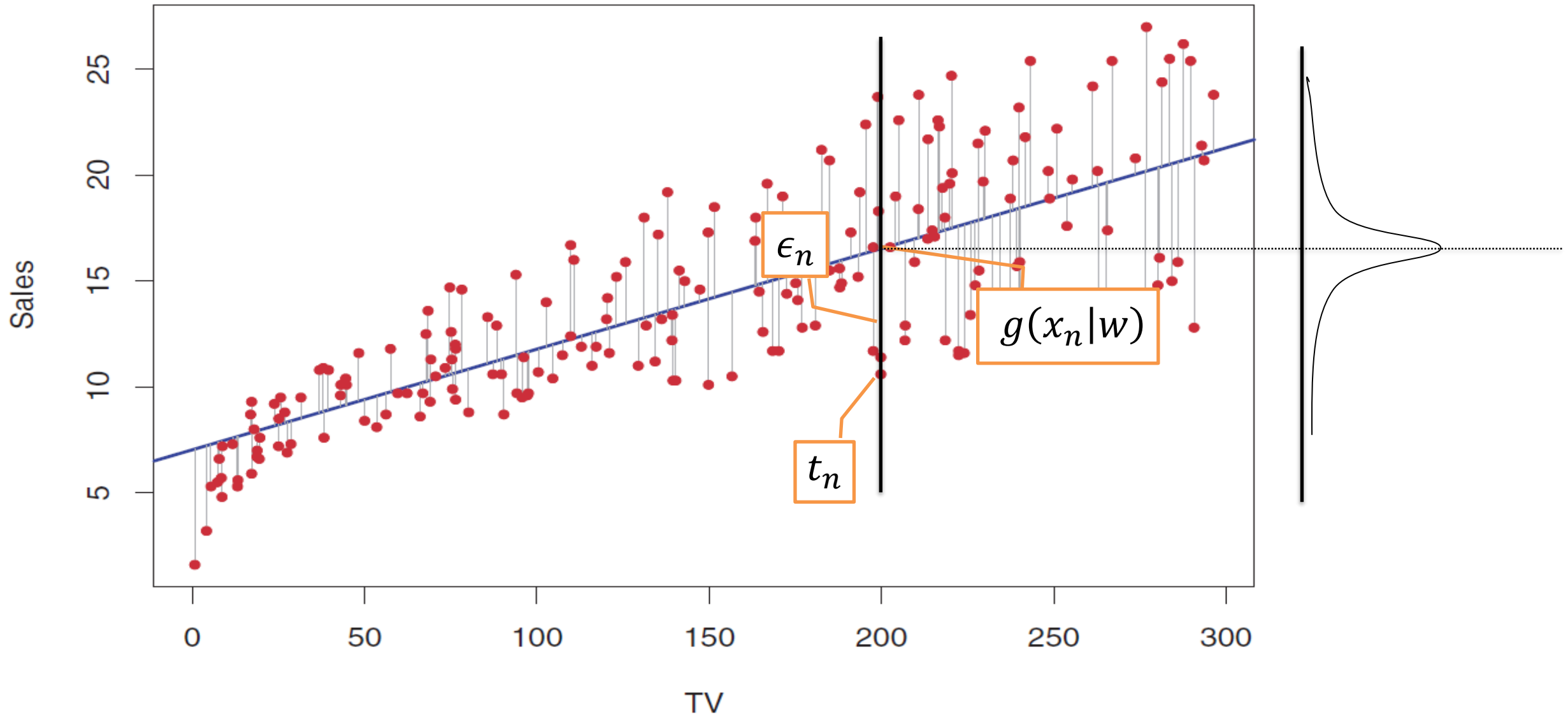
Neural Networks for Regression



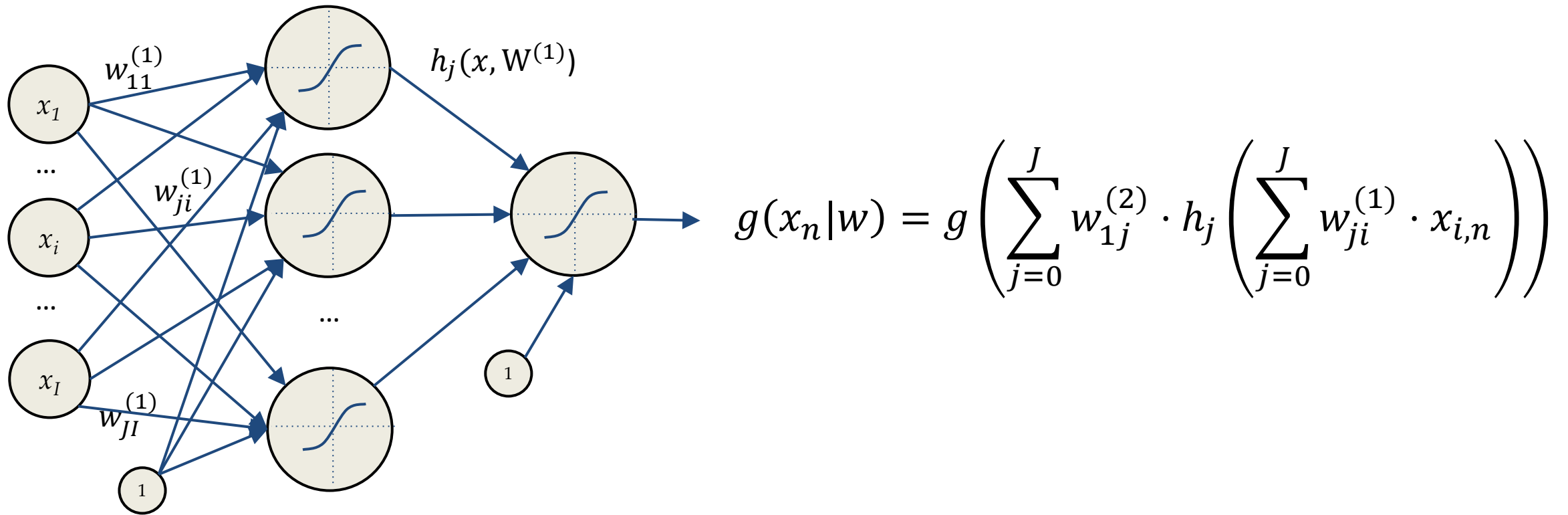
Goal: approximate a target function t having N observations

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

Statistical Learning Framework



Neural Networks for Regression



Goal: approximate a target function t having N observations

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

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 likelihood $L = P(\text{Data}|\theta)$ 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 $\mathbf{l} = \log P(\text{Data}|\theta)$ for the data

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

Maximum Likelihood Estimation for Regression

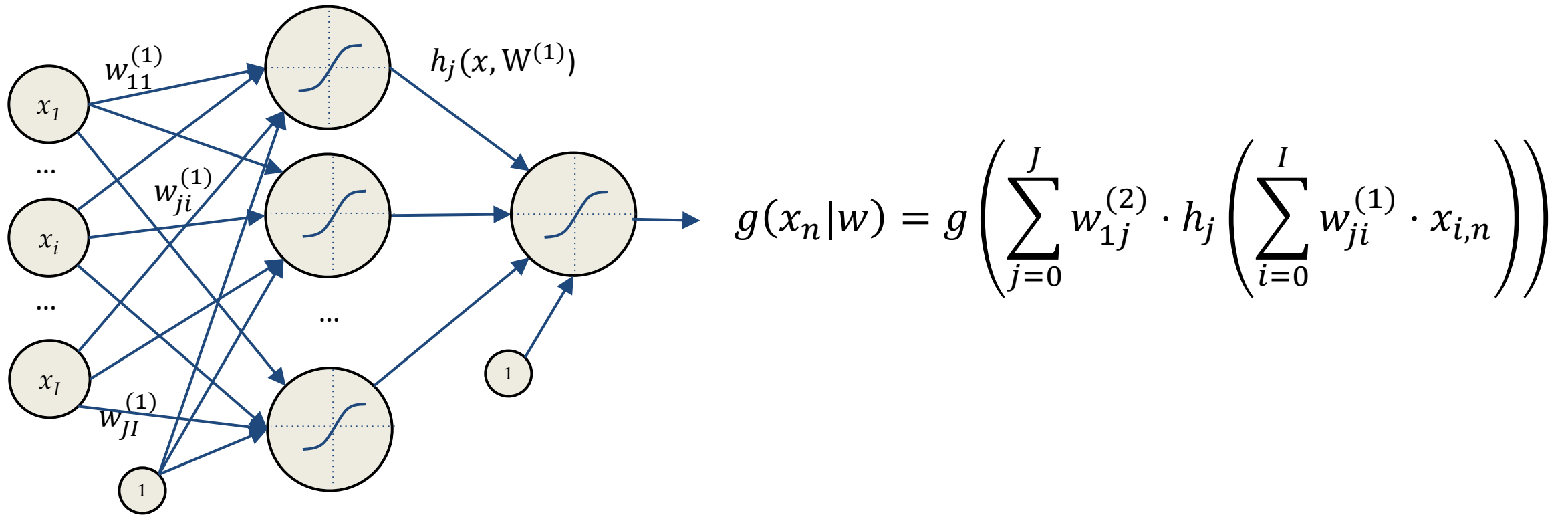
We have i.i.d. samples coming from a Gaussian with known σ^2

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

Look for the weights which maximise the loglikelihood

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

Neural Networks for Classification



Goal: approximate a posterior probability t having N observations

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

Maximum Likelihood Estimation for Classification

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

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

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

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

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

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

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

What about perceptron

How to Chose the Error Function?

We have observed different error functions so far

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

$$E(w) = - \sum_n t_n \log g(x_n|w) + (1 - t_n) \log(1 - g(x_n|w))$$

Sum of Squared Errors

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!

This requires lots of trial and errors ...

As for the Perceptron ...

Hyperplanes Linear Algebra

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

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

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

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

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

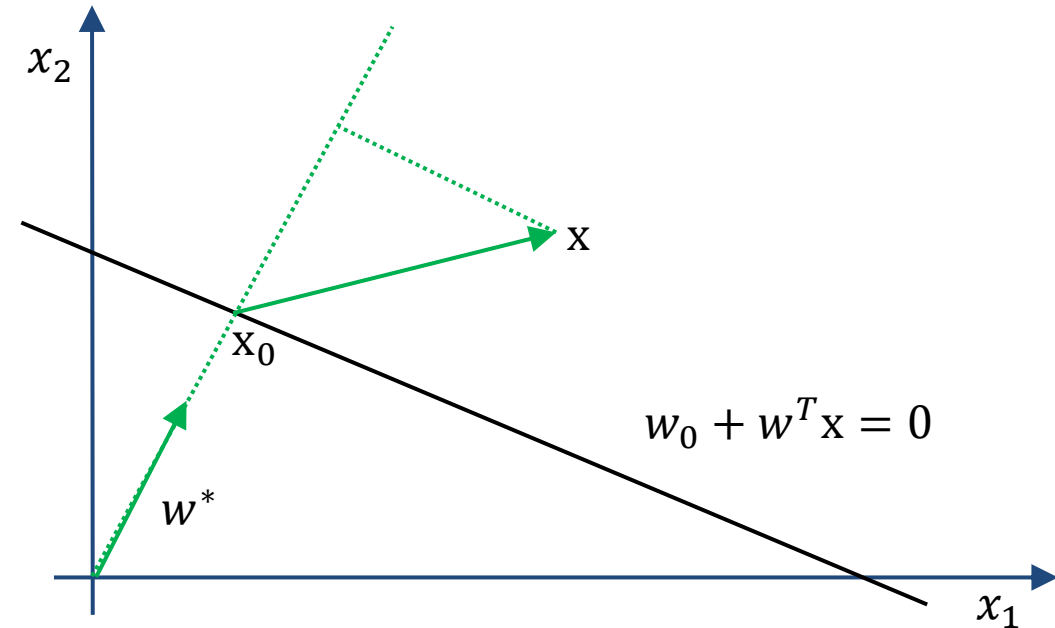
$$w^* = w / \|w\|$$

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

$$w^T x_0 = -w_0$$

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

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



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

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 $\mathbf{w}^T \mathbf{x} + w_0 < 0$
- For an output with -1 we have the opposite

The goal becomes minimizing

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

Set of points
misclassified

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

Perceptron Learning Algorithm (2/2)

Let's minimize by stochastic gradient descent 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