

# **PV021: Neural networks**

**Tomáš Brázdil**

# Course organization

Course materials:

- ▶ **Main:** The lecture
- ▶ Neural Networks and Deep Learning by Michael Nielsen  
<http://neuralnetworksanddeeplearning.com/>  
(Extremely well written online textbook (a little outdated))
- ▶ Deep learning by Ian Goodfellow, Yoshua Bengio and Aaron Courville  
<http://www.deeplearningbook.org/>  
("Classical" overview of the theory of neural networks (a little outdated))
- ▶ Probabilistic Machine Learning: An Introduction by Kevin Murphy  
<https://probml.github.io/pml-book/book1.html>  
(Great advanced ML textbook with (almost) up-to-date basic neural networks.)
- ▶ Infinitely many online tutorials on everything (to build intuition)

**Suggested:** deeplearning.ai courses by Andrew Ng

## Evaluation:

- ▶ Project (Dr. Tomáš Foltýnek)
  - ▶ implementation of a selected model + analysis of given data
  - ▶ implementation C/C++/Java/Rust **without use of any specialized libraries for data analysis and machine learning**
  - ▶ need to get over a given accuracy threshold (a gentle one, just to eliminate non-functional implementations)

## Evaluation:

- ▶ Project (Dr. Tomáš Foltýnek)
  - ▶ implementation of a selected model + analysis of given data
  - ▶ implementation C/C++/Java/Rust **without use of any specialized libraries for data analysis and machine learning**
  - ▶ need to get over a given accuracy threshold (a gentle one, just to eliminate non-functional implementations)
- ▶ Oral exam
  - ▶ I may ask about anything from the lecture! You will get a detailed manual specifying the mandatory knowledge.

**Q:** Why we cannot use specialized libraries in projects?

**Q:** Why we cannot use specialized libraries in projects?

**A:** In order to "touch" the low level implementation details of the algorithms. You should not even use libraries for linear algebra and numerical methods, so that you will be confronted with rounding errors and numerical instabilities.

**Q:** Why we cannot use specialized libraries in projects?

**A:** In order to "touch" the low level implementation details of the algorithms. You should not even use libraries for linear algebra and numerical methods, so that you will be confronted with rounding errors and numerical instabilities.

**Q:** Why should you attend this course when there are infinitely many great reasources elsewhere?

**A:** There are at least two reasons:

- ▶ You may discuss issues with me, my colleagues and other students.
- ▶ I will make you truly learn fundamentals by heart.

# Notable features of the course

- ▶ Use of mathematical notation and reasoning (mandatory for the exam)
- ▶ Sometimes goes deeper into statistical underpinnings of neural networks learning
- ▶ The project demands a complete working solution which must satisfy a prescribed performance specification



# Notable features of the course

- ▶ Use of mathematical notation and reasoning (mandatory for the exam)
- ▶ Sometimes goes deeper into statistical underpinnings of neural networks learning
- ▶ The project demands a complete working solution which must satisfy a prescribed performance specification

An unusual exam system! You can repeat the oral exam as many times as needed (only the best grade goes into IS).

# Notable features of the course

- ▶ Use of mathematical notation and reasoning (mandatory for the exam)
- ▶ Sometimes goes deeper into statistical underpinnings of neural networks learning
- ▶ The project demands a complete working solution which must satisfy a prescribed performance specification

An unusual exam system! You can repeat the oral exam as many times as needed (only the best grade goes into IS).

An example of an instruction email (from another course with the same system):

It is typically not sufficient to devote a single afternoon to the preparation for the exam.

You have to know `_everything_` (which means every single thing) starting with the slide 42 and ending with the slide 245 with notable exceptions of slides: 121 - 123, 137 - 140, 165, 167.

Proofs presented on the whiteboard are also mandatory.

# Machine learning in general

- ▶ Machine learning = construction of systems that learn their functionality from data  
(... and thus do not need to be programmed.)

# Machine learning in general

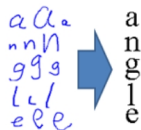
- ▶ Machine learning = construction of systems that learn their functionality from data  
(... and thus do not need to be programmed.)
  - ▶ spam filter
    - ▶ learns to recognize spam from a database of "labelled" emails
    - ▶ consequently is able to distinguish spam from ham

# Machine learning in general

- ▶ Machine learning = construction of systems that learn their functionality from data

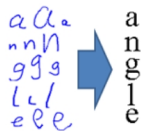
(... and thus do not need to be programmed.)

- ▶ spam filter
  - ▶ learns to recognize spam from a database of "labelled" emails
  - ▶ consequently is able to distinguish spam from ham
- ▶ handwritten text reader
  - ▶ learns from a database of handwritten letters (or text) labelled by their correct meaning
  - ▶ consequently is able to recognize text



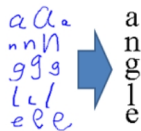
# Machine learning in general

- ▶ Machine learning = construction of systems that learn their functionality from data  
(... and thus do not need to be programmed.)
  - ▶ spam filter
    - ▶ learns to recognize spam from a database of "labelled" emails
    - ▶ consequently is able to distinguish spam from ham
  - ▶ handwritten text reader
    - ▶ learns from a database of handwritten letters (or text) labelled by their correct meaning
    - ▶ consequently is able to recognize text
  - ▶ ...
  - ▶ and lots of much much more sophisticated applications ...



# Machine learning in general

- ▶ Machine learning = construction of systems that learn their functionality from data  
(... and thus do not need to be programmed.)
  - ▶ spam filter
    - ▶ learns to recognize spam from a database of "labelled" emails
    - ▶ consequently is able to distinguish spam from ham
  - ▶ handwritten text reader
    - ▶ learns from a database of handwritten letters (or text) labelled by their correct meaning
    - ▶ consequently is able to recognize text
  - ▶ ...
  - ▶ and lots of much much more sophisticated applications ...
- ▶ Basic attributes of learning algorithms:
  - ▶ **representation**: ability to capture the inner structure of training data
  - ▶ **generalization**: ability to work properly on new data



# Machine learning in general

**Machine learning algorithms** typically construct mathematical models of given data. The models may be subsequently applied to fresh data.



# Machine learning in general

**Machine learning algorithms** typically construct mathematical models of given data. The models may be subsequently applied to fresh data.

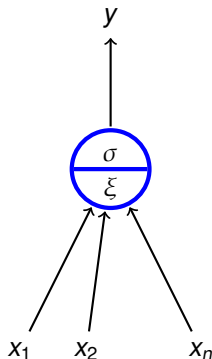
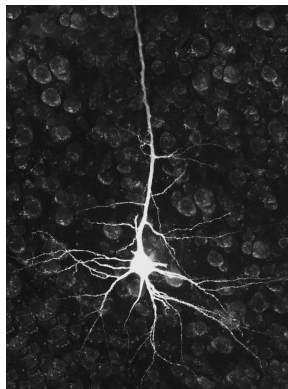
There are many types of models:

- ▶ decision trees
- ▶ support vector machines
- ▶ hidden Markov models
- ▶ Bayes networks and other graphical models
- ▶ **neural networks**
- ▶ ...

Neural networks, based on models of a (human) brain, form a natural basis for learning algorithms!

# Artificial neural networks

- ▶ **Artificial neuron** is a *rough mathematical approximation* of a biological neuron.
- ▶ **(Artificial) neural network (NN)** consists of a number of interconnected artificial neurons. "Behavior" of the network is encoded in connections between neurons.



# Why artificial neural networks?

Modelling of biological neural networks (computational neuroscience).

- ▶ simplified mathematical models help to identify important mechanisms
  - ▶ How the brain receives information?
  - ▶ How the information is stored?
  - ▶ How the brain develops?
  - ▶ ...

# Why artificial neural networks?

Modelling of biological neural networks (computational neuroscience).

- ▶ simplified mathematical models help to identify important mechanisms
  - ▶ How the brain receives information?
  - ▶ How the information is stored?
  - ▶ How the brain develops?
  - ▶ ...
- ▶ neuroscience is strongly multidisciplinary; precise mathematical descriptions help in communication among experts and in design of new experiments.

I will not spend much time on this area!

# Why artificial neural networks?

Neural networks in machine learning.

- ▶ Typically primitive models, far from their biological counterparts (but often inspired by biology).

# Why artificial neural networks?

Neural networks in machine learning.

- ▶ Typically primitive models, far from their biological counterparts (but often inspired by biology).
- ▶ Strongly oriented towards concrete application domains:
  - ▶ decision making and control - autonomous vehicles, manufacturing processes, control of natural resources
  - ▶ games - backgammon, poker, GO, Starcraft, ...
  - ▶ finance - stock prices, risk analysis
  - ▶ medicine - diagnosis, signal processing (EKG, EEG, ...), image processing (MRI, CT, WSI ...)
  - ▶ text and speech processing - machine translation, *text generation*, speech recognition
  - ▶ other signal processing - filtering, radar tracking, noise reduction
  - ▶ art - music and painting generation, deepfakes
  - ▶ ...

I will concentrate on this area!

# Important features of neural networks

- ▶ Massive parallelism
  - ▶ many slow (and "dumb") computational elements work in parallel on several levels of abstraction

# Important features of neural networks

- ▶ Massive parallelism
  - ▶ many slow (and "dumb") computational elements work in parallel on several levels of abstraction
- ▶ Learning
  - ▶ a kid learns to recognize a rabbit after seeing several rabbits



# Important features of neural networks

- ▶ Massive parallelism
  - ▶ many slow (and "dumb") computational elements work in parallel on several levels of abstraction
- ▶ Learning
  - ▶ a kid learns to recognize a rabbit after seeing several rabbits
- ▶ Generalization
  - ▶ a kid is able to recognize a new rabbit after seeing several (old) rabbits

# Important features of neural networks

- ▶ Massive parallelism
  - ▶ many slow (and "dumb") computational elements work in parallel on several levels of abstraction
- ▶ Learning
  - ▶ a kid learns to recognize a rabbit after seeing several rabbits
- ▶ Generalization
  - ▶ a kid is able to recognize a new rabbit after seeing several (old) rabbits
- ▶ Robustness
  - ▶ a blurred photo of a rabbit may still be classified as an image of a rabbit

# Important features of neural networks

- ▶ Massive parallelism
  - ▶ many slow (and "dumb") computational elements work in parallel on several levels of abstraction
- ▶ Learning
  - ▶ a kid learns to recognize a rabbit after seeing several rabbits
- ▶ Generalization
  - ▶ a kid is able to recognize a new rabbit after seeing several (old) rabbits
- ▶ Robustness
  - ▶ a blurred photo of a rabbit may still be classified as an image of a rabbit
- ▶ Graceful degradation
  - ▶ Experiments have shown that damaged neural network is still able to work quite well
  - ▶ Damaged network may re-adapt, remaining neurons may take on functionality of the damaged ones

# The aim of the course

- ▶ We will concentrate on
  - ▶ basic techniques and principles of neural networks,
  - ▶ fundamental models of neural networks and their applications.
- ▶ You should learn
  - ▶ basic models  
(multilayer perceptron, convolutional networks, recurrent networks, transformers, autoencoders and generative adversarial networks)

# The aim of the course

- ▶ We will concentrate on
  - ▶ basic techniques and principles of neural networks,
  - ▶ fundamental models of neural networks and their applications.
- ▶ You should learn
  - ▶ basic models  
(multilayer perceptron, convolutional networks, recurrent networks, transformers, autoencoders and generative adversarial networks)
  - ▶ Simple applications of these models  
(image processing, a little bit of text processing)

# The aim of the course

- ▶ We will concentrate on
  - ▶ basic techniques and principles of neural networks,
  - ▶ fundamental models of neural networks and their applications.
- ▶ You should learn
  - ▶ basic models  
(multilayer perceptron, convolutional networks, recurrent networks, transformers, autoencoders and generative adversarial networks)
  - ▶ Simple applications of these models  
(image processing, a little bit of text processing)
  - ▶ Basic learning algorithms  
(gradient descent with backpropagation)

# The aim of the course

- ▶ We will concentrate on
  - ▶ basic techniques and principles of neural networks,
  - ▶ fundamental models of neural networks and their applications.
- ▶ You should learn
  - ▶ basic models  
(multilayer perceptron, convolutional networks, recurrent networks, transformers, autoencoders and generative adversarial networks)
  - ▶ Simple applications of these models  
(image processing, a little bit of text processing)
  - ▶ Basic learning algorithms  
(gradient descent with backpropagation)
  - ▶ Basic practical training techniques  
(data preparation, setting various hyper-parameters, control of learning, improving generalization)

# The aim of the course

- ▶ We will concentrate on
  - ▶ basic techniques and principles of neural networks,
  - ▶ fundamental models of neural networks and their applications.
- ▶ You should learn
  - ▶ basic models  
(multilayer perceptron, convolutional networks, recurrent networks, transformers, autoencoders and generative adversarial networks)
  - ▶ Simple applications of these models  
(image processing, a little bit of text processing)
  - ▶ Basic learning algorithms  
(gradient descent with backpropagation)
  - ▶ Basic practical training techniques  
(data preparation, setting various hyper-parameters, control of learning, improving generalization)
  - ▶ Basic information about current implementations  
(TensorFlow-Keras, Pytorch)



# Biological neural network

- ▶ Human neural network consists of approximately  $10^{11}$  (100 billion on the short scale) neurons; a single cubic centimeter of a human brain contains almost 50 million neurons.
- ▶ Each neuron is connected with approx.  $10^4$  neurons.
- ▶ Neurons themselves are very complex systems.

# Biological neural network

- ▶ Human neural network consists of approximately  $10^{11}$  (100 billion on the short scale) neurons; a single cubic centimeter of a human brain contains almost 50 million neurons.
- ▶ Each neuron is connected with approx.  $10^4$  neurons.
- ▶ Neurons themselves are very complex systems.

Rough description of nervous system:

- ▶ External stimulus is received by *sensory receptors* (e.g. eye cells).

# Biological neural network

- ▶ Human neural network consists of approximately  $10^{11}$  (100 billion on the short scale) neurons; a single cubic centimeter of a human brain contains almost 50 million neurons.
- ▶ Each neuron is connected with approx.  $10^4$  neurons.
- ▶ Neurons themselves are very complex systems.

Rough description of nervous system:

- ▶ External stimulus is received by *sensory receptors* (e.g. eye cells).
- ▶ Information is further transferred via peripheral nervous system (PNS) to the central nervous systems (CNS) where it is processed (integrated), and subsequently, an output signal is produced.

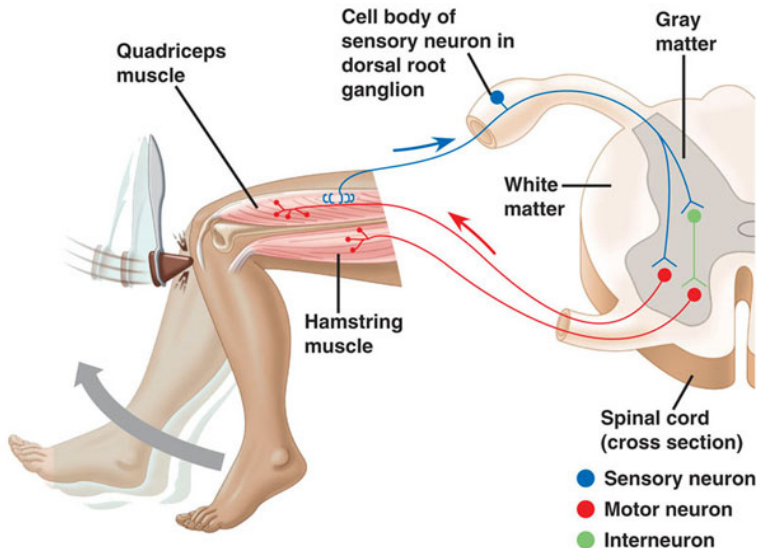
# Biological neural network

- ▶ Human neural network consists of approximately  $10^{11}$  (100 billion on the short scale) neurons; a single cubic centimeter of a human brain contains almost 50 million neurons.
- ▶ Each neuron is connected with approx.  $10^4$  neurons.
- ▶ Neurons themselves are very complex systems.

Rough description of nervous system:

- ▶ External stimulus is received by *sensory receptors* (e.g. eye cells).
- ▶ Information is further transferred via peripheral nervous system (PNS) to the central nervous systems (CNS) where it is processed (integrated), and subsequently, an output signal is produced.
- ▶ Afterwards, the output signal is transferred via PNS to *effectors* (e.g. muscle cells).

# Biological neural network



# Summation

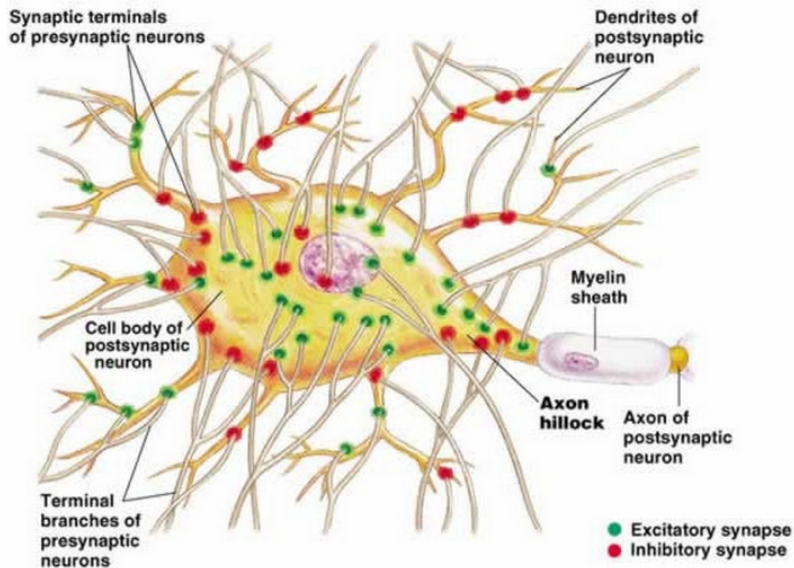
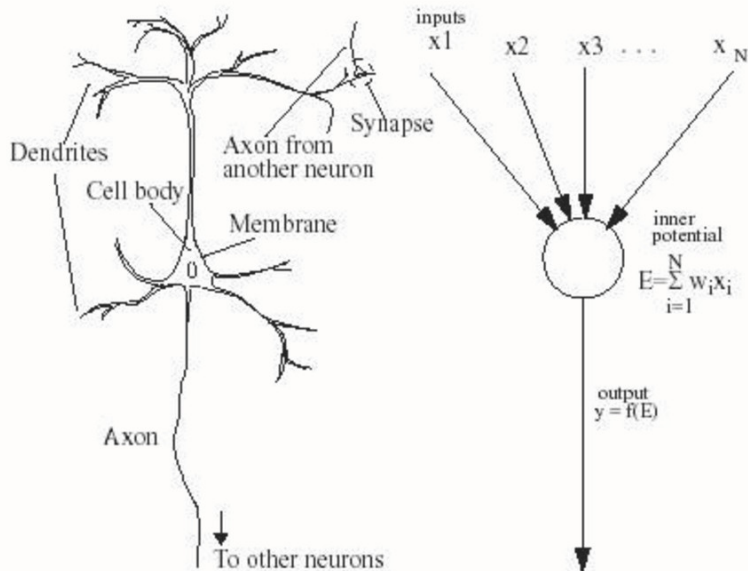


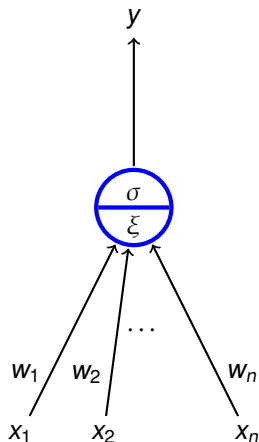
Figure 48.11(a), page 972, Campbell's *Biology*, 5th Edition

# Biological and Mathematical neurons



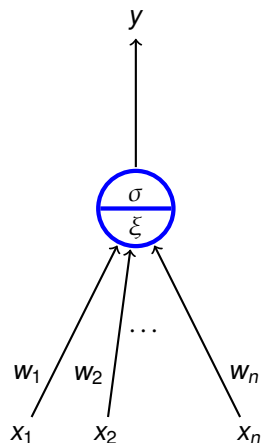
# Formal neuron (without bias)

- $x_1, \dots, x_n \in \mathbb{R}$  are **inputs**



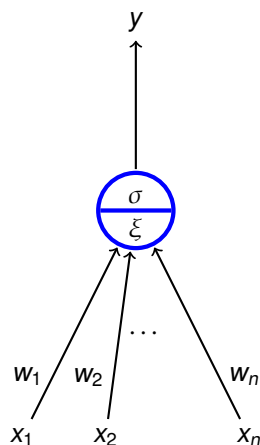


## Formal neuron (without bias)



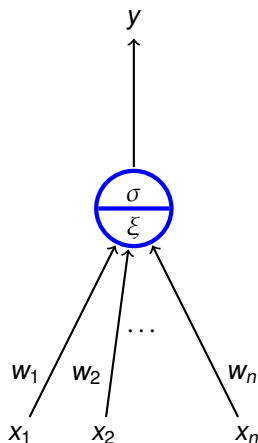
- ▶  $x_1, \dots, x_n \in \mathbb{R}$  are **inputs**
- ▶  $w_1, \dots, w_n \in \mathbb{R}$  are **weights**

# Formal neuron (without bias)



- ▶  $x_1, \dots, x_n \in \mathbb{R}$  are **inputs**
- ▶  $w_1, \dots, w_n \in \mathbb{R}$  are **weights**
- ▶  $\xi$  is an **inner potential**;  
almost always  $\xi = \sum_{i=1}^n w_i x_i$

# Formal neuron (without bias)



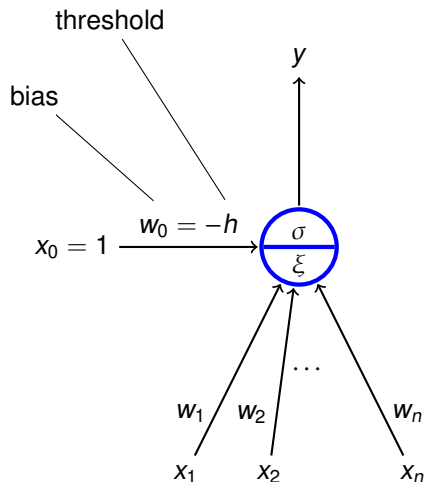
- ▶  $x_1, \dots, x_n \in \mathbb{R}$  are **inputs**
- ▶  $w_1, \dots, w_n \in \mathbb{R}$  are **weights**
- ▶  $\xi$  is an **inner potential**;  
almost always  $\xi = \sum_{i=1}^n w_i x_i$
- ▶  $y$  is an **output** given by  $y = \sigma(\xi)$   
where  $\sigma$  is an **activation function**;  
e.g. a *unit step function*

$$\sigma(\xi) = \begin{cases} 1 & \xi \geq h; \\ 0 & \xi < h. \end{cases}$$

where  $h \in \mathbb{R}$  is a *threshold*.

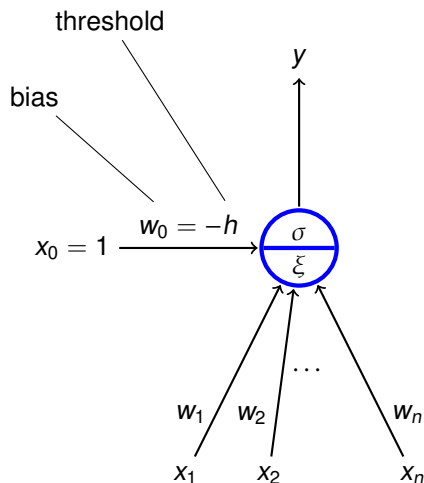
# Formal neuron (with bias)

- $x_0 = 1, x_1, \dots, x_n \in \mathbb{R}$  are **inputs**

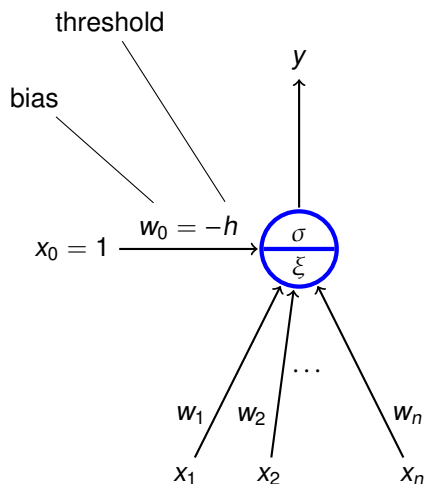


# Formal neuron (with bias)

- ▶  $x_0 = 1, x_1, \dots, x_n \in \mathbb{R}$  are **inputs**
- ▶  $w_0, w_1, \dots, w_n \in \mathbb{R}$  are **weights**

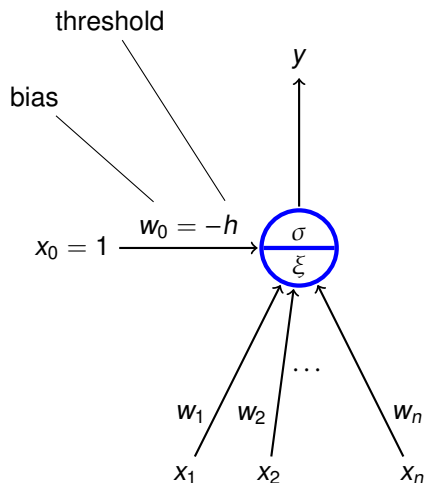


# Formal neuron (with bias)



- ▶  $x_0 = 1, x_1, \dots, x_n \in \mathbb{R}$  are **inputs**
- ▶  $w_0, w_1, \dots, w_n \in \mathbb{R}$  are **weights**
- ▶  $\xi$  is an **inner potential**;  
almost always  $\xi = w_0 + \sum_{i=1}^n w_i x_i$

# Formal neuron (with bias)

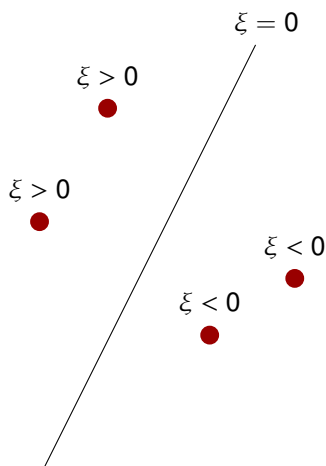


- ▶  $x_0 = 1, x_1, \dots, x_n \in \mathbb{R}$  are **inputs**
- ▶  $w_0, w_1, \dots, w_n \in \mathbb{R}$  are **weights**
- ▶  $\xi$  is an **inner potential**;  
almost always  $\xi = w_0 + \sum_{i=1}^n w_i x_i$
- ▶  $y$  is an **output** given by  $y = \sigma(\xi)$   
where  $\sigma$  is an **activation function**;  
e.g. a *unit step function*

$$\sigma(\xi) = \begin{cases} 1 & \xi \geq 0; \\ 0 & \xi < 0. \end{cases}$$

(The threshold  $h$  has been substituted with the new input  $x_0 = 1$  and the weight  $w_0 = -h$ .)

# Neuron and linear separation



- ▶ inner potential

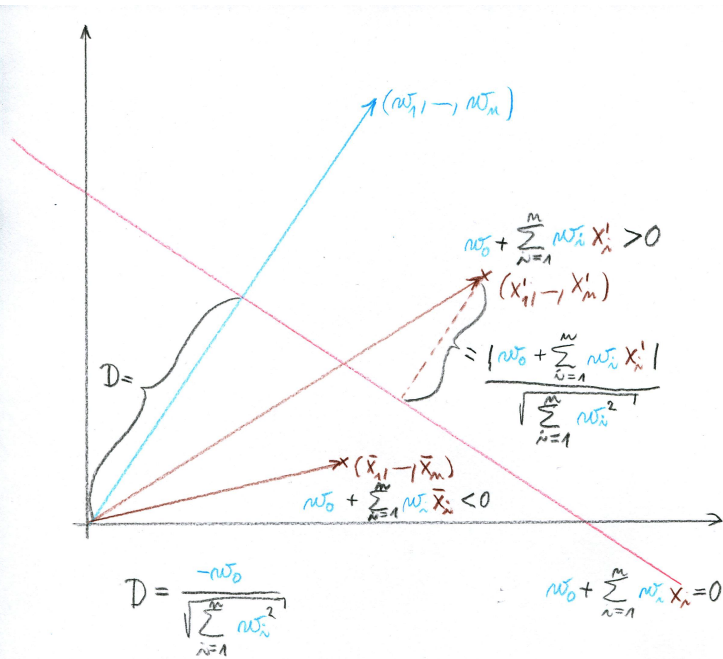
$$\xi = w_0 + \sum_{i=1}^n w_i x_i$$

determines a separation hyperplane in the  $n$ -dimensional **input space**

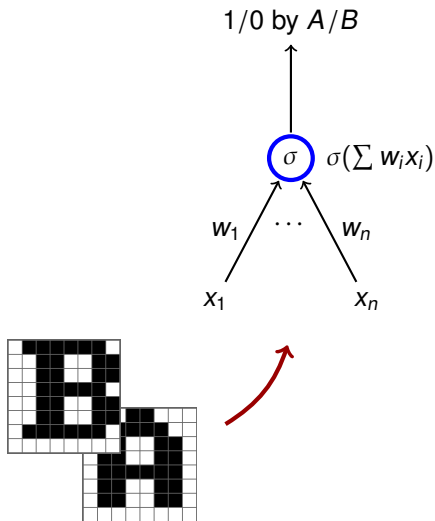
- ▶ in 2d line
- ▶ in 3d plane
- ▶ ...



# Neuron geometry

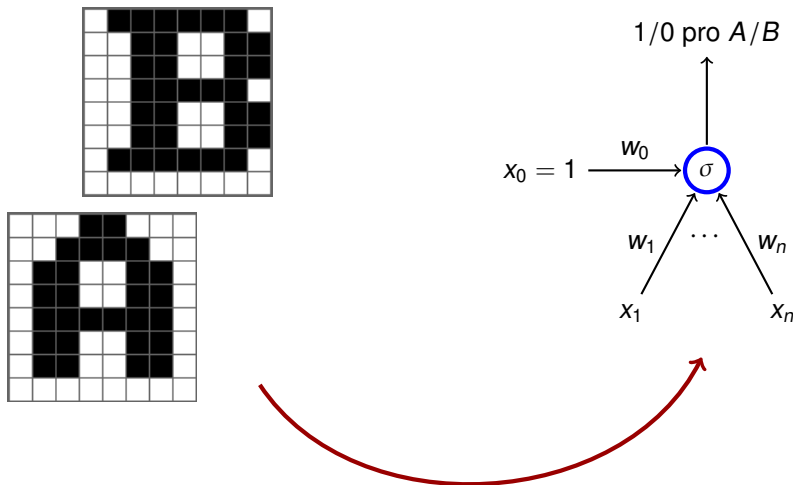


# Neuron and linear separation



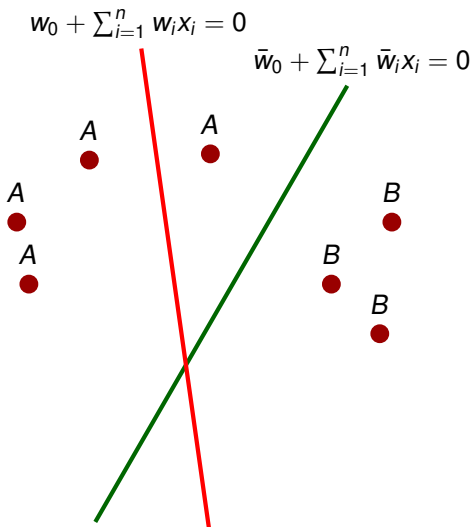
$n = 8 \cdot 8$ , i.e. the number of pixels in the images. Inputs are binary vectors of dimension  $n$  (black pixel  $\approx 1$ , white pixel  $\approx 0$ ).

# Neuron and linear separation



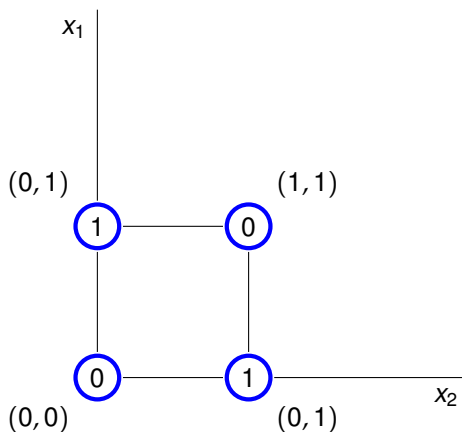
$n = 8 \cdot 8$ , i.e. the number of pixels in the images. Inputs are binary vectors of dimension  $n$  (black pixel  $\approx 1$ , white pixel  $\approx 0$ ).

# Neuron and linear separation



- ▶ Red line classifies incorrectly
- ▶ Green line classifies correctly (may be a result of a correction by a learning algorithm)

# Neuron and linear separation (XOR)



- No line separates ones from zeros.

**Neural network** consists of formal neurons interconnected in such a way that the output of one neuron is an input of several other neurons.

In order to describe a particular type of neural networks we need to specify:

- ▶ **Architecture**  
How the neurons are connected.
- ▶ **Activity**  
How the network transforms inputs to outputs.
- ▶ **Learning**  
How the weights are changed during training.

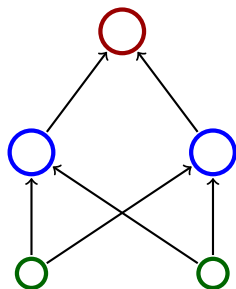
# Architecture

**Network architecture** is given as a digraph whose nodes are neurons and edges are connections.

We distinguish several categories of neurons:

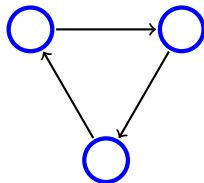
- ▶ **Output neurons**
- ▶ **Hidden neurons**
- ▶ **Input neurons**

(In general, a neuron may be both input and output; a neuron is hidden if it is neither input, nor output.)



# Architecture – Cycles

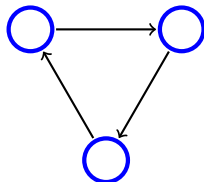
- ▶ A network is **cyclic** (recurrent) if its architecture contains a directed cycle.



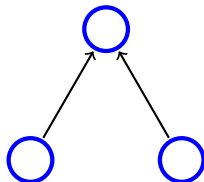


# Architecture – Cycles

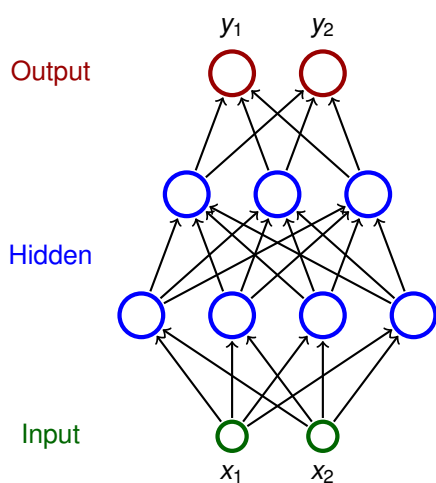
- ▶ A network is **cyclic** (recurrent) if its architecture contains a directed cycle.



- ▶ Otherwise it is **acyclic** (feed-forward)

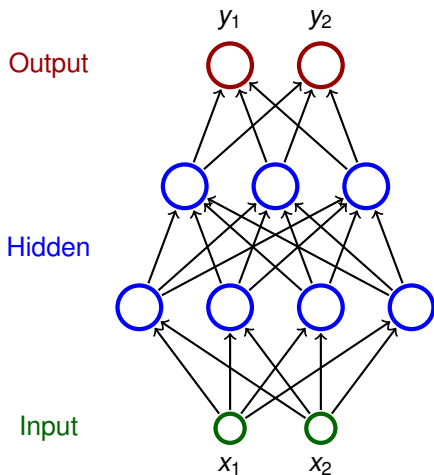


# Architecture – Multilayer Perceptron (MLP)



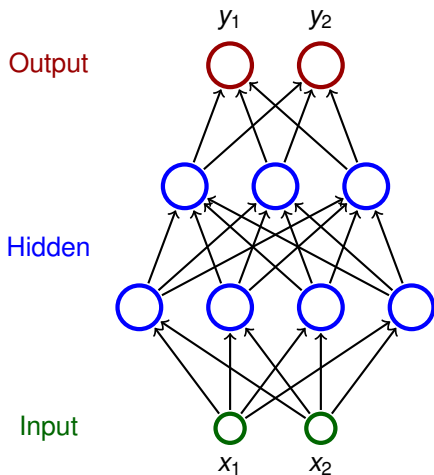
- Neurons partitioned into **layers**; one input layer, one output layer, possibly several hidden layers

# Architecture – Multilayer Perceptron (MLP)



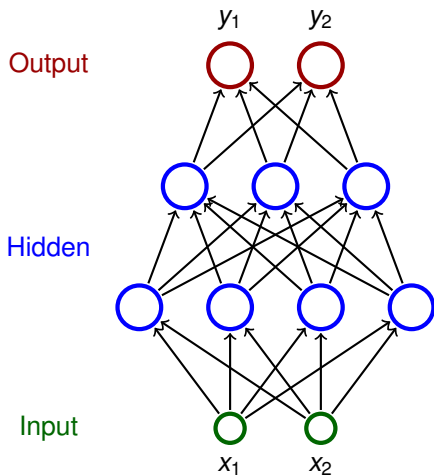
- ▶ Neurons partitioned into **layers**; one input layer, one output layer, possibly several hidden layers
- ▶ layers numbered from 0; the input layer has number 0
  - ▶ E.g. three-layer network has two hidden layers and one output layer

# Architecture – Multilayer Perceptron (MLP)



- ▶ Neurons partitioned into **layers**; one input layer, one output layer, possibly several hidden layers
- ▶ layers numbered from 0; the input layer has number 0
  - ▶ E.g. three-layer network has two hidden layers and one output layer
- ▶ Neurons in the  $i$ -th layer are connected with all neurons in the  $i + 1$ -st layer

# Architecture – Multilayer Perceptron (MLP)



- ▶ Neurons partitioned into **layers**; one input layer, one output layer, possibly several hidden layers
- ▶ layers numbered from 0; the input layer has number 0
  - ▶ E.g. three-layer network has two hidden layers and one output layer
- ▶ Neurons in the  $i$ -th layer are connected with all neurons in the  $i + 1$ -st layer
- ▶ Architecture of a MLP is typically described by numbers of neurons in individual layers (e.g. 2-4-3-2)

## Activity

Consider a network with  $n$  neurons,  $k$  input and  $\ell$  output.

## Activity

Consider a network with  $n$  neurons,  $k$  input and  $\ell$  output.

- ▶ **State** of a network is a vector of output values of all neurons.

(States of a network with  $n$  neurons are vectors of  $\mathbb{R}^n$ )

- ▶ **State-space** of a network is a set of all states.

# Activity

Consider a network with  $n$  neurons,  $k$  input and  $\ell$  output.

- ▶ **State** of a network is a vector of output values of all neurons.

(States of a network with  $n$  neurons are vectors of  $\mathbb{R}^n$ )

- ▶ **State-space** of a network is a set of all states.

- ▶ **Network input** is a vector of  $k$  real numbers, i.e. an element of  $\mathbb{R}^k$ .

- ▶ **Network input space** is a set of all network inputs.  
(sometimes we restrict ourselves to a proper subset of  $\mathbb{R}^k$ )



# Activity

Consider a network with  $n$  neurons,  $k$  input and  $\ell$  output.

- ▶ **State** of a network is a vector of output values of all neurons.

(States of a network with  $n$  neurons are vectors of  $\mathbb{R}^n$ )

- ▶ **State-space** of a network is a set of all states.

- ▶ **Network input** is a vector of  $k$  real numbers, i.e. an element of  $\mathbb{R}^k$ .

- ▶ **Network input space** is a set of all network inputs.  
(sometimes we restrict ourselves to a proper subset of  $\mathbb{R}^k$ )

- ▶ **Initial state**

Input neurons set to values from the network input  
(each component of the network input corresponds to an input neuron)

Values of the remaining neurons set to 0.

## Activity – computation of a network

- ▶ **Computation** (typically) proceeds in discrete steps.

## Activity – computation of a network

- ▶ **Computation** (typically) proceeds in discrete steps.  
In every step the following happens:

## Activity – computation of a network

- ▶ **Computation** (typically) proceeds in discrete steps.  
In every step the following happens:
  1. A set of neurons is selected according to some rule.
  2. The selected neurons change their states according to their inputs (they are simply evaluated).  
(If a neuron does not have any inputs, its value remains constant.)

## Activity – computation of a network

- ▶ **Computation** (typically) proceeds in discrete steps. In every step the following happens:
  1. A set of neurons is selected according to some rule.
  2. The selected neurons change their states according to their inputs (they are simply evaluated).  
(If a neuron does not have any inputs, its value remains constant.)

A computation is **finite** on a network input  $\vec{x}$  if the state changes only finitely many times (i.e. there is a moment in time after which the state of the network never changes). We also say that the network **stops on**  $\vec{x}$ .

## Activity – computation of a network

- ▶ **Computation** (typically) proceeds in discrete steps. In every step the following happens:
  1. A set of neurons is selected according to some rule.
  2. The selected neurons change their states according to their inputs (they are simply evaluated).

(If a neuron does not have any inputs, its value remains constant.)

A computation is **finite** on a network input  $\vec{x}$  if the state changes only finitely many times (i.e. there is a moment in time after which the state of the network never changes).

We also say that the network **stops on**  $\vec{x}$ .

- ▶ **Network output** is a vector of values of all output neurons in the network (i.e., an element of  $\mathbb{R}^\ell$ ).

Note that the network output keeps changing throughout the computation!

## Activity – computation of a network

- ▶ **Computation** (typically) proceeds in discrete steps. In every step the following happens:
  1. A set of neurons is selected according to some rule.
  2. The selected neurons change their states according to their inputs (they are simply evaluated).

(If a neuron does not have any inputs, its value remains constant.)

A computation is **finite** on a network input  $\vec{x}$  if the state changes only finitely many times (i.e. there is a moment in time after which the state of the network never changes).

We also say that the network **stops on**  $\vec{x}$ .

- ▶ **Network output** is a vector of values of all output neurons in the network (i.e., an element of  $\mathbb{R}^\ell$ ).

Note that the network output keeps changing throughout the computation!

*MLP* uses the following selection rule:

In the  $i$ -th step evaluate all neurons in the  $i$ -th layer.

# Activity – semantics of a network

## Definition

*Consider a network with  $n$  neurons,  $k$  input,  $\ell$  output.*

*Let  $A \subseteq \mathbb{R}^k$  and  $B \subseteq \mathbb{R}^\ell$ . Suppose that the network stops on every input of  $A$ .*

*Then we say that the network computes a function  $F : A \rightarrow B$  if for every network input  $\vec{x}$  the vector  $F(\vec{x}) \in B$  is the output of the network after the computation on  $\vec{x}$  stops.*



## Activity – semantics of a network

### Definition

*Consider a network with  $n$  neurons,  $k$  input,  $\ell$  output.*

*Let  $A \subseteq \mathbb{R}^k$  and  $B \subseteq \mathbb{R}^\ell$ . Suppose that the network stops on every input of  $A$ .*

*Then we say that the network computes a function  $F : A \rightarrow B$  if for every network input  $\vec{x}$  the vector  $F(\vec{x}) \in B$  is the output of the network after the computation on  $\vec{x}$  stops.*

# Activity – semantics of a network

## Definition

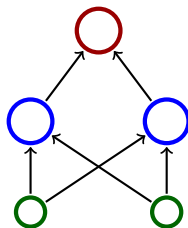
Consider a network with  $n$  neurons,  $k$  input,  $\ell$  output.

Let  $A \subseteq \mathbb{R}^k$  and  $B \subseteq \mathbb{R}^\ell$ . Suppose that the network stops on every input of  $A$ .

Then we say that the network computes a function  $F : A \rightarrow B$  if for every network input  $\vec{x}$  the vector  $F(\vec{x}) \in B$  is the output of the network after the computation on  $\vec{x}$  stops.

## Example 1

This network computes a function from  $\mathbb{R}^2$  to  $\mathbb{R}$ .



## Activity – inner potential and activation functions

In order to specify activity of the network, we need to specify how the inner potentials  $\xi$  are computed and what are the activation functions  $\sigma$ .

## Activity – inner potential and activation functions

In order to specify activity of the network, we need to specify how the inner potentials  $\xi$  are computed and what are the activation functions  $\sigma$ .

We assume (unless otherwise specified) that

$$\xi = w_0 + \sum_{i=1}^n w_i \cdot x_i$$

here  $\vec{x} = (x_1, \dots, x_n)$  are inputs of the neuron and  $\vec{w} = (w_1, \dots, w_n)$  are weights.

## Activity – inner potential and activation functions

In order to specify activity of the network, we need to specify how the inner potentials  $\xi$  are computed and what are the activation functions  $\sigma$ .

We assume (unless otherwise specified) that

$$\xi = w_0 + \sum_{i=1}^n w_i \cdot x_i$$

here  $\vec{x} = (x_1, \dots, x_n)$  are inputs of the neuron and  $\vec{w} = (w_1, \dots, w_n)$  are weights.

There are special types of neural networks where the inner potential is computed differently, e.g., as a "distance" of an input from the weight vector:

$$\xi = \|\vec{x} - \vec{w}\|$$

here  $\|\cdot\|$  is a vector norm, typically Euclidean.

## Activity – inner potential and activation functions

There are many activation functions, typical examples:

- ▶ Unit step function

$$\sigma(\xi) = \begin{cases} 1 & \xi \geq 0; \\ 0 & \xi < 0. \end{cases}$$

# Activity – inner potential and activation functions

There are many activation functions, typical examples:

- ▶ Unit step function

$$\sigma(\xi) = \begin{cases} 1 & \xi \geq 0; \\ 0 & \xi < 0. \end{cases}$$

- ▶ (Logistic) sigmoid

$$\sigma(\xi) = \frac{1}{1 + e^{-\lambda \cdot \xi}} \quad \text{here } \lambda \in \mathbb{R} \text{ is a } \textit{steepness} \text{ parameter.}$$

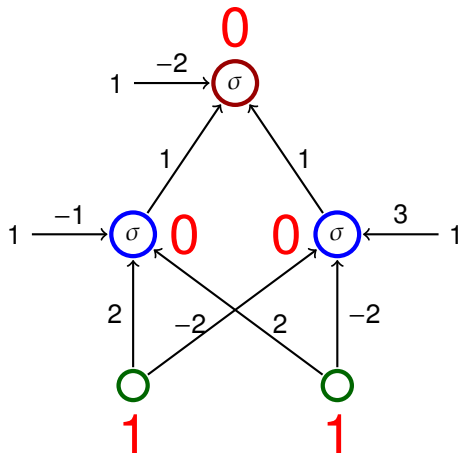
- ▶ Hyperbolic tangens

$$\sigma(\xi) = \frac{1 - e^{-\xi}}{1 + e^{-\xi}}$$

- ▶ ReLU

$$\sigma(\xi) = \max(\xi, 0)$$

# Activity – XOR



- Activation function is a unit step function

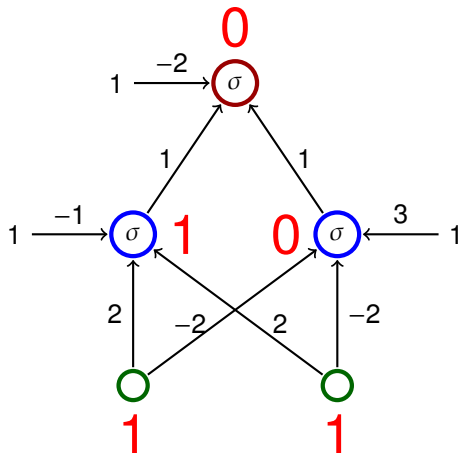
$$\sigma(\xi) = \begin{cases} 1 & \xi \geq 0; \\ 0 & \xi < 0. \end{cases}$$

- The network computes  $XOR(x_1, x_2)$

$x_1$	$x_2$	$y$
1	1	0
1	0	1
0	1	1
0	0	0



# Activity – XOR



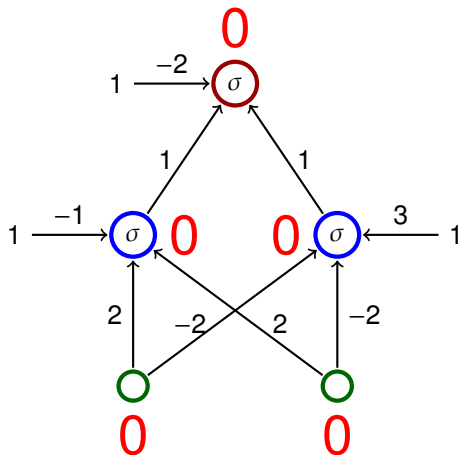
- Activation function is a unit step function

$$\sigma(\xi) = \begin{cases} 1 & \xi \geq 0; \\ 0 & \xi < 0. \end{cases}$$

- The network computes  $XOR(x_1, x_2)$

$x_1$	$x_2$	$y$
1	1	0
1	0	1
0	1	1
0	0	0

# Activity – XOR



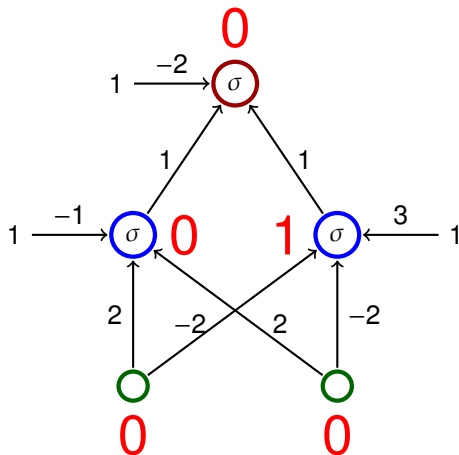
- Activation function is a unit step function

$$\sigma(\xi) = \begin{cases} 1 & \xi \geq 0; \\ 0 & \xi < 0. \end{cases}$$

- The network computes  $XOR(x_1, x_2)$

$x_1$	$x_2$	$y$
1	1	0
1	0	1
0	1	1
0	0	0

# Activity – XOR



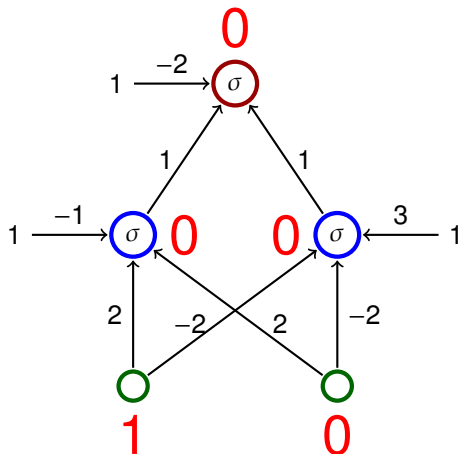
- Activation function is a unit step function

$$\sigma(\xi) = \begin{cases} 1 & \xi \geq 0; \\ 0 & \xi < 0. \end{cases}$$

- The network computes  $XOR(x_1, x_2)$

$x_1$	$x_2$	$y$
1	1	0
1	0	1
0	1	1
0	0	0

# Activity – XOR



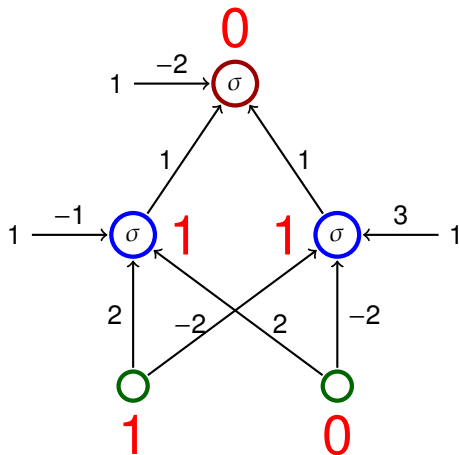
- Activation function is a unit step function

$$\sigma(\xi) = \begin{cases} 1 & \xi \geq 0; \\ 0 & \xi < 0. \end{cases}$$

- The network computes  $XOR(x_1, x_2)$

$x_1$	$x_2$	$y$
1	1	0
1	0	1
0	1	1
0	0	0

# Activity – XOR



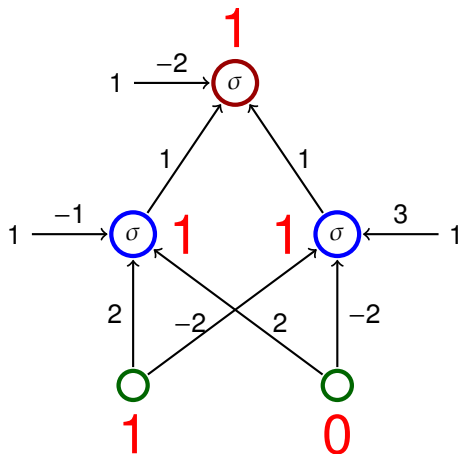
- Activation function is a unit step function

$$\sigma(\xi) = \begin{cases} 1 & \xi \geq 0; \\ 0 & \xi < 0. \end{cases}$$

- The network computes  $XOR(x_1, x_2)$

$x_1$	$x_2$	$y$
1	1	0
1	0	1
0	1	1
0	0	0

# Activity – XOR



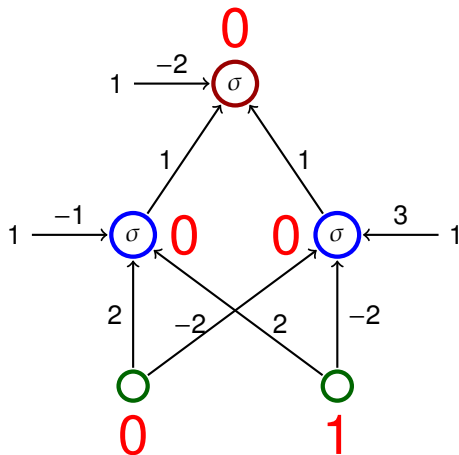
- Activation function is a unit step function

$$\sigma(\xi) = \begin{cases} 1 & \xi \geq 0; \\ 0 & \xi < 0. \end{cases}$$

- The network computes  $XOR(x_1, x_2)$

$x_1$	$x_2$	$y$
1	1	0
1	0	1
0	1	1
0	0	0

# Activity – XOR



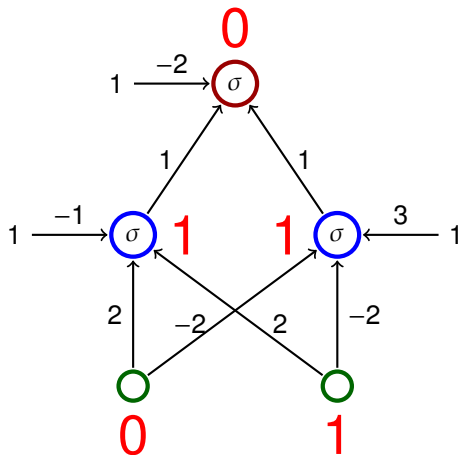
- Activation function is a unit step function

$$\sigma(\xi) = \begin{cases} 1 & \xi \geq 0; \\ 0 & \xi < 0. \end{cases}$$

- The network computes  $XOR(x_1, x_2)$

$x_1$	$x_2$	$y$
1	1	0
1	0	1
0	1	1
0	0	0

# Activity – XOR



- Activation function is a unit step function

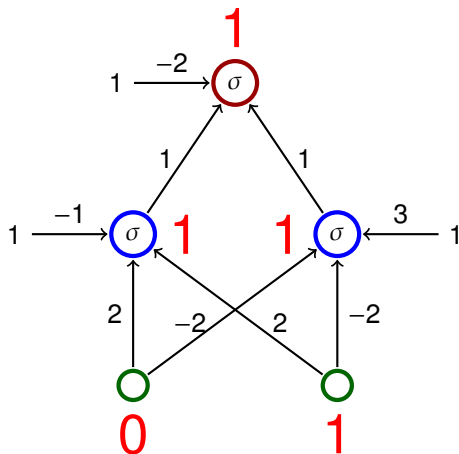
$$\sigma(\xi) = \begin{cases} 1 & \xi \geq 0; \\ 0 & \xi < 0. \end{cases}$$

- The network computes  $XOR(x_1, x_2)$

$x_1$	$x_2$	$y$
1	1	0
1	0	1
0	1	1
0	0	0



# Activity – XOR



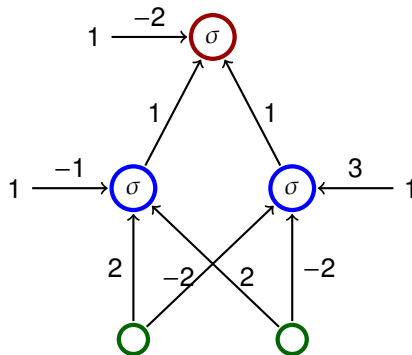
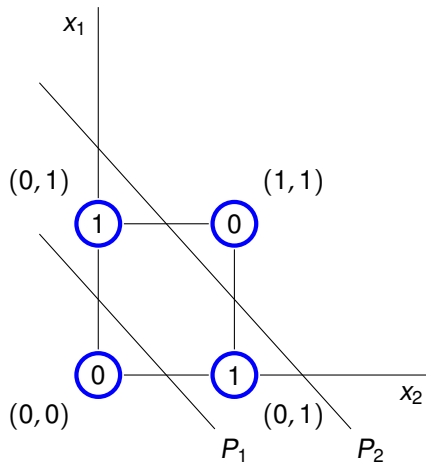
- Activation function is a unit step function

$$\sigma(\xi) = \begin{cases} 1 & \xi \geq 0; \\ 0 & \xi < 0. \end{cases}$$

- The network computes  $XOR(x_1, x_2)$

$x_1$	$x_2$	$y$
1	1	0
1	0	1
0	1	1
0	0	0

# Activity – MLP and linear separation



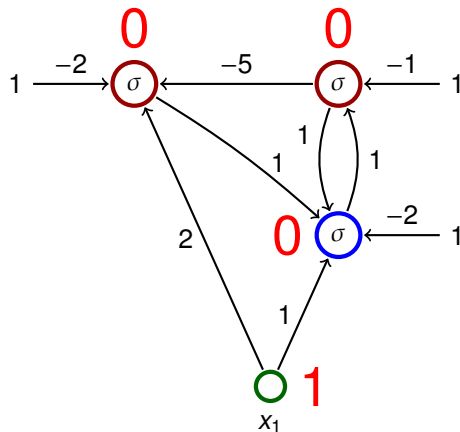
- ▶ The line  $P_1$  is given by  $-1 + 2x_1 + 2x_2 = 0$
- ▶ The line  $P_2$  is given by  $3 - 2x_1 - 2x_2 = 0$

## Activity – example

The activation function is the unit step function

$$\sigma(\xi) = \begin{cases} 1 & \xi \geq 0; \\ 0 & \xi < 0. \end{cases}$$

The input is equal to 1

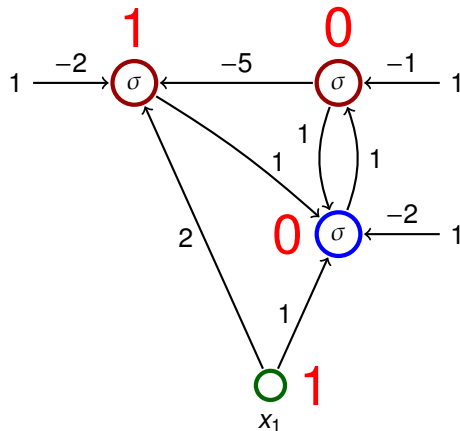


## Activity – example

The activation function is the unit step function

$$\sigma(\xi) = \begin{cases} 1 & \xi \geq 0; \\ 0 & \xi < 0. \end{cases}$$

The input is equal to 1

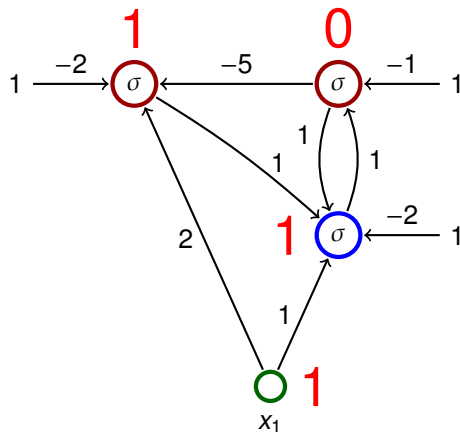


## Activity – example

The activation function is the unit step function

$$\sigma(\xi) = \begin{cases} 1 & \xi \geq 0; \\ 0 & \xi < 0. \end{cases}$$

The input is equal to 1

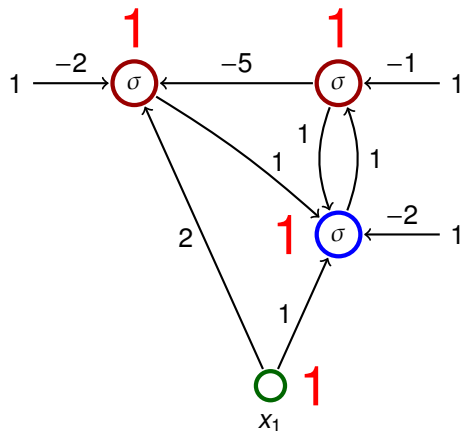


## Activity – example

The activation function is the unit step function

$$\sigma(\xi) = \begin{cases} 1 & \xi \geq 0; \\ 0 & \xi < 0. \end{cases}$$

The input is equal to 1

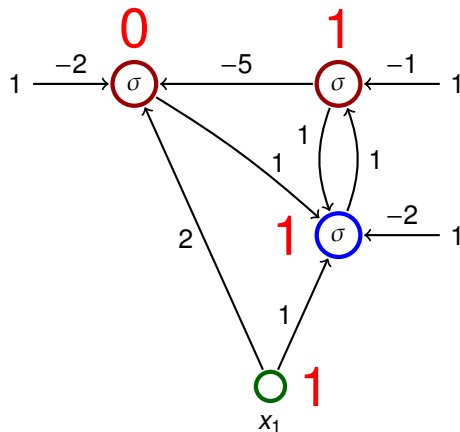


## Activity – example

The activation function is the unit step function

$$\sigma(\xi) = \begin{cases} 1 & \xi \geq 0; \\ 0 & \xi < 0. \end{cases}$$

The input is equal to 1



Consider a network with  $n$  neurons,  $k$  input and  $\ell$  output.



Consider a network with  $n$  neurons,  $k$  input and  $\ell$  output.

- ▶ **Configuration** of a network is a vector of all values of weights.  
(Configurations of a network with  $m$  connections are elements of  $\mathbb{R}^m$ )
- ▶ **Weight-space** of a network is a set of all configurations.

Consider a network with  $n$  neurons,  $k$  input and  $\ell$  output.

- ▶ **Configuration** of a network is a vector of all values of weights.  
(Configurations of a network with  $m$  connections are elements of  $\mathbb{R}^m$ )
- ▶ **Weight-space** of a network is a set of all configurations.
- ▶ **initial configuration**  
weights can be initialized randomly or using some sophisticated algorithm

**Learning rule** for weight adaptation.

(the goal is to find a configuration in which the network computes a desired function)

**Learning rule** for weight adaptation.

(the goal is to find a configuration in which the network computes a desired function)

- ▶ Supervised learning
  - ▶ The desired function is described using *training examples* that are pairs of the form (input, output).
  - ▶ Learning algorithm searches for a configuration which "corresponds" to the training examples, typically by minimizing an error function.

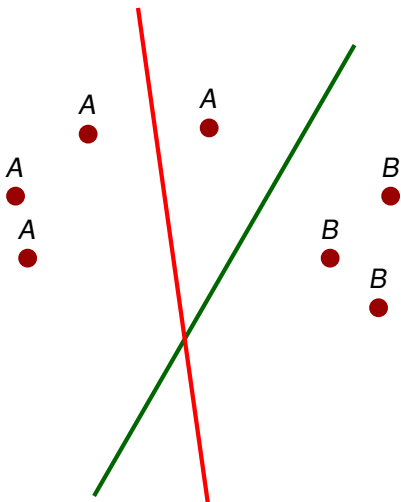
## **Learning rule** for weight adaptation.

(the goal is to find a configuration in which the network computes a desired function)

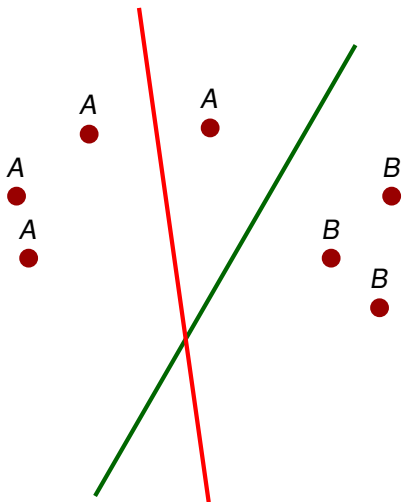
- ▶ Supervised learning
  - ▶ The desired function is described using *training examples* that are pairs of the form (input, output).
  - ▶ Learning algorithm searches for a configuration which "corresponds" to the training examples, typically by minimizing an error function.
- ▶ Unsupervised learning
  - ▶ The training set contains only inputs.
  - ▶ The goal is to determine distribution of the inputs (clustering, deep belief networks, etc.)

# Supervised learning – illustration

- classification in the plane using a single neuron

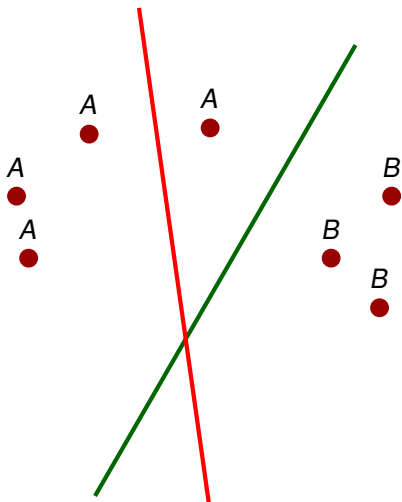


# Supervised learning – illustration



- ▶ classification in the plane using a single neuron
- ▶ training examples are of the form (point, value) where the value is either 1, or 0 depending on whether the point is either *A*, or *B*

# Supervised learning – illustration



- ▶ classification in the plane using a single neuron
- ▶ training examples are of the form (point, value) where the value is either 1, or 0 depending on whether the point is either A, or B
- ▶ the algorithm considers examples one after another
- ▶ whenever an incorrectly classified point is considered, the learning algorithm turns the line in the direction of the point



# Summary – Advantages of neural networks

- ▶ Massive parallelism
  - ▶ neurons can be evaluated in parallel

# Summary – Advantages of neural networks

- ▶ Massive parallelism
  - ▶ neurons can be evaluated in parallel
- ▶ Learning
  - ▶ many sophisticated learning algorithms used to "program" neural networks

# Summary – Advantages of neural networks

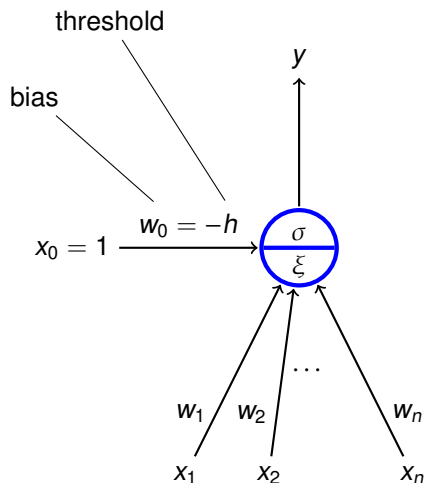
- ▶ Massive parallelism
  - ▶ neurons can be evaluated in parallel
- ▶ Learning
  - ▶ many sophisticated learning algorithms used to "program" neural networks
- ▶ generalization and robustness
  - ▶ information is encoded in a distributed manner in weights
  - ▶ "close" inputs typically get similar values

# Summary – Advantages of neural networks

- ▶ Massive parallelism
  - ▶ neurons can be evaluated in parallel
- ▶ Learning
  - ▶ many sophisticated learning algorithms used to "program" neural networks
- ▶ generalization and robustness
  - ▶ information is encoded in a distributed manner in weights
  - ▶ "close" inputs typically get similar values
- ▶ Graceful degradation
  - ▶ damage typically causes only a decrease in precision of results

# Expressive power of neural networks

# Formal neuron (with bias)



- ▶  $x_0 = 1, x_1, \dots, x_n \in \mathbb{R}$  are **inputs**
- ▶  $w_0, w_1, \dots, w_n \in \mathbb{R}$  are **weights**
- ▶  $\xi$  is an **inner potential**;  
almost always  $\xi = w_0 + \sum_{i=1}^n w_i x_i$
- ▶  $y$  is an **output** given by  $y = \sigma(\xi)$   
where  $\sigma$  is an **activation function**;  
e.g. a *unit step function*

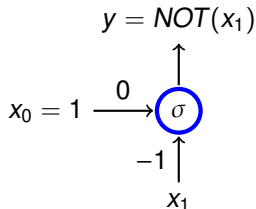
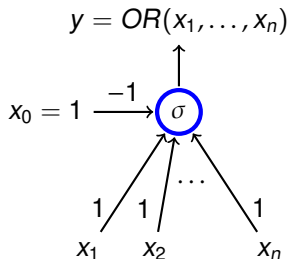
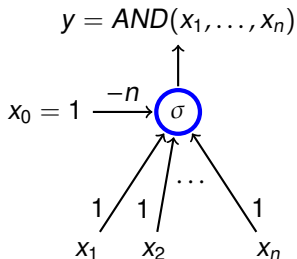
$$\sigma(\xi) = \begin{cases} 1 & \xi \geq 0; \\ 0 & \xi < 0. \end{cases}$$

# Boolean functions

Activation function: *unit step function*  $\sigma(\xi) = \begin{cases} 1 & \xi \geq 0; \\ 0 & \xi < 0. \end{cases}$

# Boolean functions

Activation function: *unit step function*  $\sigma(\xi) = \begin{cases} 1 & \xi \geq 0; \\ 0 & \xi < 0. \end{cases}$





# Boolean functions

## Theorem

*Let  $\sigma$  be the unit step function. Two layer MLPs, where each neuron has  $\sigma$  as the activation function, are able to compute all functions of the form  $F : \{0, 1\}^n \rightarrow \{0, 1\}$ .*

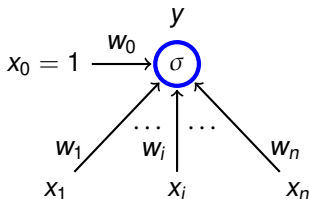
# Boolean functions

## Theorem

Let  $\sigma$  be the unit step function. Two layer MLPs, where each neuron has  $\sigma$  as the activation function, are able to compute all functions of the form  $F : \{0, 1\}^n \rightarrow \{0, 1\}$ .

## Proof.

- ▶ Given a vector  $\vec{v} = (v_1, \dots, v_n) \in \{0, 1\}^n$ , consider a neuron  $N_{\vec{v}}$  whose output is 1 iff the input is  $\vec{v}$ :

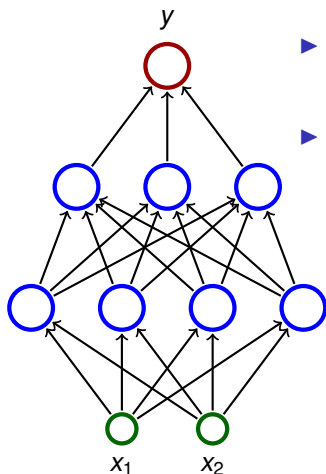


$$w_0 = -\sum_{i=1}^n v_i$$

$$w_i = \begin{cases} 1 & v_i = 1 \\ -1 & v_i = 0 \end{cases}$$

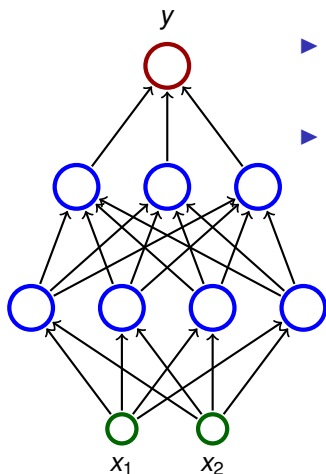
- ▶ Now let us connect all outputs of all neurons  $N_{\vec{v}}$  satisfying  $F(\vec{v}) = 1$  using a neuron implementing OR. □

# Non-linear separation



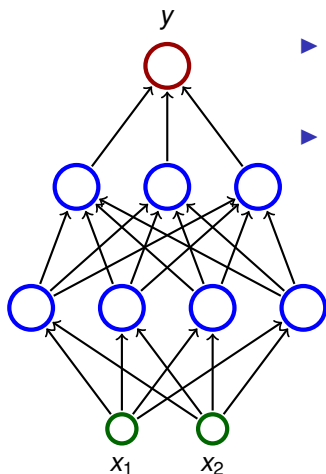
- ▶ Consider a three layer network; each neuron has the unit step activation function.
- ▶ The network divides the input space in two subspaces according to the output (0 or 1).

# Non-linear separation



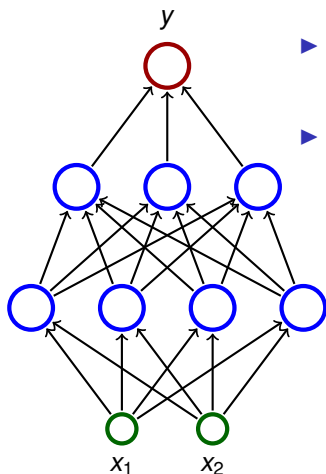
- ▶ Consider a three layer network; each neuron has the unit step activation function.
- ▶ The network divides the input space in two subspaces according to the output (0 or 1).
  - ▶ The first (hidden) layer divides the input space into half-spaces.

# Non-linear separation



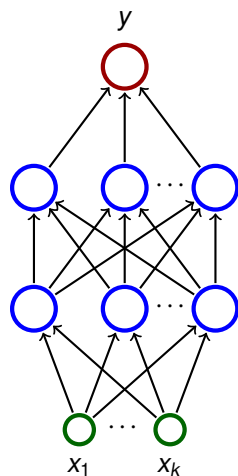
- ▶ Consider a three layer network; each neuron has the unit step activation function.
- ▶ The network divides the input space in two subspaces according to the output (0 or 1).
  - ▶ The first (hidden) layer divides the input space into half-spaces.
  - ▶ The second layer may e.g. make intersections of the half-spaces  $\Rightarrow$  convex sets.

# Non-linear separation



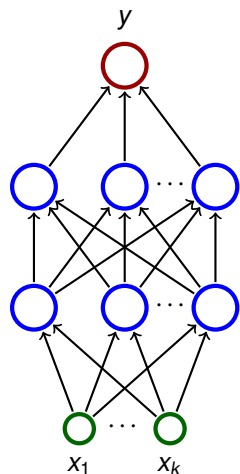
- ▶ Consider a three layer network; each neuron has the unit step activation function.
- ▶ The network divides the input space in two subspaces according to the output (0 or 1).
  - ▶ The first (hidden) layer divides the input space into half-spaces.
  - ▶ The second layer may e.g. make intersections of the half-spaces  $\Rightarrow$  convex sets.
  - ▶ The third layer may e.g. make unions of some convex sets.

# Non-linear separation – illustration



- ▶ Consider three layer networks; each neuron has the unit step activation function.
- ▶ Three layer nets are capable of "approximating" any "reasonable" subset  $A$  of the input space  $\mathbb{R}^k$ .

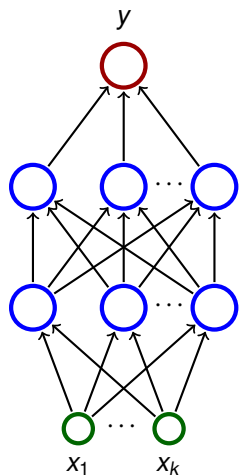
# Non-linear separation – illustration



- ▶ Consider three layer networks; each neuron has the unit step activation function.
- ▶ Three layer nets are capable of "approximating" any "reasonable" subset  $A$  of the input space  $\mathbb{R}^k$ .
  - ▶ Cover  $A$  with hypercubes (in 2D squares, in 3D cubes, ...)

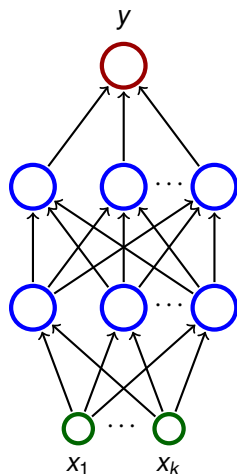


# Non-linear separation – illustration



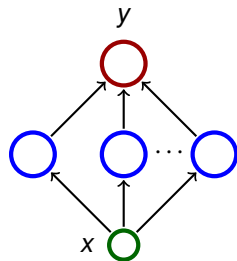
- ▶ Consider three layer networks; each neuron has the unit step activation function.
- ▶ Three layer nets are capable of "approximating" any "reasonable" subset  $A$  of the input space  $\mathbb{R}^k$ .
  - ▶ Cover  $A$  with hypercubes (in 2D squares, in 3D cubes, ...)
  - ▶ Each hypercube  $K$  can be separated using a two layer network  $N_K$  (i.e. a function computed by  $N_K$  gives 1 for points in  $K$  and 0 for the rest).

# Non-linear separation – illustration



- ▶ Consider three layer networks; each neuron has the unit step activation function.
- ▶ Three layer nets are capable of "approximating" any "reasonable" subset  $A$  of the input space  $\mathbb{R}^k$ .
  - ▶ Cover  $A$  with hypercubes (in 2D squares, in 3D cubes, ...)
  - ▶ Each hypercube  $K$  can be separated using a two layer network  $N_K$  (i.e. a function computed by  $N_K$  gives 1 for points in  $K$  and 0 for the rest).
  - ▶ Finally, connect outputs of the nets  $N_K$  satisfying  $K \cap A \neq \emptyset$  using a neuron implementing *OR*.

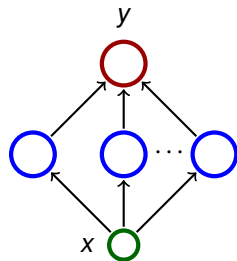
# Power of ReLU



Consider a two layer network

- ▶ with a single input and single output;
- ▶ hidden neurons with the ReLU activation:  
 $\sigma(\xi) = \max(\xi, 0)$ ;
- ▶ the output neuron with identity activation:  
 $\sigma(\xi) = \xi$  (linear model)

# Power of ReLU

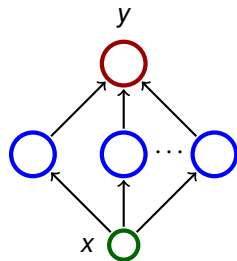


Consider a two layer network

- ▶ with a single input and single output;
- ▶ hidden neurons with the ReLU activation:  
 $\sigma(\xi) = \max(\xi, 0)$ ;
- ▶ the output neuron with identity activation:  
 $\sigma(\xi) = \xi$  (linear model)

For every continuous function  $f : [0, 1] \rightarrow [0, 1]$  and  $\varepsilon > 0$  there is a network of the above type computing a function  $F : [0, 1] \rightarrow \mathbb{R}$  such that  $|f(x) - F(x)| \leq \varepsilon$  for all  $x \in [0, 1]$ .

# Power of ReLU



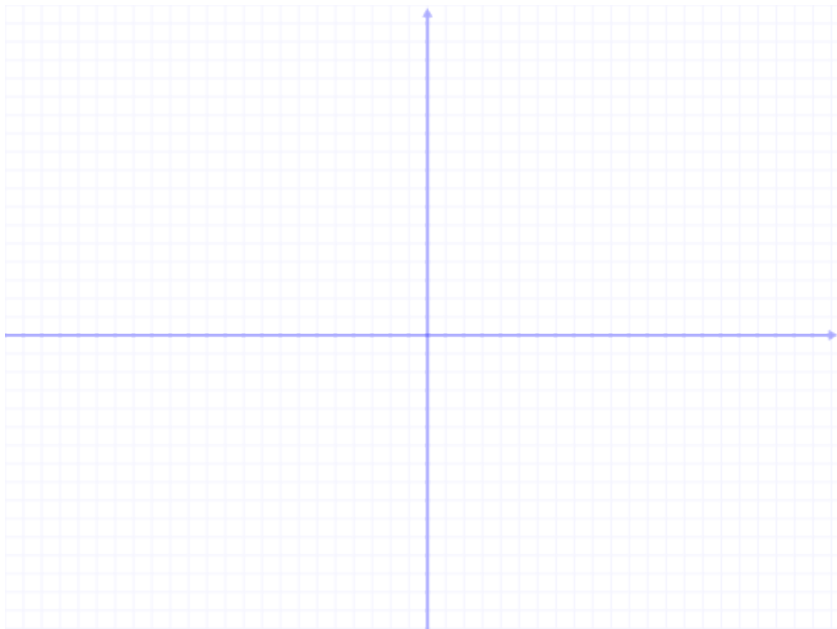
Consider a two layer network

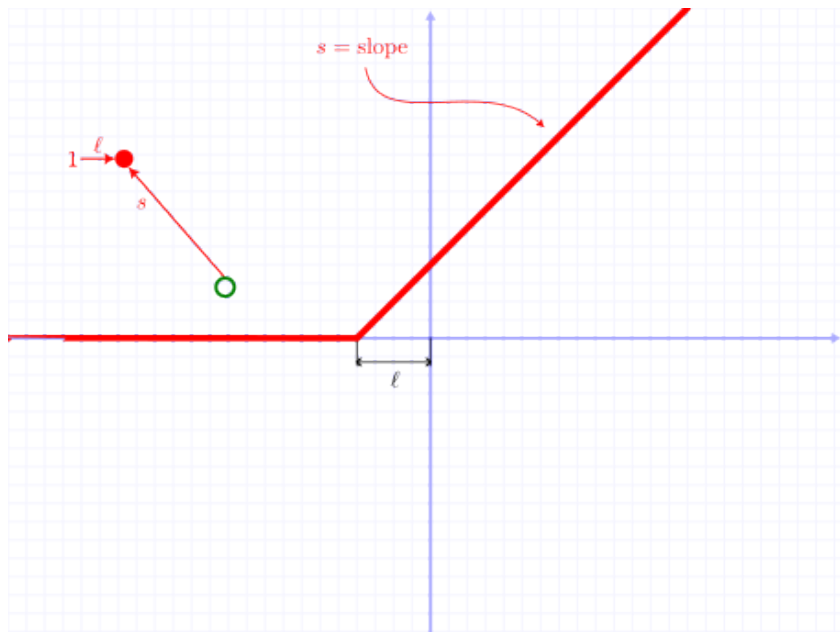
- ▶ with a single input and single output;
- ▶ hidden neurons with the ReLU activation:  
 $\sigma(\xi) = \max(\xi, 0)$ ;
- ▶ the output neuron with identity activation:  
 $\sigma(\xi) = \xi$  (linear model)

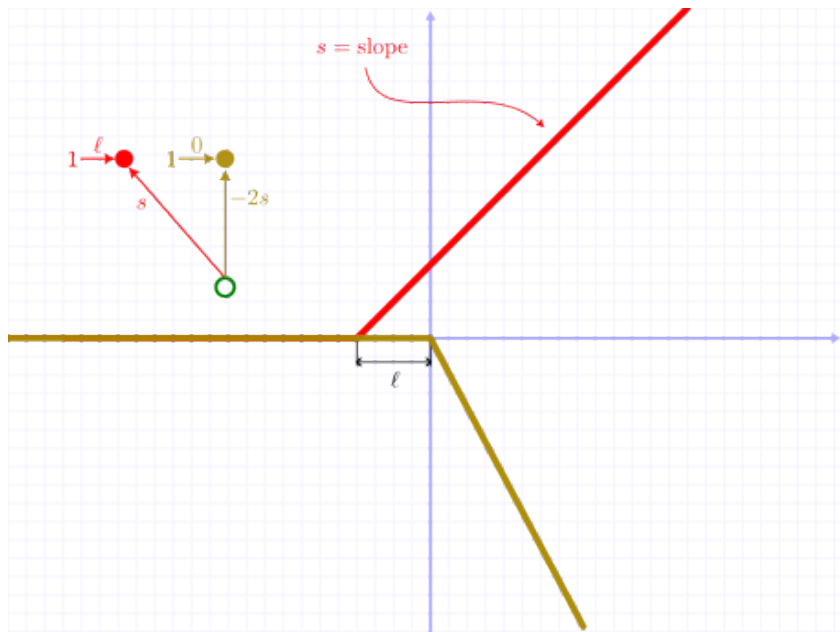
For every continuous function  $f : [0, 1] \rightarrow [0, 1]$  and  $\varepsilon > 0$  there is a network of the above type computing a function  $F : [0, 1] \rightarrow \mathbb{R}$  such that  $|f(x) - F(x)| \leq \varepsilon$  for all  $x \in [0, 1]$ .

For every open subset  $A \subseteq [0, 1]$  there is a network of the above type such that for "most"  $x \in [0, 1]$  we have that  $x \in A$  iff the network's output is  $> 0$  for the input  $x$ .

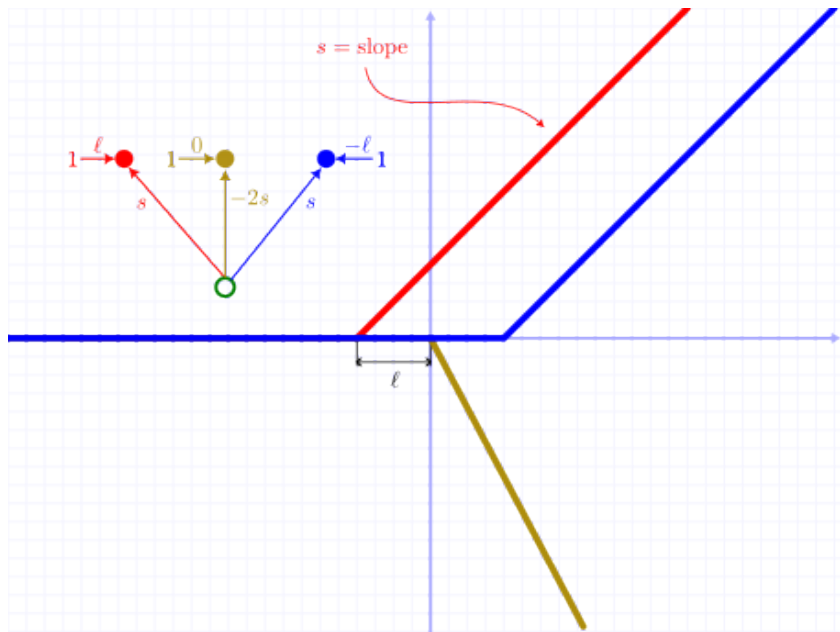
Just consider a continuous function  $f$  where  $f(x)$  is the minimum difference between  $x$  and a point on the boundary of  $A$ . Then uniformly approximate  $f$  using the networks.

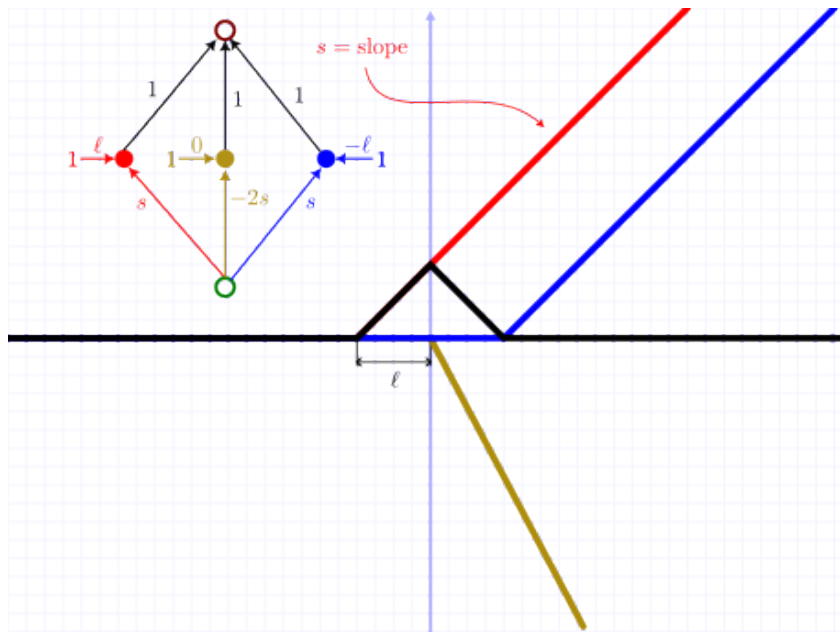


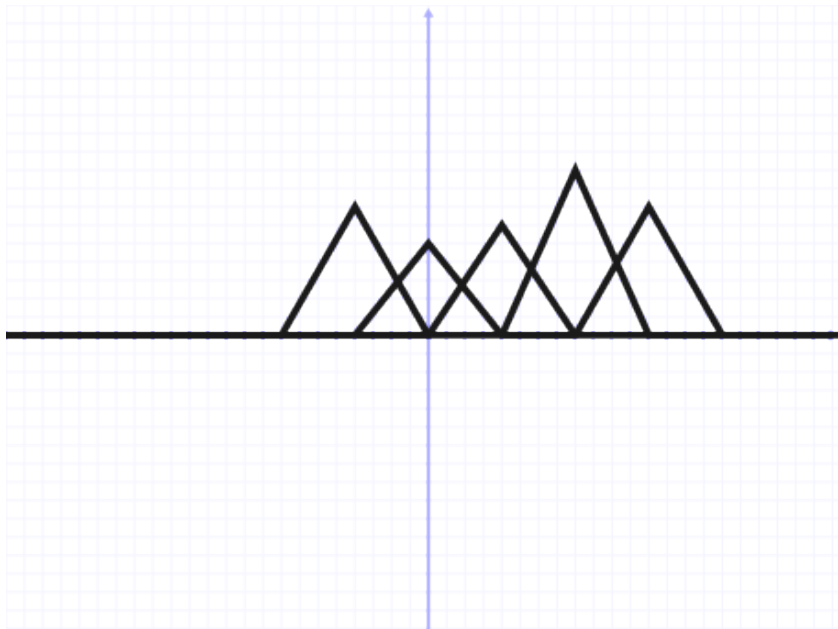


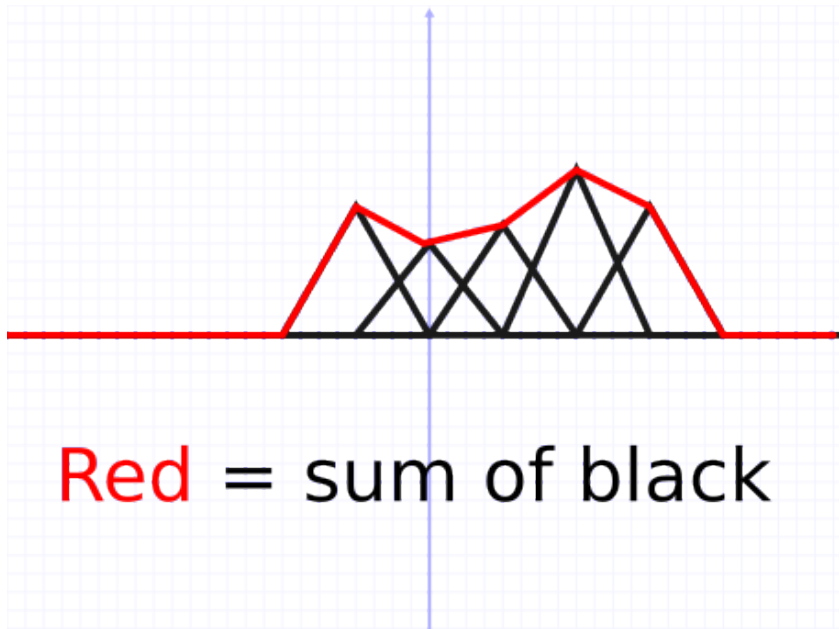












# Non-linear separation - sigmoid

## Theorem (Cybenko 1989 - informal version)

Let  $\sigma$  be a continuous function which is sigmoidal, i.e. satisfies

$$\sigma(x) = \begin{cases} 1 & \text{for } x \rightarrow +\infty \\ 0 & \text{for } x \rightarrow -\infty \end{cases}$$

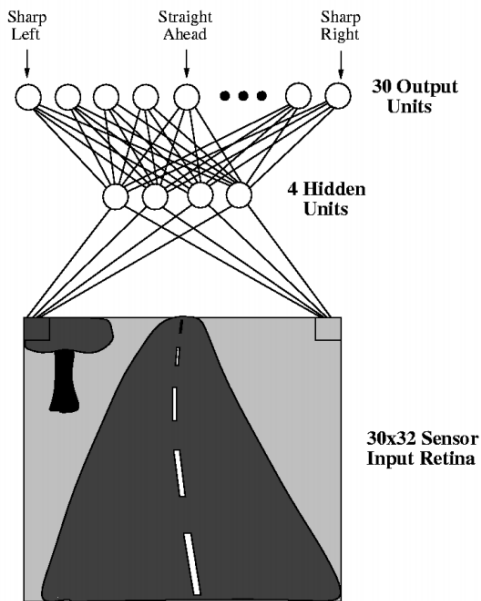
For every "reasonable" set  $A \subseteq [0, 1]^n$ , there is a **two layer network** where each hidden neuron has the activation function  $\sigma$  (output neurons are linear), that satisfies the following:

For "most" vectors  $\vec{v} \in [0, 1]^n$  we have that  $\vec{v} \in A$  iff the network output is  $> 0$  for the input  $\vec{v}$ .

For mathematically oriented:

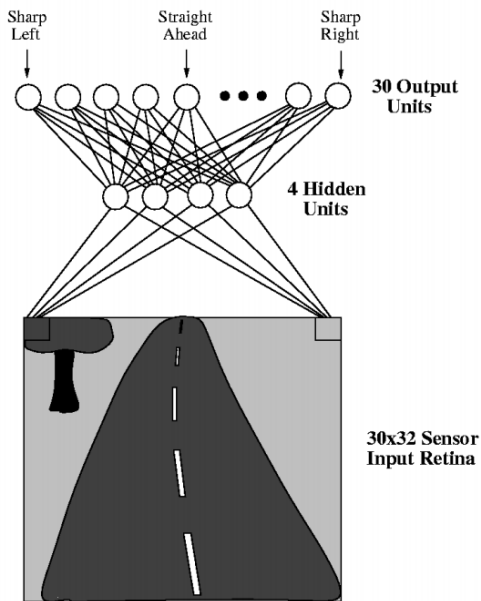
- ▶ "reasonable" means Lebesgue measurable
- ▶ "most" means that the set of incorrectly classified vectors has the Lebesgue measure smaller than a given  $\varepsilon > 0$

# Non-linear separation - practical illustration



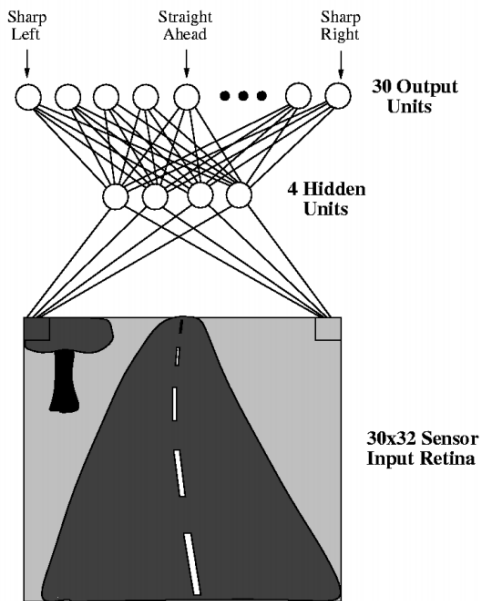
► ALVINN drives a car

# Non-linear separation - practical illustration



- ▶ ALVINN drives a car
- ▶ The net has  $30 \times 32 = 960$  inputs (the input space is thus  $\mathbb{R}^{960}$ )

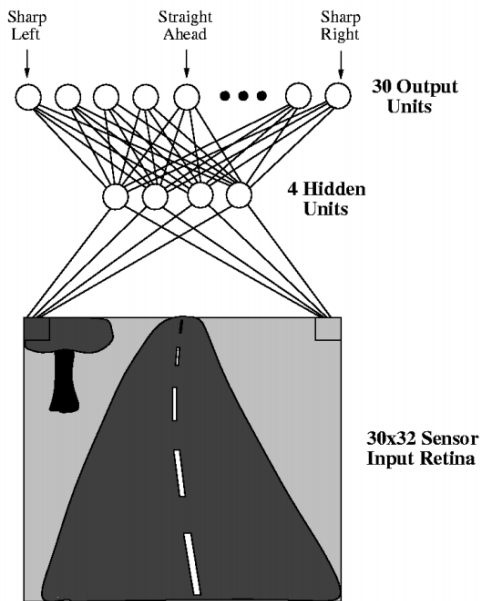
# Non-linear separation - practical illustration



- ▶ ALVINN drives a car
- ▶ The net has  $30 \times 32 = 960$  inputs (the input space is thus  $\mathbb{R}^{960}$ )
- ▶ Input values correspond to shades of gray of pixels.



# Non-linear separation - practical illustration



- ▶ ALVINN drives a car
- ▶ The net has  $30 \times 32 = 960$  inputs (the input space is thus  $\mathbb{R}^{960}$ )
- ▶ Input values correspond to shades of gray of pixels.
- ▶ Output neurons "classify" images of the road based on their "curvature".

# Function approximation - two-layer networks

## Theorem (Cybenko 1989)

*Let  $\sigma$  be a continuous function which is sigmoidal, i.e., is increasing and satisfies*

$$\sigma(x) = \begin{cases} 1 & \text{for } x \rightarrow +\infty \\ 0 & \text{for } x \rightarrow -\infty \end{cases}$$

*For every continuous function  $f : [0, 1]^n \rightarrow [0, 1]$  and every  $\varepsilon > 0$  there is a function  $F : [0, 1]^n \rightarrow [0, 1]$  computed by a **two layer network** where each hidden neuron has the activation function  $\sigma$  (output neurons are linear), that satisfies the following*

$$|f(\vec{v}) - F(\vec{v})| < \varepsilon \quad \text{for every } \vec{v} \in [0, 1]^n.$$

# Neural networks and computability

- ▶ Consider recurrent networks (i.e., containing cycles)

# Neural networks and computability

- ▶ Consider recurrent networks (i.e., containing cycles)
  - ▶ with real weights (in general);

# Neural networks and computability

- ▶ Consider recurrent networks (i.e., containing cycles)
  - ▶ with real weights (in general);
  - ▶ one input neuron and one output neuron (the network computes a function  $F : A \rightarrow \mathbb{R}$  where  $A \subseteq \mathbb{R}$  contains all inputs on which the network stops);

# Neural networks and computability

- ▶ Consider recurrent networks (i.e., containing cycles)
  - ▶ with real weights (in general);
  - ▶ one input neuron and one output neuron (the network computes a function  $F : A \rightarrow \mathbb{R}$  where  $A \subseteq \mathbb{R}$  contains all inputs on which the network stops);
  - ▶ parallel activity rule (output values of all neurons are recomputed in every step);

# Neural networks and computability

- ▶ Consider recurrent networks (i.e., containing cycles)
  - ▶ with real weights (in general);
  - ▶ one input neuron and one output neuron (the network computes a function  $F : A \rightarrow \mathbb{R}$  where  $A \subseteq \mathbb{R}$  contains all inputs on which the network stops);
  - ▶ parallel activity rule (output values of all neurons are recomputed in every step);
  - ▶ activation function

$$\sigma(\xi) = \begin{cases} 1 & \xi \geq 1; \\ \xi & 0 \leq \xi \leq 1; \\ 0 & \xi < 0. \end{cases}$$

# Neural networks and computability

- ▶ Consider recurrent networks (i.e., containing cycles)
  - ▶ with real weights (in general);
  - ▶ one input neuron and one output neuron (the network computes a function  $F : A \rightarrow \mathbb{R}$  where  $A \subseteq \mathbb{R}$  contains all inputs on which the network stops);
  - ▶ parallel activity rule (output values of all neurons are recomputed in every step);
  - ▶ activation function

$$\sigma(\xi) = \begin{cases} 1 & \xi \geq 1; \\ \xi & 0 \leq \xi \leq 1; \\ 0 & \xi < 0. \end{cases}$$

- ▶ We encode words  $\omega \in \{0, 1\}^+$  into numbers as follows:

$$\delta(\omega) = \sum_{i=1}^{|\omega|} \frac{\omega(i)}{2^i} + \frac{1}{2^{|\omega|+1}}$$

E.g.  $\omega = 11001$  gives  $\delta(\omega) = \frac{1}{2} + \frac{1}{2^2} + \frac{1}{2^5} + \frac{1}{2^6}$   
(= 0.110011 in binary form).



# Neural networks and computability

A network **recognizes** a language  $L \subseteq \{0, 1\}^+$  if it computes a function  $F : A \rightarrow \mathbb{R}$  ( $A \subseteq \mathbb{R}$ ) such that

$$\omega \in L \text{ iff } \delta(\omega) \in A \text{ and } F(\delta(\omega)) > 0.$$

# Neural networks and computability

A network **recognizes** a language  $L \subseteq \{0, 1\}^+$  if it computes a function  $F : A \rightarrow \mathbb{R}$  ( $A \subseteq \mathbb{R}$ ) such that

$$\omega \in L \text{ iff } \delta(\omega) \in A \text{ and } F(\delta(\omega)) > 0.$$

- ▶ Recurrent networks with rational weights are equivalent to Turing machines
  - ▶ For every recursively enumerable language  $L \subseteq \{0, 1\}^+$  there is a recurrent network with rational weights and less than 1000 neurons, which recognizes  $L$ .
  - ▶ The halting problem is undecidable for networks with at least 25 neurons and rational weights.
  - ▶ There is "universal" network (equivalent of the universal Turing machine)

# Neural networks and computability

A network **recognizes** a language  $L \subseteq \{0, 1\}^+$  if it computes a function  $F : A \rightarrow \mathbb{R}$  ( $A \subseteq \mathbb{R}$ ) such that

$$\omega \in L \text{ iff } \delta(\omega) \in A \text{ and } F(\delta(\omega)) > 0.$$

- ▶ Recurrent networks with rational weights are equivalent to Turing machines
  - ▶ For every recursively enumerable language  $L \subseteq \{0, 1\}^+$  there is a recurrent network with rational weights and less than 1000 neurons, which recognizes  $L$ .
  - ▶ The halting problem is undecidable for networks with at least 25 neurons and rational weights.
  - ▶ There is "universal" network (equivalent of the universal Turing machine)
- ▶ Recurrent networks are super-Turing powerful

# Neural networks and computability

A network **recognizes** a language  $L \subseteq \{0, 1\}^+$  if it computes a function  $F : A \rightarrow \mathbb{R}$  ( $A \subseteq \mathbb{R}$ ) such that

$$\omega \in L \text{ iff } \delta(\omega) \in A \text{ and } F(\delta(\omega)) > 0.$$

- ▶ Recurrent networks with rational weights are equivalent to Turing machines
  - ▶ For every recursively enumerable language  $L \subseteq \{0, 1\}^+$  there is a recurrent network with rational weights and less than 1000 neurons, which recognizes  $L$ .
  - ▶ The halting problem is undecidable for networks with at least 25 neurons and rational weights.
  - ▶ There is "universal" network (equivalent of the universal Turing machine)
- ▶ Recurrent networks are super-Turing powerful
  - ▶ For **every** language  $L \subseteq \{0, 1\}^+$  there is a recurrent network with less than 1000 neurons which recognizes  $L$ .

# Summary of theoretical results

- ▶ Neural networks are very strong from the point of view of theory:
  - ▶ All Boolean functions can be expressed using two-layer networks.
  - ▶ Two-layer networks may approximate any continuous function.
  - ▶ Recurrent networks are at least as strong as Turing machines.

# Summary of theoretical results

- ▶ Neural networks are very strong from the point of view of theory:
  - ▶ All Boolean functions can be expressed using two-layer networks.
  - ▶ Two-layer networks may approximate any continuous function.
  - ▶ Recurrent networks are at least as strong as Turing machines.
- ▶ These results are purely theoretical!
  - ▶ "Theoretical" networks are extremely huge.
  - ▶ It is very difficult to handcraft them even for simplest problems.
- ▶ From practical point of view, the most important advantages of neural networks are: learning, generalization, robustness.

# Neural networks vs classical computers

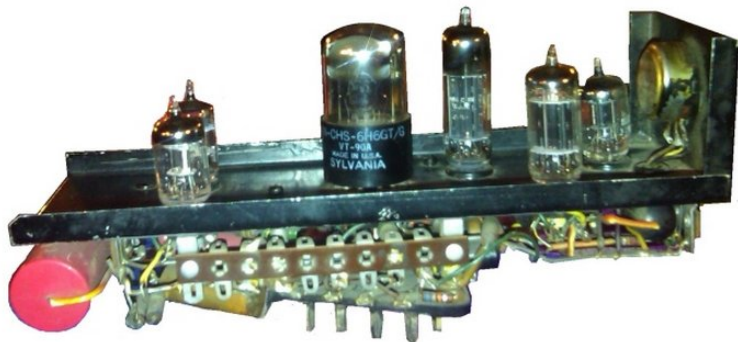
	Neural networks	"Classical" computers
Data	implicitly in weights	explicitly
Computation	naturally parallel	sequential, localized
Robustness	robust w.r.t. input corruption & damage	changing one bit may completely crash the computation
Precision	imprecise, network recalls a training example "similar" to the input	(typically) precise
Programming	learning	manual

## History & implementations



# History of neurocomputers

- ▶ 1951: SNARC (Minski et al)
  - ▶ the first implementation of neural network
  - ▶ a rat strives to exit a maze
  - ▶ 40 artificial neurons (300 vacuum tubes, engines, etc.)



# History of neurocomputers

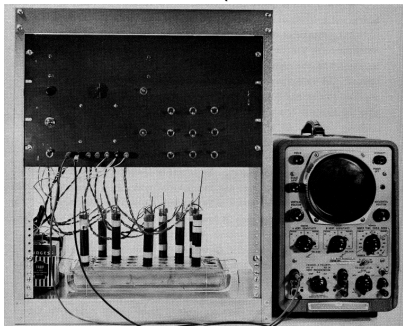
- ▶ 1957: Mark I Perceptron (Rosenblatt et al) - the first successful network for image recognition



- ▶ single layer network
- ▶ image represented by  $20 \times 20$  photocells
- ▶ intensity of pixels was treated as the input to a perceptron (basically the formal neuron), which recognized figures
- ▶ weights were implemented using potentiometers, each set by its own engine
- ▶ it was possible to arbitrarily reconnect inputs to neurons to demonstrate adaptability

# History of neurocomputers

- ▶ 1960: ADALINE (Widrow & Hof)



- ▶ single layer neural network
- ▶ weights stored in a newly invented electronic component **memistor**, which remembers history of electric current in the form of resistance.
- ▶ Widrow founded a company Memistor Corporation, which sold implementations of neural networks.
- ▶ 1960-66: several companies concerned with neural networks were founded.

# History of neurocomputers

- ▶ 1967-82: dead still after publication of a book by Minski & Papert (published 1969, title *Perceptrons*)
- ▶ 1983-end of 90s: revival of neural networks
  - ▶ many attempts at hardware implementations
    - ▶ application specific chips (ASIC)
    - ▶ programmable hardware (FPGA)
  - ▶ hw implementations typically not better than "software" implementations on universal computers (problems with weight storage, size, speed, cost of production etc.)

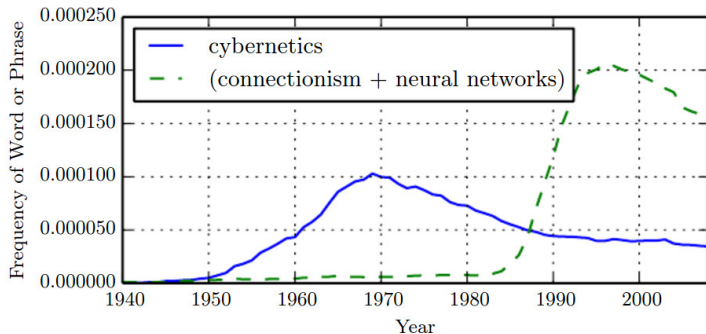
# History of neurocomputers

- ▶ 1967-82: dead still after publication of a book by Minski & Papert (published 1969, title *Perceptrons*)
- ▶ 1983-end of 90s: revival of neural networks
  - ▶ many attempts at hardware implementations
    - ▶ application specific chips (ASIC)
    - ▶ programmable hardware (FPGA)
  - ▶ hw implementations typically not better than "software" implementations on universal computers (problems with weight storage, size, speed, cost of production etc.)
- ▶ end of 90s-cca 2005: NN suppressed by other machine learning methods (support vector machines (SVM))
- ▶ 2006-now: The boom of neural networks!
  - ▶ deep networks – often better than any other method
  - ▶ GPU implementations
  - ▶ ... specialized hw implementations (Google's TPU)

# Some highlights

- ▶ Breakthrough in image recognition.  
Accuracy of image recognition improved by an order of magnitude in 5 years.
- ▶ Breakthrough in game playing.  
Superhuman results in Go and Chess almost without any human intervention. Master level in Starcraft, poker, etc.
- ▶ Breakthrough in machine translation.  
Switching to deep learning produced a 60% increase in translation accuracy compared to the phrase-based approach previously used in Google Translate (in human evaluation)
- ▶ Breakthrough in speech processing.
- ▶ Breakthrough in text generation.  
GPT-4 generates pretty realistic articles, short plays (for a theatre) have been successfully generated, etc.

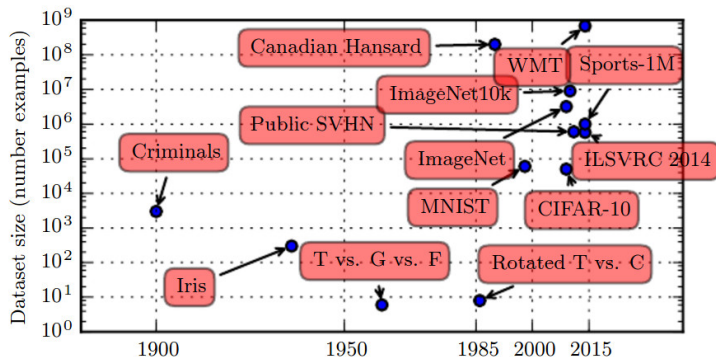
# History in waves ...



**Figure:** The figure shows two of the three historical waves of artificial neural nets research, as measured by the frequency of the phrases "cybernetics" and "connectionism" or "neural networks" according to Google Books (the third wave is too recent to appear).

# Current hardware – What do we face?

Increasing dataset size ...



... weakly-supervised pre-training using hashtags from the Instagram uses  $3.6 \times 10^9$  images.

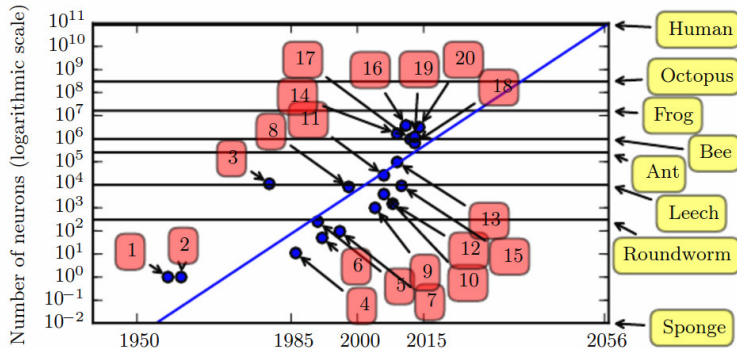
Revisiting Weakly Supervised Pre-Training of Visual Perception Models. Singh et al.

<https://arxiv.org/pdf/2201.08371.pdf>, 2022

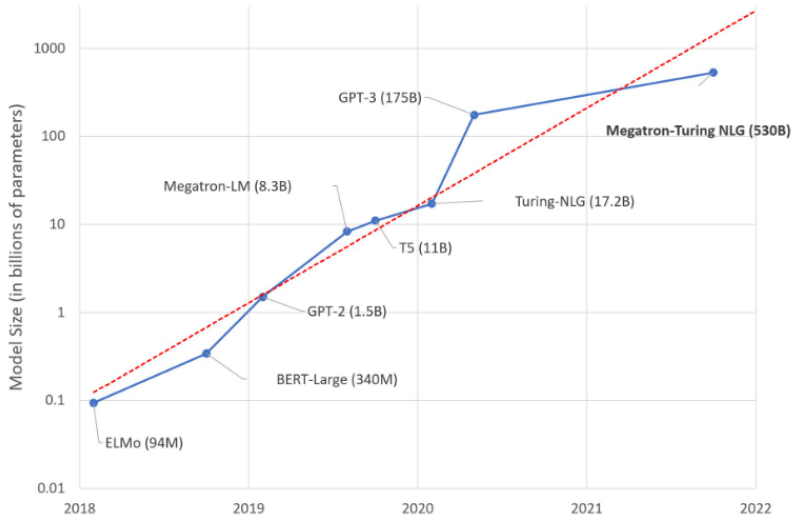


# Current hardware – What do we face?

... and thus increasing size of neural networks ...



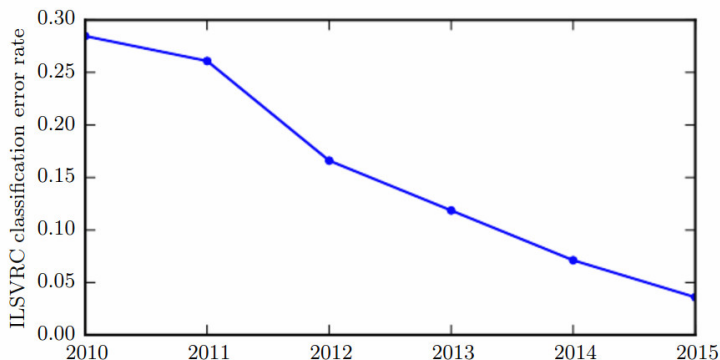
- 2. ADALINE
- 4. Early back-propagation network (Rumelhart et al., 1986b)
- 8. Image recognition: LeNet-5 (LeCun et al., 1998b)
- 10. Dimensionality reduction: Deep belief network (Hinton et al., 2006)  
... here the third "wave" of neural networks started
- 15. Digit recognition: GPU-accelerated multilayer perceptron (Ciresan et al., 2010)
- 18. Image recognition (AlexNet): Multi-GPU convolutional network (Krizhevsky et al., 2012)
- 20. Image recognition: GoogLeNet (Szegedy et al., 2014a)



**GPT-4's Scale:** GPT-4 has 1.8 trillion parameters across 120 layers, which is over 10 times larger than GPT-3.

# Current hardware – What do we face?

... as a reward we get this ...



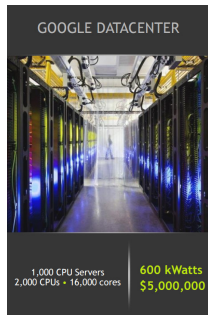
**Figure:** Since deep networks reached the scale necessary to compete in the ImageNetLarge Scale Visual Recognition Challenge, they have consistently won the competition every year, and yielded lower and lower error rates each time. Data from Russakovsky et al. (2014b) and He et al. (2015).

# Current hardware

In 2012, Google trained a large network of 1.7 billion weights and 9 layers

The task was image recognition (10 million youtube video frames)

The hw comprised a 1000 computer network (16 000 cores), computation took three days.



# Current hardware

In 2012, Google trained a large network of 1.7 billion weights and 9 layers

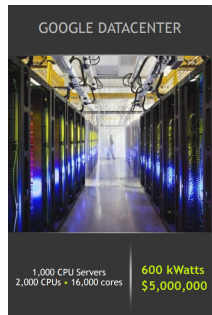
The task was image recognition (10 million youtube video frames)

The hw comprised a 1000 computer network (16 000 cores), computation took three days.

In 2014, similar task performed on Commodity Off-The-Shelf High Performance Computing (COTS HPC) technology: a cluster of GPU servers with Infiniband interconnects and MPI.

Able to train 1 billion parameter networks on just 3 machines in a couple of days.

Able to scale to 11 billion weights (approx. 6.5 times larger than the Google model) on 16 GPUs.

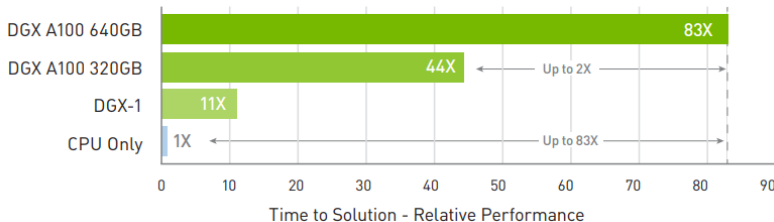


# Current hardware – NVIDIA DGX Station

- ▶ 8x GPU (Nvidia A100 80GB Tensor Core)
- ▶ 5 petaFLOPS
- ▶ System memory: 2 TB
- ▶ Network: 200 Gb/s InfiniBand



Up to 83X Higher Throughput than CPU, 2X Higher Throughput than DGX A100 320GB on Big Data Analytics Benchmark



# Deep learning in clouds

Big companies offer cloud services for deep learning:

- ▶ Amazon Web Services
- ▶ Google Cloud
- ▶ Deep Cognition
- ▶ ...

## **Advantages:**

- ▶ Do not have to care (too much) about technical problems.
- ▶ Do not have to buy and optimize highend hw/sw, networks etc.
- ▶ Scaling & virtually limitless storage.

## **Disadvantages:**

- ▶ Do not have full control.
- ▶ Performance can vary, connectivity problems.
- ▶ Have to pay for services.
- ▶ Privacy issues.

# Current software

- ▶ **TensorFlow** (Google)
  - ▶ open source software library for numerical computation using data flow graphs
  - ▶ allows implementation of most current neural networks
  - ▶ allows computation on multiple devices (CPUs, GPUs, ...)
  - ▶ Python API
  - ▶ **Keras**: a part of TensorFlow that allows easy description of most modern neural networks
- ▶ **PyTorch** (Facebook)
  - ▶ similar to TensorFlow
  - ▶ object oriented
  - ▶ ... majority of new models in research papers implemented in PyTorch

<https://www.cioinsight.com/big-data/pytorch-vs-tensorflow/>

- ▶ **Theano (dead)**:
  - ▶ The "academic" grand-daddy of deep-learning frameworks, written in Python. Strongly inspired TensorFlow (some people developing Theano moved on to develop TensorFlow).
- ▶ There are others: Caffe, Deeplearning4j, ...



# Current software – Keras

```
from keras.models import Sequential
from keras.layers import Dense, Dropout, Activation
from keras.optimizers import SGD

model = Sequential()
# Dense(64) is a fully-connected layer with 64 hidden units.
# in the first layer, you must specify the expected input data shape
# here, 20-dimensional vectors.
model.add(Dense(64, input_dim=20, init='uniform'))
model.add(Activation('tanh'))
model.add(Dropout(0.5))
model.add(Dense(64, init='uniform'))
model.add(Activation('tanh'))
model.add(Dropout(0.5))
model.add(Dense(10, init='uniform'))
model.add(Activation('softmax'))

sgd = SGD(lr=0.1, decay=1e-6, momentum=0.9, nesterov=True)
model.compile(loss='categorical_crossentropy',
              optimizer=sgd,
              metrics=['accuracy'])

model.fit(X_train, y_train,
          nb_epoch=20,
          batch_size=16)
score = model.evaluate(X_test, y_test, batch_size=16)
```

# Current software – Keras functional API

```
from keras.layers import Input, Dense
from keras.models import Model

# This returns a tensor
inputs = Input(shape=(784,))

# a layer instance is callable on a tensor, and returns a tensor
output_1 = Dense(64, activation='relu')(inputs)
output_2 = Dense(64, activation='relu')(output_1)
predictions = Dense(10, activation='softmax')(output_2)

# This creates a model that includes
# the Input layer and three Dense layers
model = Model(inputs=inputs, outputs=predictions)
model.compile(optimizer='rmsprop',
              loss='categorical_crossentropy',
              metrics=['accuracy'])
model.fit(data, labels) # starts training
```

# Current software – TensorFlow

```
41 # tf Graph input
42 X = tf.placeholder("float", [None, n_input])
43 Y = tf.placeholder("float", [None, n_classes])
44
45 # Store layers weight & bias
46 weights = {
47     'h1': tf.Variable(tf.random_normal([n_input, n_hidden_1])),
48     'h2': tf.Variable(tf.random_normal([n_hidden_1, n_hidden_2])),
49     'out': tf.Variable(tf.random_normal([n_hidden_2, n_classes]))
50 }
51 biases = {
52     'b1': tf.Variable(tf.random_normal([n_hidden_1])),
53     'b2': tf.Variable(tf.random_normal([n_hidden_2])),
54     'out': tf.Variable(tf.random_normal([n_classes]))
55 }
```

# Current software – TensorFlow

```
58 # Create model
59 def multilayer_perceptron(x):
60     # Hidden fully connected layer with 256 neurons
61     layer_1 = tf.add(tf.matmul(x, weights['h1']), biases['b1'])
62     # Hidden fully connected layer with 256 neurons
63     layer_2 = tf.add(tf.matmul(layer_1, weights['h2']), biases['b2'])
64     # Output fully connected layer with a neuron for each class
65     out_layer = tf.matmul(layer_2, weights['out']) + biases['out']
66     return out_layer
67
68 # Construct model
69 logits = multilayer_perceptron(X)
```

# Current software – PyTorch

```
36 class Net(nn.Module):
37     def __init__(self, input_size, hidden_size, num_classes):
38         super(Net, self).__init__()
39         self.fc1 = nn.Linear(input_size, hidden_size)
40         self.relu = nn.ReLU()
41         self.fc2 = nn.Linear(hidden_size, num_classes)
42
43     def forward(self, x):
44         out = self.fc1(x)
45         out = self.relu(out)
46         out = self.fc2(out)
47         return out
48
49 net = Net(input_size, hidden_size, num_classes)
```

# Other software implementations

Most "mathematical" software packages contain some support of neural networks:

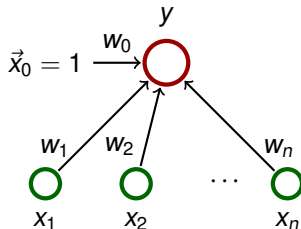
- ▶ MATLAB
- ▶ R
- ▶ STATISTICA
- ▶ Weka
- ▶ ...

The implementations are typically not on par with the previously mentioned dedicated deep-learning libraries.

## Training linear models

# Linear regression (ADALINE)

## Architecture:



$\vec{w} = (w_0, w_1, \dots, w_n)$  and  $\vec{x} = (x_0, x_1, \dots, x_n)$  where  $x_0 = 1$ .

## Activity:

- ▶ inner potential:  $\xi = w_0 + \sum_{i=1}^n w_i x_i = \sum_{i=0}^n w_i x_i = \vec{w} \cdot \vec{x}$
- ▶ activation function:  $\sigma(\xi) = \xi$
- ▶ network function:  $y[\vec{w}](\vec{x}) = \sigma(\xi) = \vec{w} \cdot \vec{x}$



# Linear regression (ADALINE)

## Learning:

- ▶ Given a **training dataset**

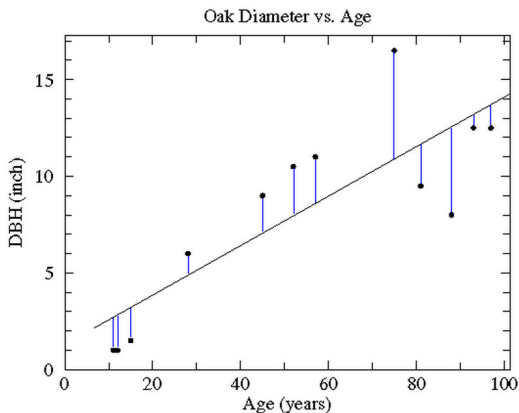
$$\mathcal{T} = \{(\vec{x}_1, d_1), (\vec{x}_2, d_2), \dots, (\vec{x}_p, d_p)\}$$

Here  $\vec{x}_k = (x_{k0}, x_{k1}, \dots, x_{kn}) \in \mathbb{R}^{n+1}$ ,  $x_{k0} = 1$ , is the  $k$ -th input, and  $d_k \in \mathbb{R}$  is the expected output.

Intuition: The network is supposed to compute an affine approximation of the function (some of) whose values are given in the training set.

# Oaks in Wisconsin

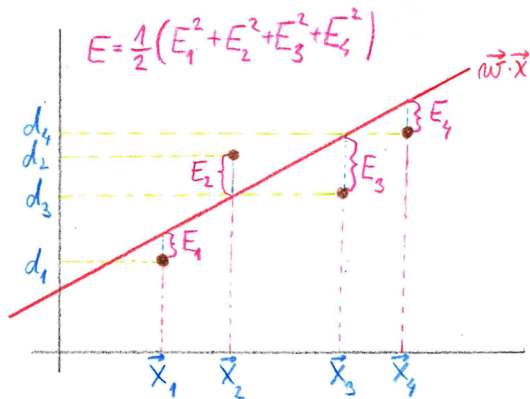
Age (years)	DBH (inch)
97	12.5
93	12.5
88	8.0
81	9.5
75	16.5
57	11.0
52	10.5
45	9.0
28	6.0
15	1.5
12	1.0
11	1.0



# Linear regression (ADALINE)

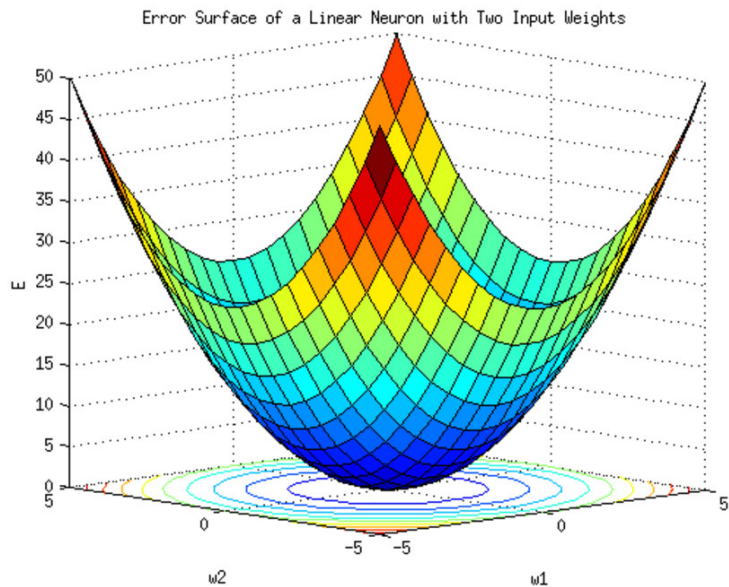
► Error function:

$$E(\vec{w}) = \frac{1}{2} \sum_{k=1}^p (\vec{w} \cdot \vec{x}_k - d_k)^2 = \frac{1}{2} \sum_{k=1}^p \left( \sum_{i=0}^n w_i x_{ki} - d_k \right)^2$$



► The goal is to find  $\vec{w}$  which minimizes  $E(\vec{w})$ .

# Error function



# Gradient of the error function

Consider **gradient** of the error function:

$$\nabla E(\vec{w}) = \left( \frac{\partial E}{\partial w_0}(\vec{w}), \dots, \frac{\partial E}{\partial w_n}(\vec{w}) \right)$$

Intuition:  $\nabla E(\vec{w})$  is a vector in the **weight space** which points in the direction of the *steepest ascent* of the error function.

Note that the vectors  $\vec{x}_k$  are just parameters of the function  $E$ , and are thus fixed!

# Gradient of the error function

Consider **gradient** of the error function:

$$\nabla E(\vec{w}) = \left( \frac{\partial E}{\partial w_0}(\vec{w}), \dots, \frac{\partial E}{\partial w_n}(\vec{w}) \right)$$

Intuition:  $\nabla E(\vec{w})$  is a vector in the **weight space** which points in the direction of the *steepest ascent* of the error function.

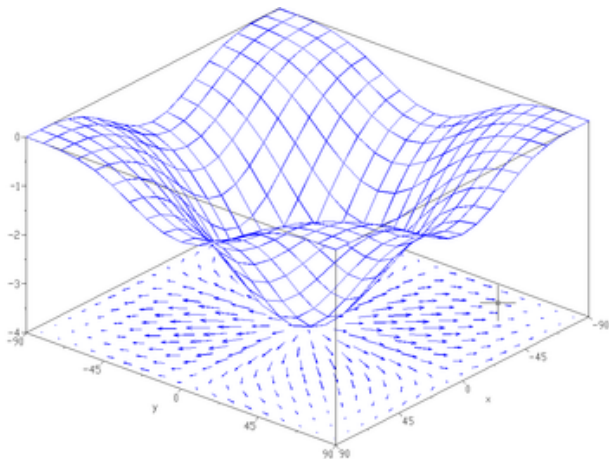
Note that the vectors  $\vec{x}_k$  are just parameters of the function  $E$ , and are thus fixed!

## Fact

If  $\nabla E(\vec{w}) = \vec{0} = (0, \dots, 0)$ , then  $\vec{w}$  is a global minimum of  $E$ .

For ADALINE, the error function  $E(\vec{w})$  is a convex paraboloid and thus has the unique global minimum.

# Gradient - illustration



Caution! This picture just illustrates the notion of gradient ... it is not the convex paraboloid  $E(\vec{w})$  !

# Gradient of the error function

$$\frac{\partial E}{\partial \mathbf{w}_\ell}(\vec{\mathbf{w}}) = \frac{1}{2} \sum_{k=1}^p \frac{\delta}{\delta \mathbf{w}_\ell} \left( \sum_{i=0}^n w_i x_{ki} - d_k \right)^2$$



# Gradient of the error function

$$\begin{aligned}\frac{\partial E}{\partial \mathbf{w}_\ell}(\vec{\mathbf{w}}) &= \frac{1}{2} \sum_{k=1}^p \frac{\delta}{\delta \mathbf{w}_\ell} \left( \sum_{i=0}^n w_i x_{ki} - d_k \right)^2 \\ &= \frac{1}{2} \sum_{k=1}^p 2 \left( \sum_{i=0}^n w_i x_{ki} - d_k \right) \frac{\delta}{\delta \mathbf{w}_\ell} \left( \sum_{i=0}^n w_i x_{ki} - d_k \right)\end{aligned}$$

# Gradient of the error function

$$\begin{aligned}\frac{\partial E}{\partial \mathbf{w}_\ell}(\vec{\mathbf{w}}) &= \frac{1}{2} \sum_{k=1}^p \frac{\delta}{\delta \mathbf{w}_\ell} \left( \sum_{i=0}^n w_i x_{ki} - d_k \right)^2 \\&= \frac{1}{2} \sum_{k=1}^p 2 \left( \sum_{i=0}^n w_i x_{ki} - d_k \right) \frac{\delta}{\delta \mathbf{w}_\ell} \left( \sum_{i=0}^n w_i x_{ki} - d_k \right) \\&= \frac{1}{2} \sum_{k=1}^p 2 \left( \sum_{i=0}^n w_i x_{ki} - d_k \right) \left( \sum_{i=0}^n \left( \frac{\delta}{\delta \mathbf{w}_\ell} w_i x_{ki} \right) - \frac{\delta E}{\delta \mathbf{w}_\ell} d_k \right)\end{aligned}$$

# Gradient of the error function

$$\begin{aligned}\frac{\partial E}{\partial \mathbf{w}_\ell}(\vec{\mathbf{w}}) &= \frac{1}{2} \sum_{k=1}^p \frac{\delta}{\delta \mathbf{w}_\ell} \left( \sum_{i=0}^n w_i x_{ki} - d_k \right)^2 \\&= \frac{1}{2} \sum_{k=1}^p 2 \left( \sum_{i=0}^n w_i x_{ki} - d_k \right) \frac{\delta}{\delta \mathbf{w}_\ell} \left( \sum_{i=0}^n w_i x_{ki} - d_k \right) \\&= \frac{1}{2} \sum_{k=1}^p 2 \left( \sum_{i=0}^n w_i x_{ki} - d_k \right) \left( \sum_{i=0}^n \left( \frac{\delta}{\delta \mathbf{w}_\ell} w_i x_{ki} \right) - \frac{\delta E}{\delta \mathbf{w}_\ell} d_k \right) \\&= \sum_{k=1}^p \left( \vec{\mathbf{w}} \cdot \vec{\mathbf{x}}_k - d_k \right) \mathbf{x}_{k\ell}\end{aligned}$$

# Gradient of the error function

$$\begin{aligned}\frac{\partial E}{\partial \mathbf{w}_\ell}(\vec{\mathbf{w}}) &= \frac{1}{2} \sum_{k=1}^p \frac{\delta}{\delta \mathbf{w}_\ell} \left( \sum_{i=0}^n \mathbf{w}_i x_{ki} - d_k \right)^2 \\&= \frac{1}{2} \sum_{k=1}^p 2 \left( \sum_{i=0}^n \mathbf{w}_i x_{ki} - d_k \right) \frac{\delta}{\delta \mathbf{w}_\ell} \left( \sum_{i=0}^n \mathbf{w}_i x_{ki} - d_k \right) \\&= \frac{1}{2} \sum_{k=1}^p 2 \left( \sum_{i=0}^n \mathbf{w}_i x_{ki} - d_k \right) \left( \sum_{i=0}^n \left( \frac{\delta}{\delta \mathbf{w}_\ell} \mathbf{w}_i x_{ki} \right) - \frac{\delta E}{\delta \mathbf{w}_\ell} d_k \right) \\&= \sum_{k=1}^p \left( \vec{\mathbf{w}} \cdot \vec{\mathbf{x}}_k - d_k \right) \mathbf{x}_{k\ell}\end{aligned}$$

Thus

$$\nabla E(\vec{\mathbf{w}}) = \left( \frac{\partial E}{\partial \mathbf{w}_0}(\vec{\mathbf{w}}), \dots, \frac{\partial E}{\partial \mathbf{w}_n}(\vec{\mathbf{w}}) \right) = \sum_{k=1}^p \left( \vec{\mathbf{w}} \cdot \vec{\mathbf{x}}_k - d_k \right) \vec{\mathbf{x}}_k$$

# Linear regression - learning

## **Batch algorithm (gradient descent):**

**Idea:** In every step "move" the weights in the direction *opposite* to the gradient.

# Linear regression - learning

## **Batch algorithm (gradient descent):**

**Idea:** In every step "move" the weights in the direction *opposite* to the gradient.

The algorithm computes a sequence of weight vectors

$\vec{w}^{(0)}, \vec{w}^{(1)}, \vec{w}^{(2)}, \dots$

- weights in  $\vec{w}^{(0)}$  are randomly initialized to values close to 0

# Linear regression - learning

## Batch algorithm (gradient descent):

**Idea:** In every step "move" the weights in the direction *opposite* to the gradient.

The algorithm computes a sequence of weight vectors

$\vec{w}^{(0)}, \vec{w}^{(1)}, \vec{w}^{(2)}, \dots$

- ▶ weights in  $\vec{w}^{(0)}$  are randomly initialized to values close to 0
- ▶ in the step  $t + 1$ , weights  $\vec{w}^{(t+1)}$  are computed as follows:

$$\begin{aligned}\vec{w}^{(t+1)} &= \vec{w}^{(t)} - \varepsilon \cdot \nabla E(\vec{w}^{(t)}) \\ &= \vec{w}^{(t)} - \varepsilon \cdot \sum_{k=1}^p (\vec{w}^{(t)} \cdot \vec{x}_k - d_k) \cdot \vec{x}_k\end{aligned}$$

Here  $k = (t \bmod p) + 1$  and  $0 < \varepsilon \leq 1$  is a *learning rate*.

# Linear regression - learning

## Batch algorithm (gradient descent):

**Idea:** In every step "move" the weights in the direction *opposite* to the gradient.

The algorithm computes a sequence of weight vectors

$\vec{w}^{(0)}, \vec{w}^{(1)}, \vec{w}^{(2)}, \dots$

- ▶ weights in  $\vec{w}^{(0)}$  are randomly initialized to values close to 0
- ▶ in the step  $t + 1$ , weights  $\vec{w}^{(t+1)}$  are computed as follows:

$$\begin{aligned}\vec{w}^{(t+1)} &= \vec{w}^{(t)} - \varepsilon \cdot \nabla E(\vec{w}^{(t)}) \\ &= \vec{w}^{(t)} - \varepsilon \cdot \sum_{k=1}^p (\vec{w}^{(t)} \cdot \vec{x}_k - d_k) \cdot \vec{x}_k\end{aligned}$$

Here  $k = (t \bmod p) + 1$  and  $0 < \varepsilon \leq 1$  is a *learning rate*.

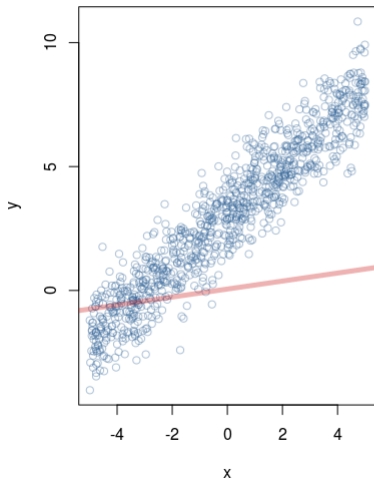
## Proposition

For sufficiently small  $\varepsilon > 0$  the sequence  $\vec{w}^{(0)}, \vec{w}^{(1)}, \vec{w}^{(2)}, \dots$  converges (componentwise) to the global minimum of  $E$  (i.e. to the vector  $\vec{w}$  satisfying  $\nabla E(\vec{w}) = \vec{0}$ ).

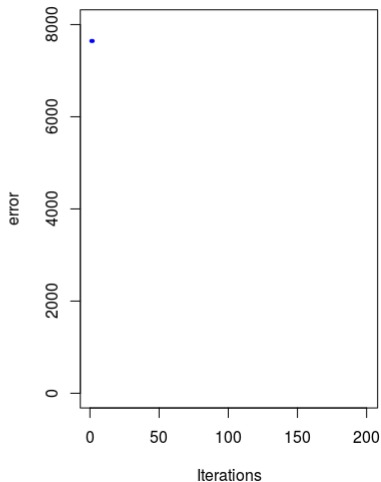


# Linear regression - animation

Linear regression by gradient descent

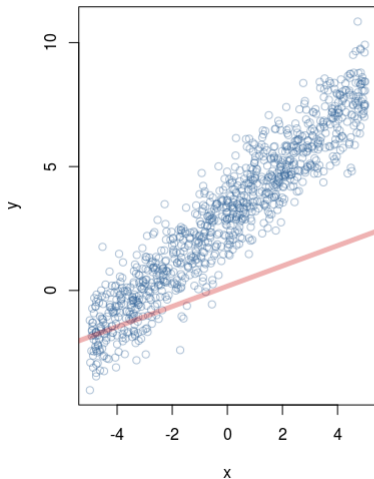


Error function

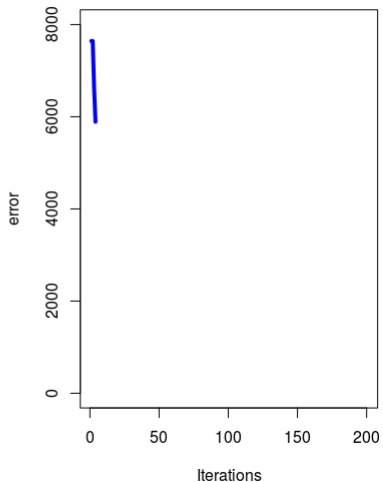


# Linear regression - animation

Linear regression by gradient descent

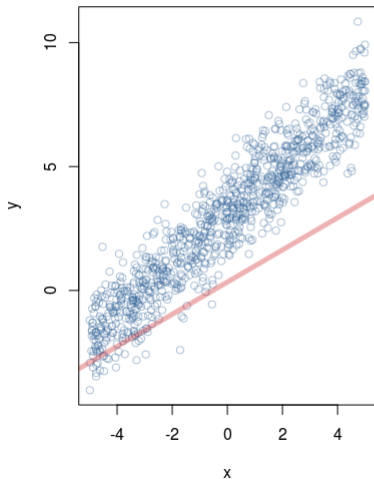


Error function

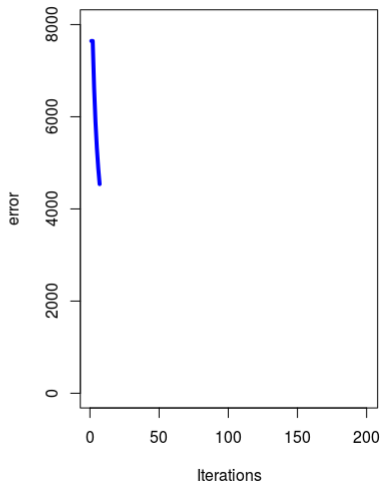


# Linear regression - animation

Linear regression by gradient descent

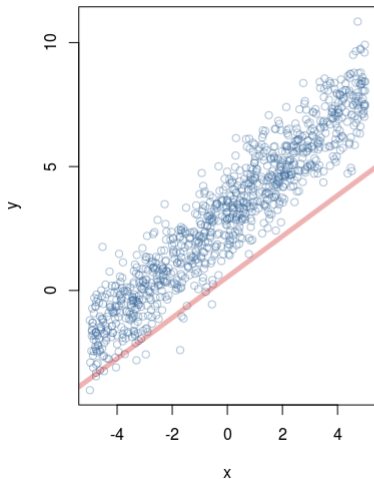


Error function

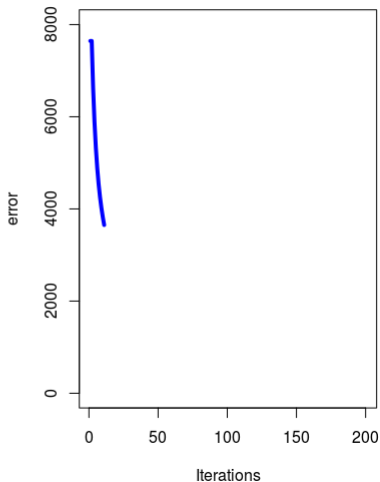


# Linear regression - animation

Linear regression by gradient descent

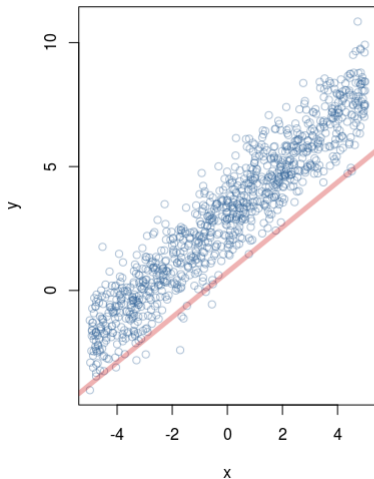


Error function

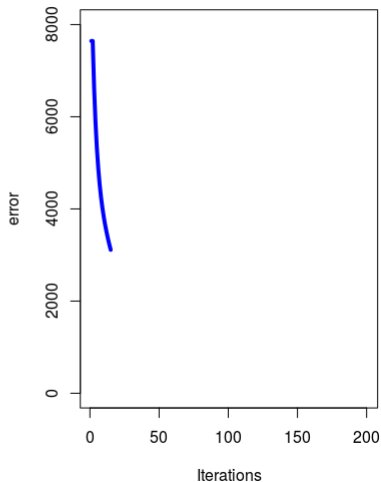


# Linear regression - animation

Linear regression by gradient descent

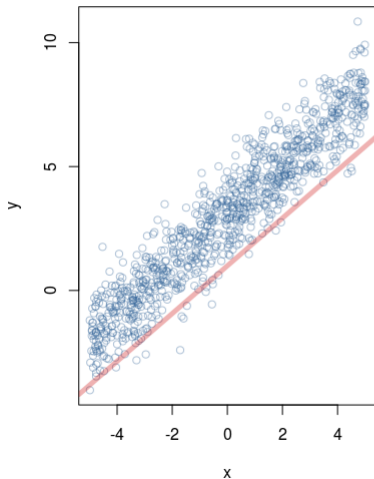


Error function

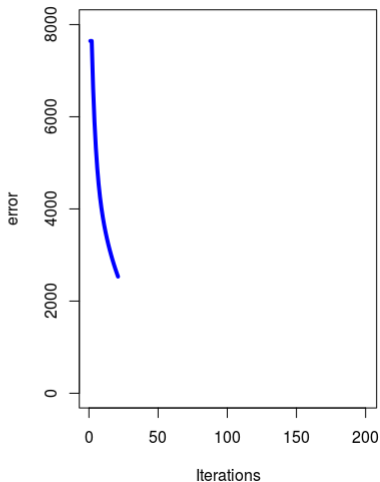


# Linear regression - animation

Linear regression by gradient descent

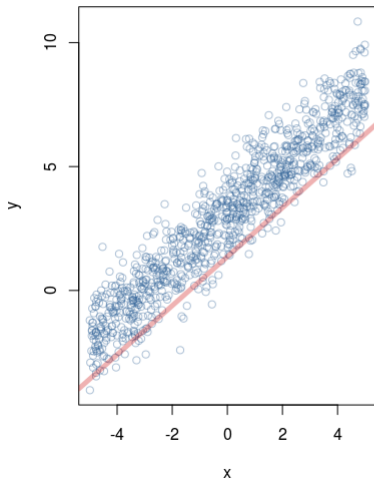


Error function

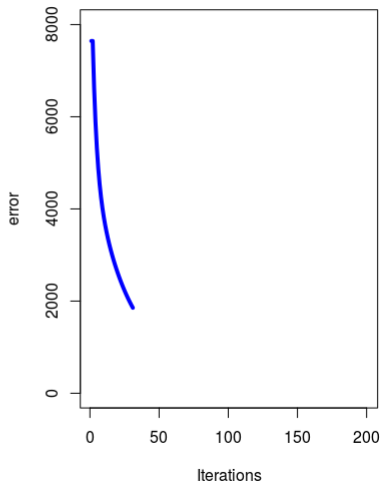


# Linear regression - animation

Linear regression by gradient descent

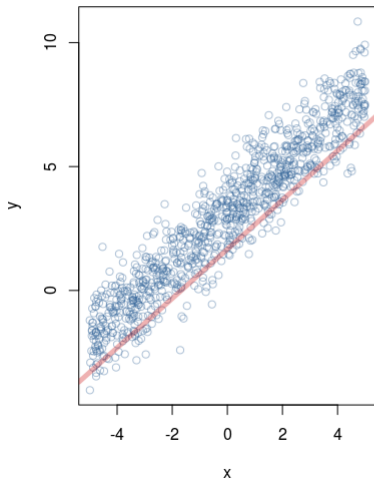


Error function

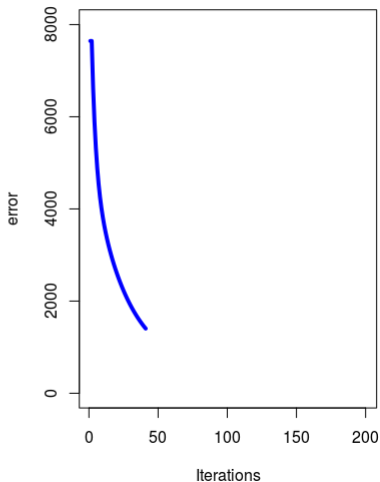


# Linear regression - animation

Linear regression by gradient descent



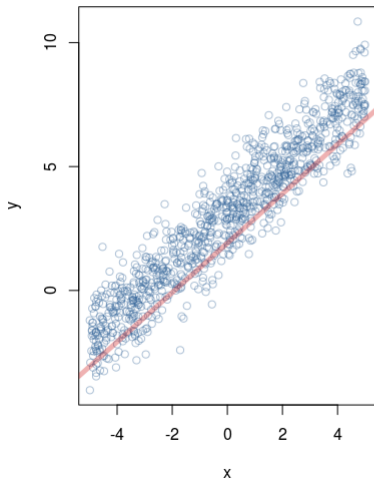
Error function



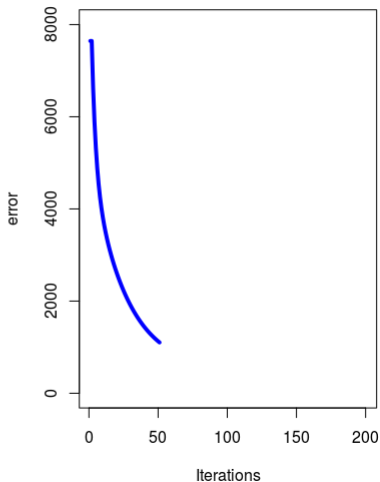


# Linear regression - animation

Linear regression by gradient descent

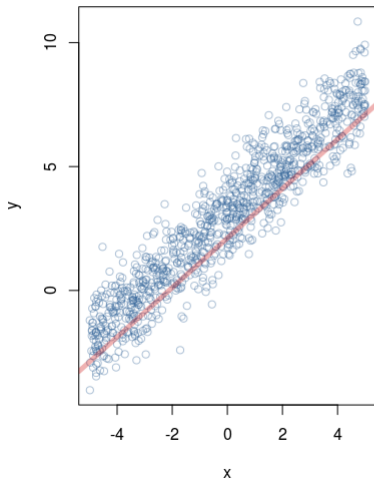


Error function

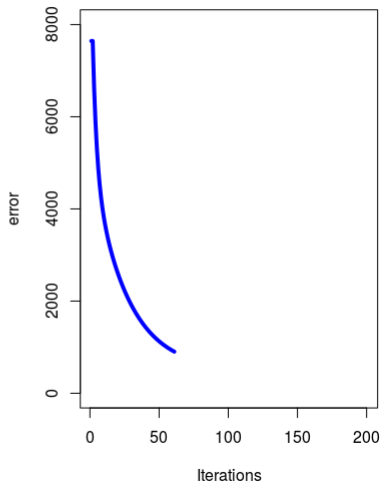


# Linear regression - animation

Linear regression by gradient descent

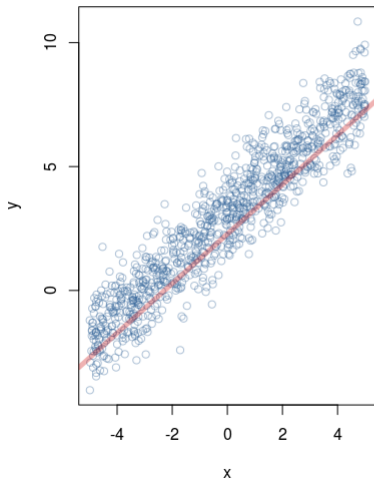


Error function

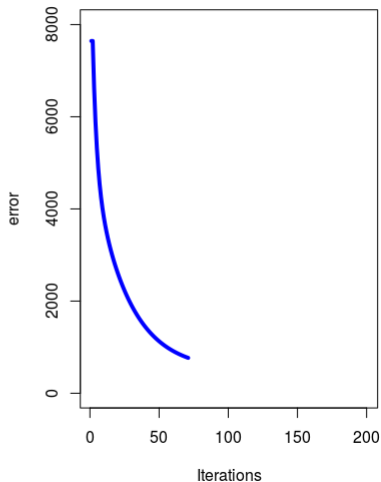


# Linear regression - animation

Linear regression by gradient descent

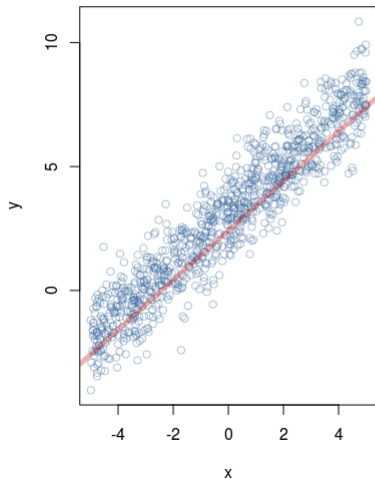


Error function

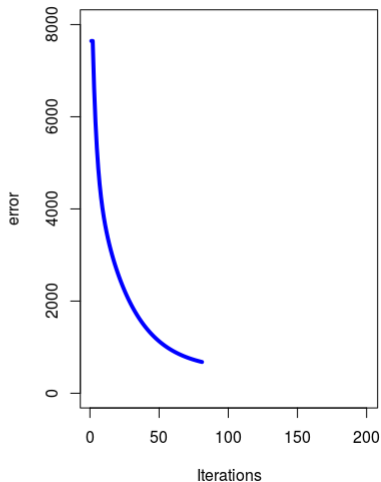


# Linear regression - animation

Linear regression by gradient descent

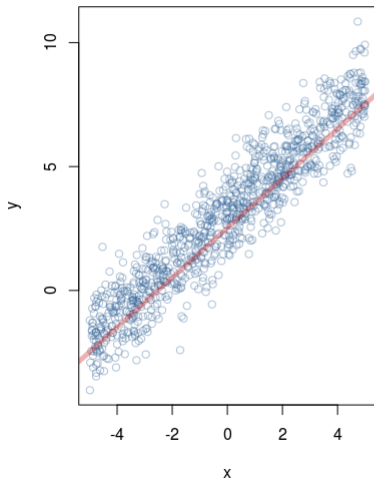


Error function

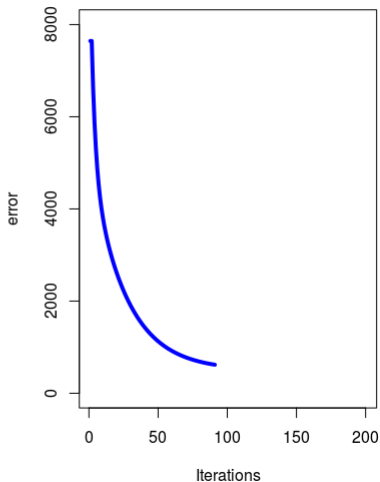


# Linear regression - animation

Linear regression by gradient descent

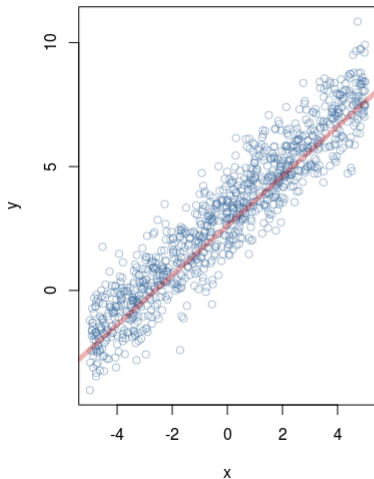


Error function

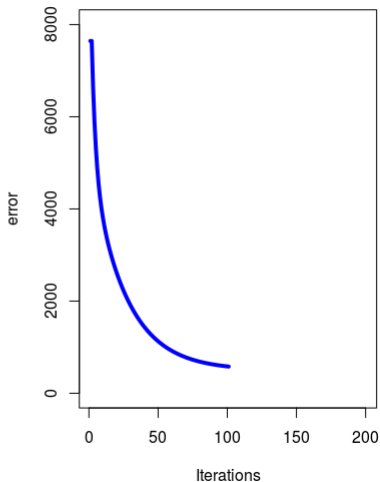


# Linear regression - animation

Linear regression by gradient descent

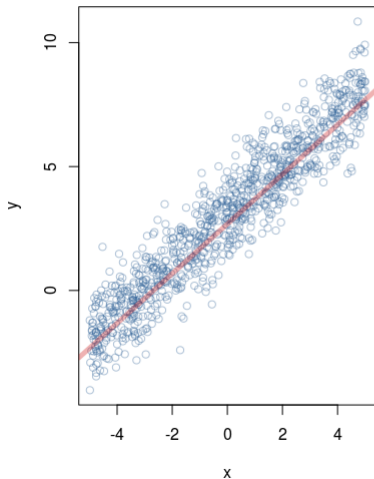


Error function

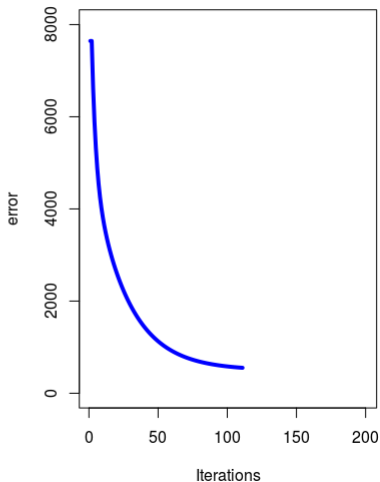


# Linear regression - animation

Linear regression by gradient descent

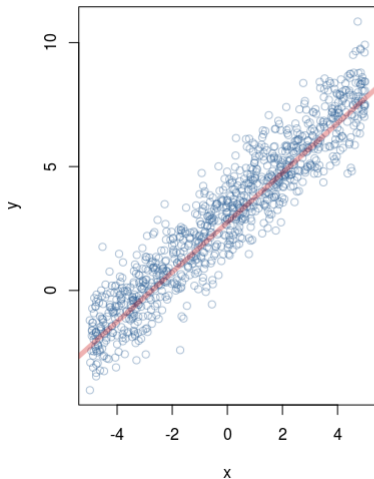


Error function

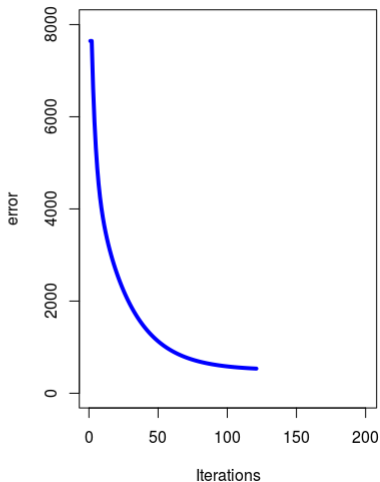


# Linear regression - animation

Linear regression by gradient descent



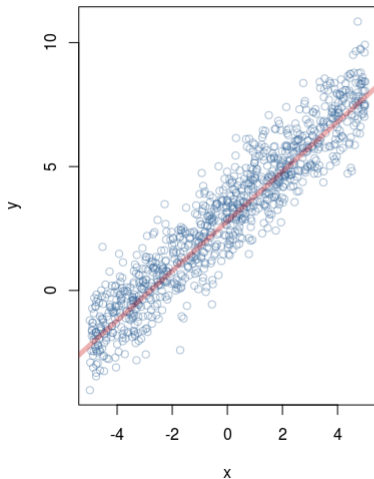
Error function



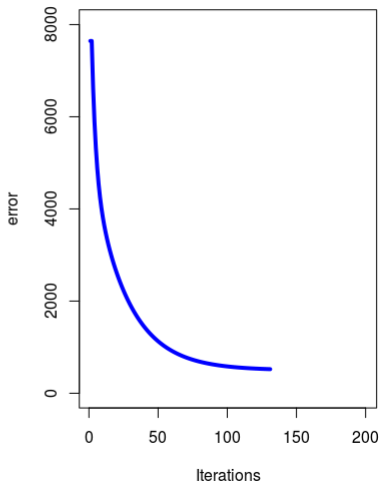


# Linear regression - animation

Linear regression by gradient descent

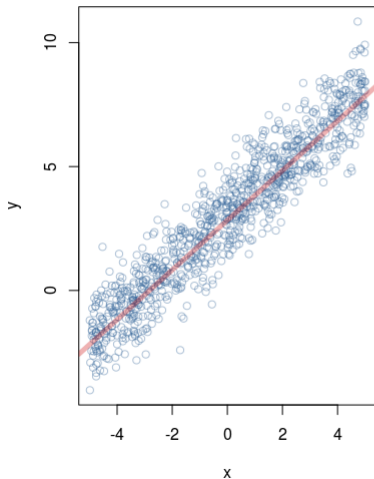


Error function

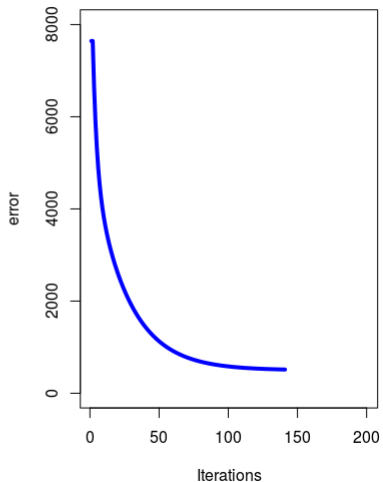


# Linear regression - animation

Linear regression by gradient descent

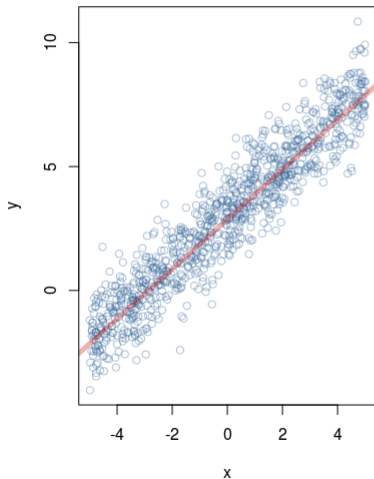


Error function

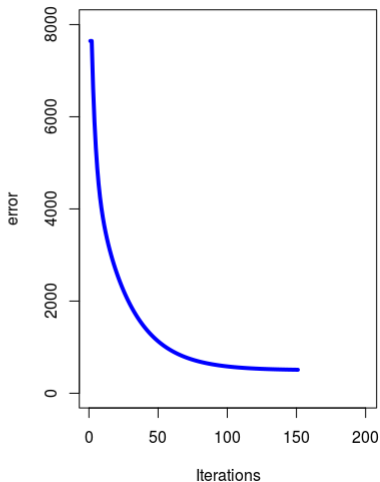


# Linear regression - animation

Linear regression by gradient descent

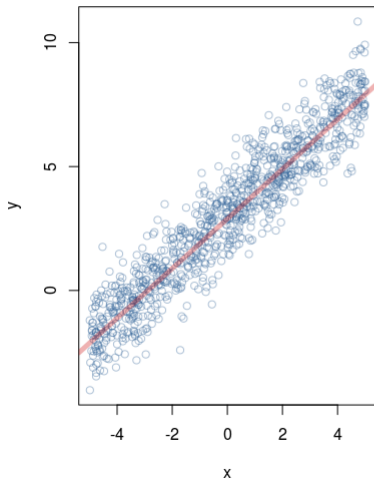


Error function

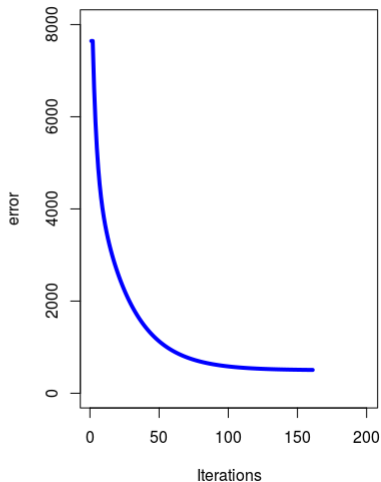


# Linear regression - animation

Linear regression by gradient descent

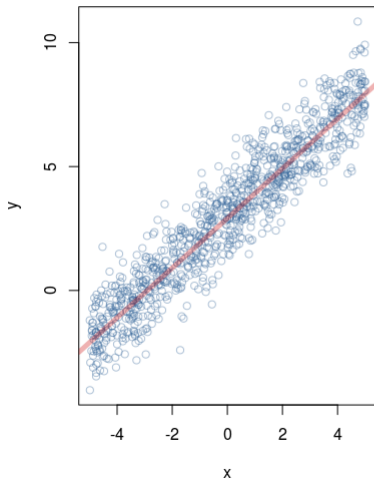


Error function

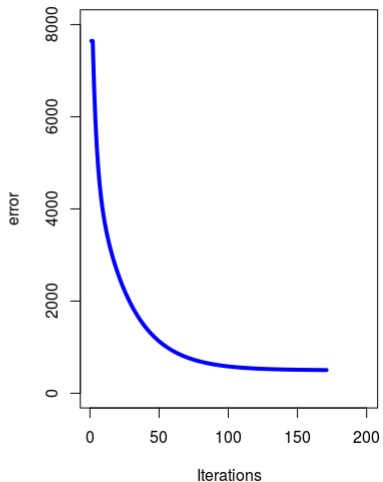


# Linear regression - animation

Linear regression by gradient descent

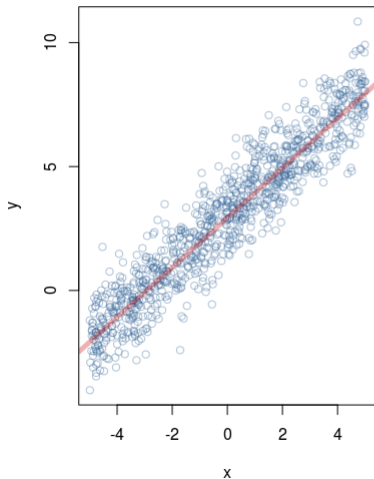


Error function

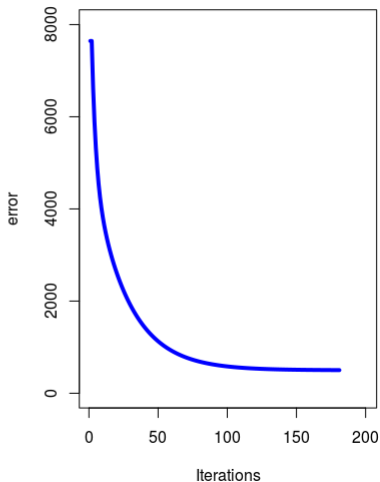


# Linear regression - animation

Linear regression by gradient descent

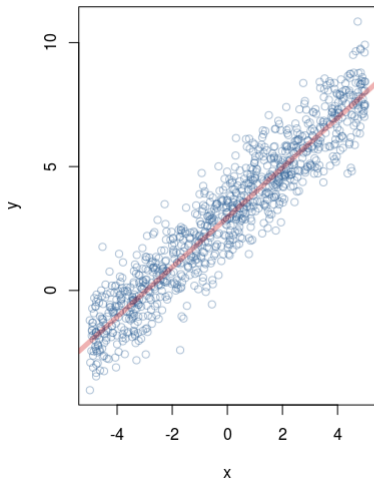


Error function

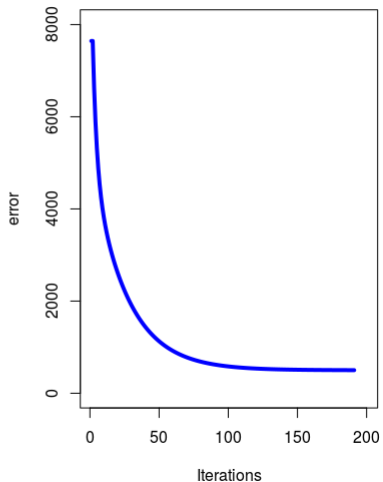


# Linear regression - animation

Linear regression by gradient descent

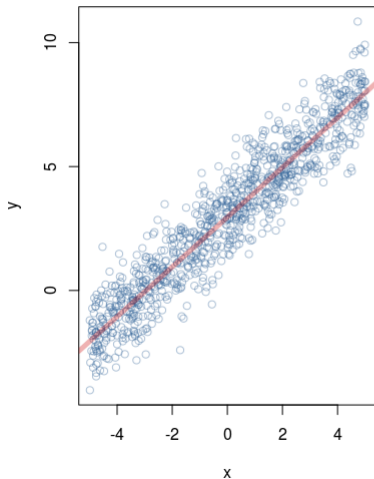


Error function

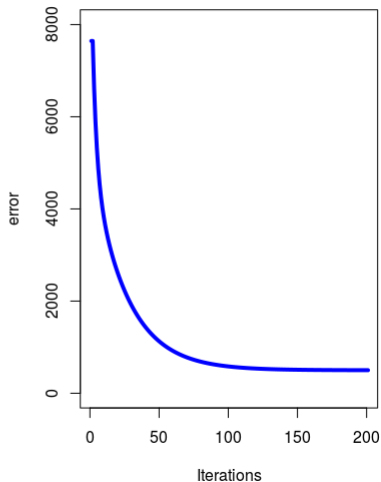


# Linear regression - animation

Linear regression by gradient descent



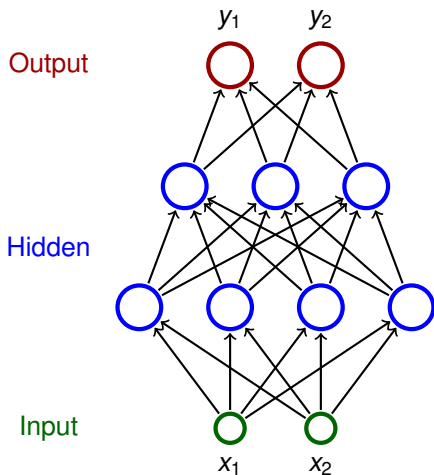
Error function





## MLP training – theory

# Architecture – Multilayer Perceptron (MLP)



- ▶ Neurons partitioned into **layers**; one input layer, one output layer, possibly several hidden layers
- ▶ layers numbered from 0; the input layer has number 0
  - ▶ E.g. three-layer network has two hidden layers and one output layer
- ▶ Neurons in the  $i$ -th layer are connected with all neurons in the  $i + 1$ -st layer
- ▶ Architecture of a MLP is typically described by numbers of neurons in individual layers (e.g. 2-4-3-2)

# MLP – architecture

## Notation:

- ▶ Denote
  - ▶  $X$  a set of *input* neurons
  - ▶  $Y$  a set of *output* neurons
  - ▶  $Z$  a set of *all* neurons ( $X, Y \subseteq Z$ )

## Notation:

- ▶ Denote
  - ▶  $X$  a set of *input* neurons
  - ▶  $Y$  a set of *output* neurons
  - ▶  $Z$  a set of *all* neurons ( $X, Y \subseteq Z$ )
- ▶ individual neurons denoted by indices  $i, j$  etc.
  - ▶  $\xi_j$  is the inner potential of the neuron  $j$  *after the computation stops*

# MLP – architecture

## Notation:

- ▶ Denote
  - ▶  $X$  a set of *input* neurons
  - ▶  $Y$  a set of *output* neurons
  - ▶  $Z$  a set of *all* neurons ( $X, Y \subseteq Z$ )
- ▶ individual neurons denoted by indices  $i, j$  etc.
  - ▶  $\xi_j$  is the inner potential of the neuron  $j$  *after the computation stops*
  - ▶  $y_j$  is the output of the neuron  $j$  *after the computation stops*

(define  $y_0 = 1$  is the value of the formal unit input)

# MLP – architecture

## Notation:

- ▶ Denote
  - ▶  $X$  a set of *input* neurons
  - ▶  $Y$  a set of *output* neurons
  - ▶  $Z$  a set of *all* neurons ( $X, Y \subseteq Z$ )
- ▶ individual neurons denoted by indices  $i, j$  etc.
  - ▶  $\xi_j$  is the inner potential of the neuron  $j$  *after the computation stops*
  - ▶  $y_j$  is the output of the neuron  $j$  *after the computation stops*

(define  $y_0 = 1$  is the value of the formal unit input)

- ▶  $w_{ji}$  is the weight of the connection **from  $i$  to  $j$**

(in particular,  $w_{j0}$  is the weight of the connection from the formal unit input, i.e.  $w_{j0} = -b_j$  where  $b_j$  is the bias of the neuron  $j$ )

## Notation:

- ▶ Denote
  - ▶  $X$  a set of *input* neurons
  - ▶  $Y$  a set of *output* neurons
  - ▶  $Z$  a set of *all* neurons ( $X, Y \subseteq Z$ )
- ▶ individual neurons denoted by indices  $i, j$  etc.
  - ▶  $\xi_j$  is the inner potential of the neuron  $j$  *after the computation stops*
  - ▶  $y_j$  is the output of the neuron  $j$  *after the computation stops*

(define  $y_0 = 1$  is the value of the formal unit input)

- ▶  $w_{ji}$  is the weight of the connection **from  $i$  to  $j$**   
(in particular,  $w_{j0}$  is the weight of the connection from the formal unit input, i.e.  $w_{j0} = -b_j$  where  $b_j$  is the bias of the neuron  $j$ )
- ▶  $j_{\leftarrow}$  is a set of all  $i$  such that  $j$  is adjacent from  $i$   
(i.e. there is an arc **to**  $j$  from  $i$ )

## Notation:

- ▶ Denote
  - ▶  $X$  a set of *input* neurons
  - ▶  $Y$  a set of *output* neurons
  - ▶  $Z$  a set of *all* neurons ( $X, Y \subseteq Z$ )
- ▶ individual neurons denoted by indices  $i, j$  etc.
  - ▶  $\xi_j$  is the inner potential of the neuron  $j$  *after the computation stops*
  - ▶  $y_j$  is the output of the neuron  $j$  *after the computation stops*

(define  $y_0 = 1$  is the value of the formal unit input)

- ▶  $w_{ji}$  is the weight of the connection **from  $i$  to  $j$**   
(in particular,  $w_{j0}$  is the weight of the connection from the formal unit input, i.e.  $w_{j0} = -b_j$  where  $b_j$  is the bias of the neuron  $j$ )
- ▶  $j_{\leftarrow}$  is a set of all  $i$  such that  $j$  is adjacent from  $i$   
(i.e. there is an arc **to**  $j$  from  $i$ )
- ▶  $j_{\rightarrow}$  is a set of all  $i$  such that  $j$  is adjacent to  $i$   
(i.e. there is an arc **from**  $j$  to  $i$ )



# MLP – activity

## Activity:

- ▶ inner potential of neuron  $j$ :

$$\xi_j = \sum_{i \in j_{\leftarrow}} w_{ji} y_i$$

# MLP – activity

## Activity:

- ▶ inner potential of neuron  $j$ :

$$\xi_j = \sum_{i \in j_{\leftarrow}} w_{ji} y_i$$

- ▶ activation function  $\sigma_j$  for neuron  $j$  (arbitrary differentiable)

## Activity:

- ▶ inner potential of neuron  $j$ :

$$\xi_j = \sum_{i \in j_{\leftarrow}} w_{ji} y_i$$

- ▶ activation function  $\sigma_j$  for neuron  $j$  (arbitrary differentiable)
- ▶ State of non-input neuron  $j \in Z \setminus X$  after the computation stops:

$$y_j = \sigma_j(\xi_j)$$

( $y_j$  depends on the configuration  $\vec{w}$  and the input  $\vec{x}$ , so we sometimes write  $y_j(\vec{w}, \vec{x})$ )

# MLP – activity

## Activity:

- ▶ inner potential of neuron  $j$ :

$$\xi_j = \sum_{i \in j_{\leftarrow}} w_{ji} y_i$$

- ▶ activation function  $\sigma_j$  for neuron  $j$  (arbitrary differentiable)
- ▶ State of non-input neuron  $j \in Z \setminus X$  after the computation stops:

$$y_j = \sigma_j(\xi_j)$$

( $y_j$  depends on the configuration  $\vec{w}$  and the input  $\vec{x}$ , so we sometimes write  $y_j(\vec{w}, \vec{x})$ )

- ▶ The network computes a function  $\mathbb{R}^{|X|}$  to  $\mathbb{R}^{|Y|}$ . Layer-wise computation: First, all input neurons are assigned values of the input. In the  $\ell$ -th step, all neurons of the  $\ell$ -th layer are evaluated.

## Learning:

- ▶ Given a **training dataset**  $\mathcal{T}$  of the form

$$\left\{ \left( \vec{x}_k, \vec{d}_k \right) \mid k = 1, \dots, p \right\}$$

Here, every  $\vec{x}_k \in \mathbb{R}^{|X|}$  is an *input vector* and every  $\vec{d}_k \in \mathbb{R}^{|Y|}$  is the desired network output. For every  $j \in Y$ , denote by  $d_{kj}$  the desired output of the neuron  $j$  for a given network input  $\vec{x}_k$  (the vector  $\vec{d}_k$  can be written as  $(d_{kj})_{j \in Y}$ ).

## Learning:

- ▶ Given a **training dataset**  $\mathcal{T}$  of the form

$$\left\{ \left( \vec{x}_k, \vec{d}_k \right) \mid k = 1, \dots, p \right\}$$

Here, every  $\vec{x}_k \in \mathbb{R}^{|X|}$  is an *input vector* and every  $\vec{d}_k \in \mathbb{R}^{|Y|}$  is the desired network output. For every  $j \in Y$ , denote by  $d_{kj}$  the desired output of the neuron  $j$  for a given network input  $\vec{x}_k$  (the vector  $\vec{d}_k$  can be written as  $(d_{kj})_{j \in Y}$ ).

- ▶ **Error function:**

$$E(\vec{w}) = \sum_{k=1}^p E_k(\vec{w})$$

where

$$E_k(\vec{w}) = \frac{1}{2} \sum_{j \in Y} \left( y_j(\vec{w}, \vec{x}_k) - d_{kj} \right)^2$$

# MLP – learning algorithm

## Batch algorithm (gradient descent):

The algorithm computes a sequence of weight vectors

$\vec{w}^{(0)}, \vec{w}^{(1)}, \vec{w}^{(2)}, \dots$

- ▶ weights in  $\vec{w}^{(0)}$  are randomly initialized to values close to 0
- ▶ in the step  $t + 1$  (here  $t = 0, 1, 2 \dots$ ), weights  $\vec{w}^{(t+1)}$  are computed as follows:

$$w_{ji}^{(t+1)} = w_{ji}^{(t)} + \Delta w_{ji}^{(t)}$$

# MLP – learning algorithm

## Batch algorithm (gradient descent):

The algorithm computes a sequence of weight vectors

$\vec{w}^{(0)}, \vec{w}^{(1)}, \vec{w}^{(2)}, \dots$

- ▶ weights in  $\vec{w}^{(0)}$  are randomly initialized to values close to 0
- ▶ in the step  $t + 1$  (here  $t = 0, 1, 2 \dots$ ), weights  $\vec{w}^{(t+1)}$  are computed as follows:

$$w_{ji}^{(t+1)} = w_{ji}^{(t)} + \Delta w_{ji}^{(t)}$$

where

$$\Delta w_{ji}^{(t)} = -\varepsilon(t) \cdot \frac{\partial E}{\partial w_{ji}}(\vec{w}^{(t)})$$

is a *weight update* of  $w_{ji}$  in step  $t + 1$  and  $0 < \varepsilon(t) \leq 1$  is a *learning rate* in step  $t + 1$ .



# MLP – learning algorithm

## Batch algorithm (gradient descent):

The algorithm computes a sequence of weight vectors

$\vec{w}^{(0)}, \vec{w}^{(1)}, \vec{w}^{(2)}, \dots$

- ▶ weights in  $\vec{w}^{(0)}$  are randomly initialized to values close to 0
- ▶ in the step  $t + 1$  (here  $t = 0, 1, 2 \dots$ ), weights  $\vec{w}^{(t+1)}$  are computed as follows:

$$w_{ji}^{(t+1)} = w_{ji}^{(t)} + \Delta w_{ji}^{(t)}$$

where

$$\Delta w_{ji}^{(t)} = -\varepsilon(t) \cdot \frac{\partial E}{\partial w_{ji}}(\vec{w}^{(t)})$$

is a *weight update* of  $w_{ji}$  in step  $t + 1$  and  $0 < \varepsilon(t) \leq 1$  is a *learning rate* in step  $t + 1$ .

Note that  $\frac{\partial E}{\partial w_{ji}}(\vec{w}^{(t)})$  is a component of the gradient  $\nabla E$ , i.e. the weight update can be written as  $\vec{w}^{(t+1)} = \vec{w}^{(t)} - \varepsilon(t) \cdot \nabla E(\vec{w}^{(t)})$ .

# MLP – error function gradient

For every  $w_{ji}$  we have

$$\frac{\partial E}{\partial w_{ji}} = \sum_{k=1}^p \frac{\partial E_k}{\partial w_{ji}}$$

# MLP – error function gradient

For every  $w_{ji}$  we have

$$\frac{\partial E}{\partial w_{ji}} = \sum_{k=1}^p \frac{\partial E_k}{\partial w_{ji}}$$

where for every  $k = 1, \dots, p$  holds

$$\frac{\partial E_k}{\partial w_{ji}} = \frac{\partial E_k}{\partial y_j} \cdot \sigma'_j(\xi_j) \cdot y_i$$

# MLP – error function gradient

For every  $w_{ji}$  we have

$$\frac{\partial E}{\partial w_{ji}} = \sum_{k=1}^p \frac{\partial E_k}{\partial w_{ji}}$$

where for every  $k = 1, \dots, p$  holds

$$\frac{\partial E_k}{\partial w_{ji}} = \frac{\partial E_k}{\partial y_j} \cdot \sigma'_j(\xi_j) \cdot y_i$$

and for every  $j \in Z \setminus X$  we get

$$\frac{\partial E_k}{\partial y_j} = y_j - d_{kj} \quad \text{for } j \in Y$$

# MLP – error function gradient

For every  $w_{ji}$  we have

$$\frac{\partial E}{\partial w_{ji}} = \sum_{k=1}^p \frac{\partial E_k}{\partial w_{ji}}$$

where for every  $k = 1, \dots, p$  holds

$$\frac{\partial E_k}{\partial w_{ji}} = \frac{\partial E_k}{\partial y_j} \cdot \sigma'_j(\xi_j) \cdot y_i$$

and for every  $j \in Z \setminus X$  we get

$$\frac{\partial E_k}{\partial y_j} = y_j - d_{kj} \quad \text{for } j \in Y$$

$$\frac{\partial E_k}{\partial y_j} = \sum_{r \in \vec{J}} \frac{\partial E_k}{\partial y_r} \cdot \sigma'_r(\xi_r) \cdot w_{rj} \quad \text{for } j \in Z \setminus (Y \cup X)$$

(Here all  $y_j$  are in fact  $y_j(\vec{w}, \vec{x}_k)$ ).

## MLP – error function gradient (history)

- ▶ If  $y_j = \sigma_j(\xi_j) = \frac{1}{1+e^{-\xi_j}}$  for all  $j \in Z$ , then

$$\sigma'_j(\xi_j) = y_j(1 - y_j)$$

# MLP – error function gradient (history)

- ▶ If  $y_j = \sigma_j(\xi_j) = \frac{1}{1+e^{-\xi_j}}$  for all  $j \in Z$ , then

$$\sigma'_j(\xi_j) = y_j(1 - y_j)$$

and thus for all  $j \in Z \setminus X$ :

$$\frac{\partial E_k}{\partial y_j} = y_j - d_{kj} \quad \text{for } j \in Y$$

$$\frac{\partial E_k}{\partial y_j} = \sum_{r \in J} \frac{\partial E_k}{\partial y_r} \cdot y_r(1 - y_r) \cdot w_{rj} \quad \text{for } j \in Z \setminus (Y \cup X)$$

# MLP – computing the gradient

Compute  $\frac{\partial E}{\partial w_{ji}} = \sum_{k=1}^p \frac{\partial E_k}{\partial w_{ji}}$  as follows:



# MLP – computing the gradient

Compute  $\frac{\partial E}{\partial w_{ji}} = \sum_{k=1}^p \frac{\partial E_k}{\partial w_{ji}}$  as follows:

Initialize  $\mathcal{E}_{ji} := 0$

(By the end of the computation:  $\mathcal{E}_{ji} = \frac{\partial E}{\partial w_{ji}}$ )

# MLP – computing the gradient

Compute  $\frac{\partial E}{\partial w_{ji}} = \sum_{k=1}^p \frac{\partial E_k}{\partial w_{ji}}$  as follows:

Initialize  $\mathcal{E}_{ji} := 0$

(By the end of the computation:  $\mathcal{E}_{ji} = \frac{\partial E}{\partial w_{ji}}$ )

For every  $k = 1, \dots, p$  do:

# MLP – computing the gradient

Compute  $\frac{\partial E}{\partial w_{ji}} = \sum_{k=1}^p \frac{\partial E_k}{\partial w_{ji}}$  as follows:

Initialize  $\mathcal{E}_{ji} := 0$

(By the end of the computation:  $\mathcal{E}_{ji} = \frac{\partial E}{\partial w_{ji}}$ )

For every  $k = 1, \dots, p$  do:

1. **forward pass:** compute  $y_j = y_j(\vec{w}, \vec{x}_k)$  for all  $j \in Z$

# MLP – computing the gradient

Compute  $\frac{\partial E}{\partial w_{ji}} = \sum_{k=1}^p \frac{\partial E_k}{\partial w_{ji}}$  as follows:

Initialize  $\mathcal{E}_{ji} := 0$

(By the end of the computation:  $\mathcal{E}_{ji} = \frac{\partial E}{\partial w_{ji}}$ )

For every  $k = 1, \dots, p$  do:

1. **forward pass:** compute  $y_j = y_j(\vec{w}, \vec{x}_k)$  for all  $j \in Z$
2. **backward pass:** compute  $\frac{\partial E_k}{\partial y_j}$  for all  $j \in Z$  using *backpropagation* (see the next slide!)

# MLP – computing the gradient

Compute  $\frac{\partial E}{\partial w_{ji}} = \sum_{k=1}^p \frac{\partial E_k}{\partial w_{ji}}$  as follows:

Initialize  $\mathcal{E}_{ji} := 0$

(By the end of the computation:  $\mathcal{E}_{ji} = \frac{\partial E}{\partial w_{ji}}$ )

For every  $k = 1, \dots, p$  do:

1. **forward pass:** compute  $y_j = y_j(\vec{w}, \vec{x}_k)$  for all  $j \in Z$
2. **backward pass:** compute  $\frac{\partial E_k}{\partial y_j}$  for all  $j \in Z$  using *backpropagation* (see the next slide!)
3. compute  $\frac{\partial E_k}{\partial w_{ji}}$  for all  $w_{ji}$  using

$$\frac{\partial E_k}{\partial w_{ji}} := \frac{\partial E_k}{\partial y_j} \cdot \sigma'_j(\xi_j) \cdot y_i$$

# MLP – computing the gradient

Compute  $\frac{\partial E}{\partial w_{ji}} = \sum_{k=1}^p \frac{\partial E_k}{\partial w_{ji}}$  as follows:

Initialize  $\mathcal{E}_{ji} := 0$

(By the end of the computation:  $\mathcal{E}_{ji} = \frac{\partial E}{\partial w_{ji}}$ )

For every  $k = 1, \dots, p$  do:

1. **forward pass:** compute  $y_j = y_j(\vec{w}, \vec{x}_k)$  for all  $j \in Z$
2. **backward pass:** compute  $\frac{\partial E_k}{\partial y_j}$  for all  $j \in Z$  using *backpropagation* (see the next slide!)
3. compute  $\frac{\partial E_k}{\partial w_{ji}}$  for all  $w_{ji}$  using

$$\frac{\partial E_k}{\partial w_{ji}} := \frac{\partial E_k}{\partial y_j} \cdot \sigma'_j(\xi_j) \cdot y_i$$

4.  $\mathcal{E}_{ji} := \mathcal{E}_{ji} + \frac{\partial E_k}{\partial w_{ji}}$

The resulting  $\mathcal{E}_{ji}$  equals  $\frac{\partial E}{\partial w_{ji}}$ .

# MLP – backpropagation

Compute  $\frac{\partial E_k}{\partial y_j}$  for all  $j \in Z$  as follows:

# MLP – backpropagation

Compute  $\frac{\partial E_k}{\partial y_j}$  for all  $j \in Z$  as follows:

- ▶ if  $j \in Y$ , then  $\frac{\partial E_k}{\partial y_j} = y_j - d_{kj}$



# MLP – backpropagation

Compute  $\frac{\partial E_k}{\partial y_j}$  for all  $j \in Z$  as follows:

- ▶ if  $j \in Y$ , then  $\frac{\partial E_k}{\partial y_j} = y_j - d_{kj}$
- ▶ if  $j \in Z \setminus Y \cup X$ , then assuming that  $j$  is in the  $\ell$ -th layer and assuming that  $\frac{\partial E_k}{\partial y_r}$  has already been computed for all neurons in the  $\ell + 1$ -st layer, compute

$$\frac{\partial E_k}{\partial y_j} = \sum_{r \in j^{\rightarrow}} \frac{\partial E_k}{\partial y_r} \cdot \sigma'_r(\xi_r) \cdot w_{rj}$$

(This works because all neurons of  $r \in j^{\rightarrow}$  belong to the  $\ell + 1$ -st layer.)

# Complexity of the batch algorithm

Computation of  $\frac{\partial E}{\partial \mathbf{w}_{ji}}(\vec{\mathbf{w}}^{(t-1)})$  stops in time linear in the size of the network plus the size of the training set.

(assuming unit cost of operations including computation of  $\sigma'_r(\xi_r)$  for given  $\xi_r$ )

# Complexity of the batch algorithm

Computation of  $\frac{\partial E}{\partial w_{ji}}(\vec{w}^{(t-1)})$  stops in time linear in the size of the network plus the size of the training set.

(assuming unit cost of operations including computation of  $\sigma'_r(\xi_r)$  for given  $\xi_r$ )

**Proof sketch:** The algorithm does the following  $p$  times:

# Complexity of the batch algorithm

Computation of  $\frac{\partial E}{\partial w_{ji}}(\vec{w}^{(t-1)})$  stops in time linear in the size of the network plus the size of the training set.

(assuming unit cost of operations including computation of  $\sigma'_r(\xi_r)$  for given  $\xi_r$ )

**Proof sketch:** The algorithm does the following  $p$  times:

1. forward pass, i.e. computes  $y_j(\vec{w}, \vec{x}_k)$

# Complexity of the batch algorithm

Computation of  $\frac{\partial E}{\partial w_{ji}}(\vec{w}^{(t-1)})$  stops in time linear in the size of the network plus the size of the training set.

(assuming unit cost of operations including computation of  $\sigma'_r(\xi_r)$  for given  $\xi_r$ )

**Proof sketch:** The algorithm does the following  $p$  times:

1. forward pass, i.e. computes  $y_j(\vec{w}, \vec{x}_k)$
2. backpropagation, i.e. computes  $\frac{\partial E_k}{\partial y_j}$

# Complexity of the batch algorithm

Computation of  $\frac{\partial E}{\partial w_{ji}}(\vec{w}^{(t-1)})$  stops in time linear in the size of the network plus the size of the training set.

(assuming unit cost of operations including computation of  $\sigma'_r(\xi_r)$  for given  $\xi_r$ )

**Proof sketch:** The algorithm does the following  $p$  times:

1. forward pass, i.e. computes  $y_j(\vec{w}, \vec{x}_k)$
2. backpropagation, i.e. computes  $\frac{\partial E_k}{\partial y_j}$
3. computes  $\frac{\partial E_k}{\partial w_{ji}}$  and adds it to  $\mathcal{E}_{ji}$  (a constant time operation in the unit cost framework)

# Complexity of the batch algorithm

Computation of  $\frac{\partial E}{\partial w_{ji}}(\vec{w}^{(t-1)})$  stops in time linear in the size of the network plus the size of the training set.

(assuming unit cost of operations including computation of  $\sigma'_r(\xi_r)$  for given  $\xi_r$ )

**Proof sketch:** The algorithm does the following  $p$  times:

1. forward pass, i.e. computes  $y_j(\vec{w}, \vec{x}_k)$
2. backpropagation, i.e. computes  $\frac{\partial E_k}{\partial y_j}$
3. computes  $\frac{\partial E_k}{\partial w_{ji}}$  and adds it to  $\mathcal{E}_{ji}$  (a constant time operation in the unit cost framework)

The steps 1. - 3. take linear time.

# Complexity of the batch algorithm

Computation of  $\frac{\partial E}{\partial w_{ji}}(\vec{w}^{(t-1)})$  stops in time linear in the size of the network plus the size of the training set.

(assuming unit cost of operations including computation of  $\sigma'_r(\xi_r)$  for given  $\xi_r$ )

**Proof sketch:** The algorithm does the following  $p$  times:

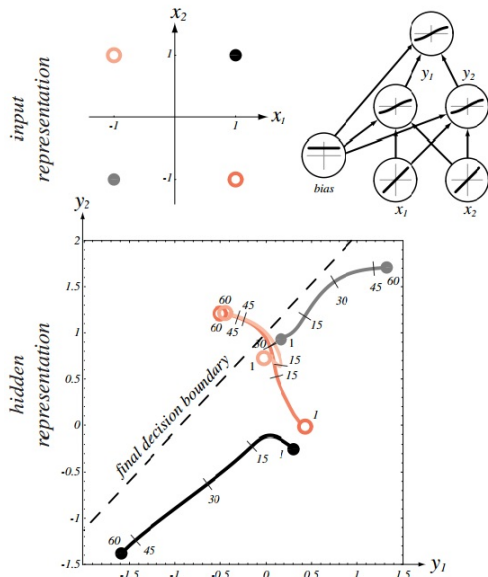
1. forward pass, i.e. computes  $y_j(\vec{w}, \vec{x}_k)$
2. backpropagation, i.e. computes  $\frac{\partial E_k}{\partial y_j}$
3. computes  $\frac{\partial E_k}{\partial w_{ji}}$  and adds it to  $\mathcal{E}_{ji}$  (a constant time operation in the unit cost framework)

The steps 1. - 3. take linear time.

Note that the speed of convergence of the gradient descent cannot be estimated ...



# Illustration of the gradient descent – XOR



Source: Pattern Classification (2nd Edition); Richard O. Duda, Peter E. Hart, David G. Stork

# MLP – learning algorithm

## Online algorithm:

The algorithm computes a sequence of weight vectors

$\vec{w}^{(0)}, \vec{w}^{(1)}, \vec{w}^{(2)}, \dots$

- ▶ weights in  $\vec{w}^{(0)}$  are randomly initialized to values close to 0
- ▶ in the step  $t + 1$  (here  $t = 0, 1, 2 \dots$ ), weights  $\vec{w}^{(t+1)}$  are computed as follows:

$$w_{ji}^{(t+1)} = w_{ji}^{(t)} + \Delta w_{ji}^{(t)}$$

where

$$\Delta w_{ji}^{(t)} = -\varepsilon(t) \cdot \frac{\partial E_k}{\partial w_{ji}}(w_{ji}^{(t)})$$

is the *weight update* of  $w_{ji}$  in the step  $t + 1$  and  $0 < \varepsilon(t) \leq 1$   
is the *learning rate* in the step  $t + 1$ .

There are other variants determined by selection of the training examples used for the error computation (more on this later).

- ▶ weights in  $\vec{w}^{(0)}$  are randomly initialized to values close to 0
- ▶ in the step  $t + 1$  (here  $t = 0, 1, 2 \dots$ ), weights  $\vec{w}^{(t+1)}$  are computed as follows:
  - ▶ Choose (randomly) a set of training examples  $T \subseteq \{1, \dots, p\}$
  - ▶ Compute

$$\vec{w}^{(t+1)} = \vec{w}^{(t)} + \Delta \vec{w}^{(t)}$$

where

$$\Delta \vec{w}^{(t)} = -\varepsilon(t) \cdot \sum_{k \in T} \nabla E_k(\vec{w}^{(t)})$$

- ▶  $0 < \varepsilon(t) \leq 1$  is a *learning rate* in step  $t + 1$
- ▶  $\nabla E_k(\vec{w}^{(t)})$  is the gradient of the error of the example  $k$

Note that the random choice of the minibatch is typically implemented by randomly shuffling all data and then choosing minibatches sequentially.

# Output activations and error functions

## Regression:

- ▶ The output activation is typically the identity  $y_i = \sigma(\xi_i) = \xi_i$ .

# Output activations and error functions

## Regression:

- ▶ The output activation is typically the identity  $y_i = \sigma(\xi_i) = \xi_i$ .
- ▶ A **training dataset**

$$\left\{ \left( \vec{x}_k, \vec{d}_k \right) \mid k = 1, \dots, p \right\}$$

Here, every  $\vec{x}_k \in \mathbb{R}^{|X|}$  is an *input vector* and every  $\vec{d}_k \in \mathbb{R}^{|Y|}$  is the desired network output. For every  $i \in Y$ , denote by  $d_{ki}$  the desired output of the neuron  $i$  for a given network input  $\vec{x}_k$  (the vector  $\vec{d}_k$  can be written as  $(d_{ki})_{i \in Y}$ ).

# Output activations and error functions

## Regression:

- ▶ The output activation is typically the identity  $y_i = \sigma(\xi_i) = \xi_i$ .
- ▶ A **training dataset**

$$\left\{ \left( \vec{x}_k, \vec{d}_k \right) \mid k = 1, \dots, p \right\}$$

Here, every  $\vec{x}_k \in \mathbb{R}^{|X|}$  is an *input vector* and every  $\vec{d}_k \in \mathbb{R}^{|Y|}$  is the desired network output. For every  $i \in Y$ , denote by  $d_{ki}$  the desired output of the neuron  $i$  for a given network input  $\vec{x}_k$  (the vector  $\vec{d}_k$  can be written as  $(d_{ki})_{i \in Y}$ ).

- ▶ The error function *mean squared error (mse)*:

$$E(\vec{w}) = \frac{1}{p} \sum_{k=1}^p E_k(\vec{w})$$

where

$$E_k(\vec{w}) = \frac{1}{2} \sum_{i \in Y} \left( y_i(\vec{w}, \vec{x}_k) - d_{ki} \right)^2$$

# Output activations and error functions

## Classification

- The output activation function *softmax*:

$$y_i = \sigma_i(\xi_i) = \frac{e^{\xi_i}}{\sum_{j \in Y} e^{\xi_j}}$$

# Output activations and error functions

## Classification

- ▶ The output activation function *softmax*:

$$y_i = \sigma_i(\xi_i) = \frac{e^{\xi_i}}{\sum_{j \in Y} e^{\xi_j}}$$

- ▶ **A training dataset**

$$\left\{ \left( \vec{x}_k, \vec{d}_k \right) \mid k = 1, \dots, p \right\}$$

Here, every  $\vec{x}_k \in \mathbb{R}^{|X|}$  is an *input vector* and every  $\vec{d}_k \in \{0, 1\}^{|Y|}$  is the desired network output. For every  $i \in Y$ , denote by  $d_{ki}$  the desired output of the neuron  $i$  for a given network input  $\vec{x}_k$  (the vector  $\vec{d}_k$  can be written as  $(d_{ki})_{i \in Y}$ ).



# Output activations and error functions

## Classification

- ▶ The output activation function *softmax*:

$$y_i = \sigma_i(\xi_i) = \frac{e^{\xi_i}}{\sum_{j \in Y} e^{\xi_j}}$$

- ▶ A training dataset

$$\left\{ \left( \vec{x}_k, \vec{d}_k \right) \mid k = 1, \dots, p \right\}$$

Here, every  $\vec{x}_k \in \mathbb{R}^{|X|}$  is an *input vector* and every  $\vec{d}_k \in \{0, 1\}^{|Y|}$  is the desired network output. For every  $i \in Y$ , denote by  $d_{ki}$  the desired output of the neuron  $i$  for a given network input  $\vec{x}_k$  (the vector  $\vec{d}_k$  can be written as  $(d_{ki})_{i \in Y}$ ).

- ▶ The error function (*categorical*) *cross entropy*:

$$E(\vec{w}) = -\frac{1}{p} \sum_{k=1}^p \sum_{i \in Y} d_{ki} \log(y_i(\vec{w}, \vec{x}_k))$$

# Gradient with Softmax & Cross-Entropy

Assume that  $V$  is the layer just below the output layer  $Y$ .

$$\begin{aligned}E(\vec{w}) &= -\frac{1}{p} \sum_{k=1}^p \sum_{i \in Y} d_{ki} \log(y_i(\vec{w}, \vec{x}_k)) \\&= -\frac{1}{p} \sum_{k=1}^p \sum_{i \in Y} d_{ki} \log\left(\frac{e^{\xi_i}}{\sum_{j \in Y} e^{\xi_j}}\right) \\&= -\frac{1}{p} \sum_{k=1}^p \sum_{i \in Y} d_{ki} \left( \xi_i - \log\left(\sum_{j \in Y} e^{\xi_j}\right) \right) \\&= -\frac{1}{p} \sum_{k=1}^p \sum_{i \in Y} d_{ki} \left( \sum_{\ell \in V} w_{i\ell} y_{\ell} - \log\left(\sum_{j \in Y} e^{\sum_{\ell \in V} w_{j\ell} y_{\ell}}\right) \right)\end{aligned}$$

Now compute the derivatives  $\frac{\delta E}{\delta y_{\ell}}$  for  $\ell \in V$ .

# Output activations and error functions

## Binary classification

Assume a single output neuron  $o \in Y = \{o\}$ .

- ▶ The output activation function *logistic sigmoid*:

$$\sigma_o(\xi_o) = \frac{e^{\xi_o}}{e^{\xi_o} + 1} = \frac{1}{1 + e^{-\xi_o}}$$

# Output activations and error functions

## Binary classification

Assume a single output neuron  $o \in Y = \{o\}$ .

- ▶ The output activation function *logistic sigmoid*:

$$\sigma_o(\xi_o) = \frac{e^{\xi_o}}{e^{\xi_o} + 1} = \frac{1}{1 + e^{-\xi_o}}$$

- ▶ **A training dataset**

$$\mathcal{T} = \{(\vec{x}_1, d_1), (\vec{x}_2, d_2), \dots, (\vec{x}_p, d_p)\}$$

Here  $\vec{x}_k = (x_{k0}, x_{k1}, \dots, x_{kn}) \in \mathbb{R}^{n+1}$ ,  $x_{k0} = 1$ , is the  $k$ -th input, and  $d_k \in \{0, 1\}$  is the desired output.

# Output activations and error functions

## Binary classification

Assume a single output neuron  $o \in Y = \{o\}$ .

- ▶ The output activation function *logistic sigmoid*:

$$\sigma_o(\xi_o) = \frac{e^{\xi_o}}{e^{\xi_o} + 1} = \frac{1}{1 + e^{-\xi_o}}$$

- ▶ A training dataset

$$\mathcal{T} = \{(\vec{x}_1, d_1), (\vec{x}_2, d_2), \dots, (\vec{x}_p, d_p)\}$$

Here  $\vec{x}_k = (x_{k0}, x_{k1}, \dots, x_{kn}) \in \mathbb{R}^{n+1}$ ,  $x_{k0} = 1$ , is the  $k$ -th input, and  $d_k \in \{0, 1\}$  is the desired output.

- ▶ The error function (*Binary*) *cross-entropy*:

$$E(\vec{w}) = \sum_{k=1}^p -(d_k \log(y_o(\vec{w}, \vec{x}_k)) + (1 - d_k) \log(1 - y_o(\vec{w}, \vec{x}_k)))$$

## But what is the meaning of the sigmoid?

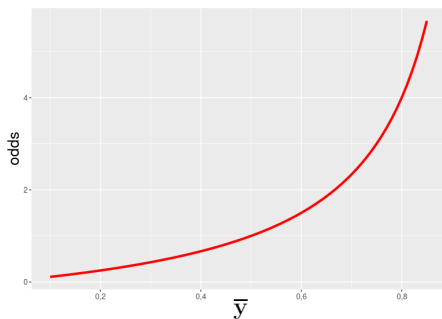
The model gives a probability  $y_o$  of the class 1 given an input  $\vec{x}$ .  
But why do we model such a probability using  $1/(1 + e^{-\xi_o})$  ?

## But what is the meaning of the sigmoid?

The model gives a probability  $y_o$  of the class 1 given an input  $\vec{x}$ .  
But why do we model such a probability using  $1/(1 + e^{-\xi_o})$  ?

Let  $\bar{y}$  be the "true" probability of the class 1 to be modeled.  
What about **odds** of the class 1?

$$\text{odds}(\bar{y}) = \bar{y}/1 - \bar{y}$$



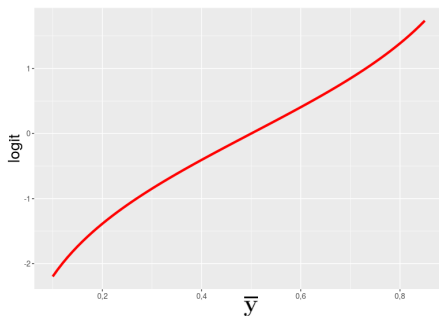
... stretches from 0 to  $\infty$

## But what is the meaning of the sigmoid?

The model gives a probability  $y_o$  of the class 1 given an input  $\vec{x}$ .  
But why do we model such a probability using  $1/(1 + e^{-\xi_o})$  ?

Let  $\bar{y}$  be the "true" probability of the class 1 to be modeled.  
What about **log odds (aka logit)** of the class 1?

$$\text{logit}(\bar{y}) = \log(\bar{y}/(1 - \bar{y}))$$



... stretches from  $-\infty$  to  $\infty$



## But what is the meaning of the sigmoid?

Assume that  $\bar{y}$  is the probability of the class 1. Put

$$\log(\bar{y}/(1 - \bar{y})) = \xi_o$$

(here  $\xi_o$  is the inner potential of the output neuron).

## But what is the meaning of the sigmoid?

Assume that  $\bar{y}$  is the probability of the class 1. Put

$$\log(\bar{y}/(1 - \bar{y})) = \xi_o$$

(here  $\xi_o$  is the inner potential of the output neuron). Then

$$\log((1 - \bar{y})/\bar{y}) = -\xi_o$$

## But what is the meaning of the sigmoid?

Assume that  $\bar{y}$  is the probability of the class 1. Put

$$\log(\bar{y}/(1 - \bar{y})) = \xi_o$$

(here  $\xi_o$  is the inner potential of the output neuron). Then

$$\log((1 - \bar{y})/\bar{y}) = -\xi_o$$

and

$$(1 - \bar{y})/\bar{y} = e^{-\xi_o}$$

## But what is the meaning of the sigmoid?

Assume that  $\bar{y}$  is the probability of the class 1. Put

$$\log(\bar{y}/(1 - \bar{y})) = \xi_o$$

(here  $\xi_o$  is the inner potential of the output neuron). Then

$$\log((1 - \bar{y})/\bar{y}) = -\xi_o$$

and

$$(1 - \bar{y})/\bar{y} = e^{-\xi_o}$$

and

$$\bar{y} = \frac{1}{1 + e^{-\xi_o}}$$

That is, modeling the probability using the *classification model* (with the logistic output activation) corresponds to modeling log-odds using the *regression model* (with the identity output activation).

# Log likelihood is your friend!

What is the statistical meaning of the cross-entropy?

- ▶ Let's have a "coin" (sides 0 and 1).

# Log likelihood is your friend!

What is the statistical meaning of the cross-entropy?

- ▶ Let's have a "coin" (sides 0 and 1).
- ▶ The probability of 1 is  $\bar{y}$  and is unknown!

# Log likelihood is your friend!

What is the statistical meaning of the cross-entropy?

- ▶ Let's have a "coin" (sides 0 and 1).
- ▶ The probability of 1 is  $\bar{y}$  and is unknown!
- ▶ You have tossed the coin 5 times and got a training dataset:

$$\mathcal{T} = \{1, 1, 0, 0, 1\} = \{d_1, \dots, d_5\}$$

Consider this to be a very special case where the input dimension is 0

# Log likelihood is your friend!

What is the statistical meaning of the cross-entropy?

- ▶ Let's have a "coin" (sides 0 and 1).
- ▶ The probability of 1 is  $\bar{y}$  and is unknown!
- ▶ You have tossed the coin 5 times and got a training dataset:

$$\mathcal{T} = \{1, 1, 0, 0, 1\} = \{d_1, \dots, d_5\}$$

Consider this to be a very special case where the input dimension is 0

- ▶ What is the best model  $y$  of  $\bar{y}$  based on the data?



# Log likelihood is your friend!

What is the statistical meaning of the cross-entropy?

- ▶ Let's have a "coin" (sides 0 and 1).
- ▶ The probability of 1 is  $\bar{y}$  and is unknown!
- ▶ You have tossed the coin 5 times and got a training dataset:

$$\mathcal{T} = \{1, 1, 0, 0, 1\} = \{d_1, \dots, d_5\}$$

Consider this to be a very special case where the input dimension is 0

- ▶ What is the best model  $y$  of  $\bar{y}$  based on the data?

**Answer:** The one that generates the data with maximum probability!

# Log likelihood is your friend!

Keep in mind our dataset:

$$\mathcal{T} = \{1, 1, 0, 0, 1\} = \{d_1, \dots, d_5\}$$

# Log likelihood is your friend!

Keep in mind our dataset:

$$\mathcal{T} = \{1, 1, 0, 0, 1\} = \{d_1, \dots, d_5\}$$

Assume that the data was generated by independent trials, then the probability of getting exactly  $\mathcal{T}$  from our model is

$$L = y \cdot y \cdot (1 - y) \cdot (1 - y) \cdot y$$

How to maximize this w.r.t.  $y$ ?

# Log likelihood is your friend!

Keep in mind our dataset:

$$\mathcal{T} = \{1, 1, 0, 0, 1\} = \{d_1, \dots, d_5\}$$

Assume that the data was generated by independent trials, then the probability of getting exactly  $\mathcal{T}$  from our model is

$$L = y \cdot y \cdot (1 - y) \cdot (1 - y) \cdot y$$

How to maximize this w.r.t.  $y$ ?

Maximize

$$LL = \log(L) = \log(y) + \log(y) + \log(1 - y) + \log(1 - y) + \log(y)$$

# Log likelihood is your friend!

Keep in mind our dataset:

$$\mathcal{T} = \{1, 1, 0, 0, 1\} = \{d_1, \dots, d_5\}$$

Assume that the data was generated by independent trials, then the probability of getting exactly  $\mathcal{T}$  from our model is

$$L = y \cdot y \cdot (1 - y) \cdot (1 - y) \cdot y$$

How to maximize this w.r.t.  $y$ ?

Maximize

$$LL = \log(L) = \log(y) + \log(y) + \log(1 - y) + \log(1 - y) + \log(y)$$

But then

$$-LL = -1 \cdot \log(y) - 1 \cdot \log(y) - (1 - 0) \cdot \log(1 - y) - (1 - 0) \cdot \log(1 - y) - 1 \cdot \log(y)$$

i.e.  $-LL$  is the cross-entropy.

## Let the coin depend on the input

Consider our model giving a probability  $y_o(\vec{w}, \vec{x})$  given input  $\vec{x}$ .

## Let the coin depend on the input

Consider our model giving a probability  $y_o(\vec{w}, \vec{x})$  given input  $\vec{x}$ . Recall that the training dataset is

$$\mathcal{T} = \{(\vec{x}_1, d_1), (\vec{x}_2, d_2), \dots, (\vec{x}_p, d_p)\}$$

Here  $\vec{x}_k = (x_{k0}, x_{k1}, \dots, x_{kn}) \in \mathbb{R}^{n+1}$ ,  $x_{k0} = 1$ , is the  $k$ -th input, and  $d_k \in \{0, 1\}$  is the expected output.

## Let the coin depend on the input

Consider our model giving a probability  $y_o(\vec{w}, \vec{x})$  given input  $\vec{x}$ . Recall that the training dataset is

$$\mathcal{T} = \{(\vec{x}_1, d_1), (\vec{x}_2, d_2), \dots, (\vec{x}_p, d_p)\}$$

Here  $\vec{x}_k = (x_{k0}, x_{k1}, \dots, x_{kn}) \in \mathbb{R}^{n+1}$ ,  $x_{k0} = 1$ , is the  $k$ -th input, and  $d_k \in \{0, 1\}$  is the expected output.

The *likelihood*:

$$L(\vec{w}) = \prod_{k=1}^p \left(y_o(\vec{w}, \vec{x}_k)\right)^{d_k} \cdot \left(1 - y_o(\vec{w}, \vec{x}_k)\right)^{(1-d_k)}$$

$$\log(L) =$$

$$\sum_{k=1}^p \left(d_k \cdot \log(y_o(\vec{w}, \vec{x}_k)) + (1 - d_k) \cdot \log(1 - y_o(\vec{w}, \vec{x}_k))\right)$$

and thus  $-\log(L)$  = the cross-entropy.

**Minimizing the cross-entropy maximizes the log-likelihood (and vice versa).**



## Squared Error vs Logistic Output Activation

Consider a single neuron model  $y = \sigma(w \cdot x) = 1/(1 + e^{-w \cdot x})$  where  $w \in \mathbb{R}$  is the weight (ignore the bias).

A training dataset  $\mathcal{T} = \{(x, d)\}$  where  $x \in \mathbb{R}$  and  $d \in \{0, 1\}$ .

## Squared Error vs Logistic Output Activation

Consider a single neuron model  $y = \sigma(w \cdot x) = 1/(1 + e^{-w \cdot x})$  where  $w \in \mathbb{R}$  is the weight (ignore the bias).

A training dataset  $\mathcal{T} = \{(x, d)\}$  where  $x \in \mathbb{R}$  and  $d \in \{0, 1\}$ .

Squared error  $E(w) = \frac{1}{2}(y - d)^2$ .

$$\frac{\delta E}{\delta w} = (y - d) \cdot y \cdot (1 - y) \cdot x$$

# Squared Error vs Logistic Output Activation

Consider a single neuron model  $y = \sigma(w \cdot x) = 1/(1 + e^{-w \cdot x})$  where  $w \in \mathbb{R}$  is the weight (ignore the bias).

A training dataset  $\mathcal{T} = \{(x, d)\}$  where  $x \in \mathbb{R}$  and  $d \in \{0, 1\}$ .

Squared error  $E(w) = \frac{1}{2}(y - d)^2$ .

$$\frac{\delta E}{\delta w} = (y - d) \cdot y \cdot (1 - y) \cdot x$$

Thus

- ▶ If  $d = 1$  and  $y \approx 0$ , then  $\frac{\delta E}{\delta w} \approx 0$
- ▶ If  $d = 0$  and  $y \approx 1$ , then  $\frac{\delta E}{\delta w} \approx 0$

The gradient of  $E$  is small even though *the model is wrong!*

## Squared Error vs Logistic Output Activation

Consider a single neuron model  $y = \sigma(w \cdot x) = 1/(1 + e^{-w \cdot x})$  where  $w \in \mathbb{R}$  is the weight (ignore the bias).

A training dataset  $\mathcal{T} = \{(x, d)\}$  where  $x \in \mathbb{R}$  and  $d \in \{0, 1\}$ .

Cross-entropy error  $E(w) = -d \cdot \log(y) - (1 - d) \cdot \log(1 - y)$ .

## Squared Error vs Logistic Output Activation

Consider a single neuron model  $y = \sigma(w \cdot x) = 1/(1 + e^{-w \cdot x})$  where  $w \in \mathbb{R}$  is the weight (ignore the bias).

A training dataset  $\mathcal{T} = \{(x, d)\}$  where  $x \in \mathbb{R}$  and  $d \in \{0, 1\}$ .

Cross-entropy error  $E(w) = -d \cdot \log(y) - (1 - d) \cdot \log(1 - y)$ .

For  $d = 1$

$$\frac{\delta E}{\delta w} = -\frac{1}{y} \cdot y \cdot (1 - y) \cdot x = -(1 - y) \cdot x$$

which is close to  $-x$  for  $y \approx 0$ .

# Squared Error vs Logistic Output Activation

Consider a single neuron model  $y = \sigma(w \cdot x) = 1/(1 + e^{-w \cdot x})$  where  $w \in \mathbb{R}$  is the weight (ignore the bias).

A training dataset  $\mathcal{T} = \{(x, d)\}$  where  $x \in \mathbb{R}$  and  $d \in \{0, 1\}$ .

Cross-entropy error  $E(w) = -d \cdot \log(y) - (1 - d) \cdot \log(1 - y)$ .

For  $d = 1$

$$\frac{\delta E}{\delta w} = -\frac{1}{y} \cdot y \cdot (1 - y) \cdot x = -(1 - y) \cdot x$$

which is close to  $-x$  for  $y \approx 0$ .

For  $d = 0$

$$\frac{\delta E}{\delta w} = -\frac{1}{1 - y} \cdot (-y) \cdot (1 - y) \cdot x = y \cdot x$$

which is close to  $x$  for  $y \approx 1$ .

## MLP training – practical issues

# Practical issues of gradient descent

- ▶ Training efficiency:
  - ▶ What size of a minibatch?
  - ▶ How to choose the learning rate  $\varepsilon(t)$  and control SGD ?
  - ▶ How to pre-process the inputs?
  - ▶ How to initialize weights?
  - ▶ How to choose desired output values of the network?



# Practical issues of gradient descent

- ▶ Training efficiency:
  - ▶ What size of a minibatch?
  - ▶ How to choose the learning rate  $\varepsilon(t)$  and control SGD ?
  - ▶ How to pre-process the inputs?
  - ▶ How to initialize weights?
  - ▶ How to choose desired output values of the network?
- ▶ Quality of the resulting model:
  - ▶ When to stop training?
  - ▶ Regularization techniques.
  - ▶ How large network?

For simplicity, I will illustrate the reasoning on MLP + mse. Later we will see other topologies and error functions with different but always somewhat related issues.

# Issues in gradient descent

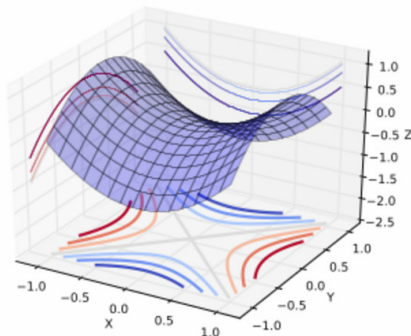
- ▶ Small networks: Lots of local minima where the descent gets stuck.
- ▶ The model identifiability problem: Swapping incoming weights of neurons  $i$  and  $j$  leaves the same network topology – **weight space symmetry**.
- ▶ Recent studies show that for sufficiently large networks all local minima have low values of the error function.

# Issues in gradient descent

- ▶ Small networks: Lots of local minima where the descent gets stuck.
- ▶ The model identifiability problem: Swapping incoming weights of neurons  $i$  and  $j$  leaves the same network topology – **weight space symmetry**.
- ▶ Recent studies show that for sufficiently large networks all local minima have low values of the error function.

## Saddle points

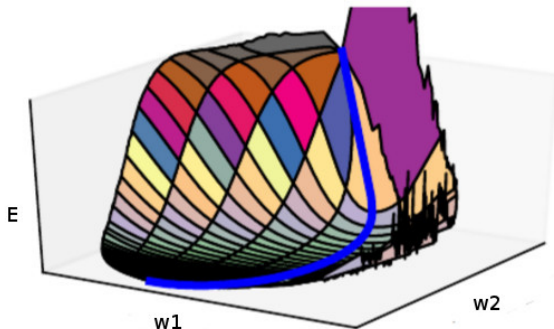
One can show (by a combinatorial argument) that larger networks have exponentially more saddle points than local minima.



# Issues in gradient descent – too slow descent

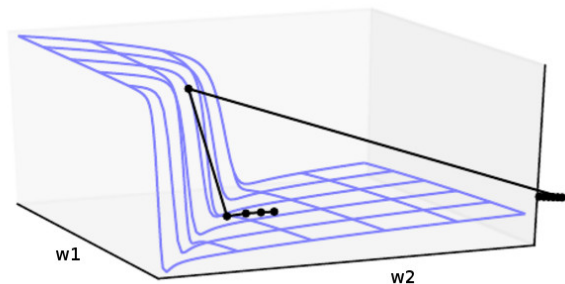
- flat regions

E.g. if the inner potentials are too large (in abs. value), then their derivative is extremely small.

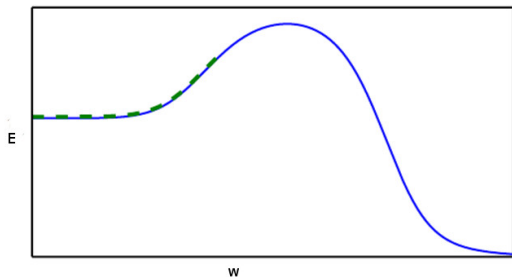


# Issues in gradient descent – too fast descent

- ▶ steep cliffs: the gradient is extremely large, descent skips important weight vectors



# Issues in gradient descent – local vs global structure



What if we initialize on the left?

# Gradient Descent in Large Networks

## Theorem

Assume (roughly),

- ▶ activation functions: "smooth" ReLU (softplus)

$$\sigma(z) = \log(1 + \exp(z))$$

*In general: Smooth, non-polynomial, analytic, Lipschitz continuous.*

- ▶ inputs  $\vec{x}_k$  of Euclidean norm equal to 1, desired values  $d_k$  such that all  $|d_k|$  are bounded by a constant,
- ▶ the number of hidden neurons per layer sufficiently large  
*(polynomial in certain numerical characteristics of inputs roughly measuring their similarity, and exponential in the depth of the network),*
- ▶ the learning rate constant and sufficiently small.

*The gradient descent converges (with high probability w.r.t. random initialization) to a global minimum with zero error at linear rate.*

*Later we get to a special type of networks called ResNet where the above result demands only polynomially many neurons per layer (w.r.t. depth).*

# Issues in computing the gradient

- ▶ vanishing and exploding gradients

$$\frac{\partial E_k}{\partial y_j} = y_j - d_{kj} \quad \text{for } j \in Y$$

$$\frac{\partial E_k}{\partial y_j} = \sum_{r \in j^{\rightarrow}} \frac{\partial E_k}{\partial y_r} \cdot \sigma'_r(\xi_r) \cdot w_{rj} \quad \text{for } j \in Z \setminus (Y \cup X)$$



# Issues in computing the gradient

- ▶ vanishing and exploding gradients

$$\frac{\partial E_k}{\partial y_j} = y_j - d_{kj} \quad \text{for } j \in Y$$

$$\frac{\partial E_k}{\partial y_j} = \sum_{r \in j \rightarrow} \frac{\partial E_k}{\partial y_r} \cdot \sigma'_r(\xi_r) \cdot w_{rj} \quad \text{for } j \in Z \setminus (Y \cup X)$$

- ▶ inexact gradient computation:
  - ▶ Minibatch gradient is only an estimate of the true gradient.
  - ▶ Note that the standard deviation of the estimate is (roughly)  $\sigma / \sqrt{m}$  where  $m$  is the size of the minibatch and  $\sigma$  is the variance of the gradient estimate for a single training example.  
(E.g. minibatch size 10 000 means 100 times more computation than the size 100 but gives only 10 times less deviation.)

## Minibatch size

- ▶ Larger batches provide a more accurate estimate of the gradient, but with less than linear returns.

# Minibatch size

- ▶ Larger batches provide a more accurate estimate of the gradient, but with less than linear returns.
- ▶ Multicore architectures are usually underutilized by extremely small batches.

# Minibatch size

- ▶ Larger batches provide a more accurate estimate of the gradient, but with less than linear returns.
- ▶ Multicore architectures are usually underutilized by extremely small batches.
- ▶ If all examples in the batch are to be processed in parallel (as is the typical case), then the amount of memory scales with the batch size. For many hardware setups this is the limiting factor in batch size.

# Minibatch size

- ▶ Larger batches provide a more accurate estimate of the gradient, but with less than linear returns.
- ▶ Multicore architectures are usually underutilized by extremely small batches.
- ▶ If all examples in the batch are to be processed in parallel (as is the typical case), then the amount of memory scales with the batch size. For many hardware setups this is the limiting factor in batch size.
- ▶ It is common (especially when using GPUs) for power of 2 batch sizes to offer better runtime. Typical power of 2 batch sizes range from 32 to 256, with 16 sometimes being attempted for large models.

# Minibatch size

- ▶ Larger batches provide a more accurate estimate of the gradient, but with less than linear returns.
- ▶ Multicore architectures are usually underutilized by extremely small batches.
- ▶ If all examples in the batch are to be processed in parallel (as is the typical case), then the amount of memory scales with the batch size. For many hardware setups this is the limiting factor in batch size.
- ▶ It is common (especially when using GPUs) for power of 2 batch sizes to offer better runtime. Typical power of 2 batch sizes range from 32 to 256, with 16 sometimes being attempted for large models.
- ▶ Small batches can offer a regularizing effect, perhaps due to the noise they add to the learning process.

It has been observed in practice that when using a larger batch there is a degradation in the quality of the model, as measured by its ability to generalize.

# Momentum

Issue in the gradient descent:

- ▶  $\nabla E(\vec{w}^{(t)})$  constantly changes direction (but the error steadily decreases).



# Momentum

Issue in the gradient descent:

- ▶  $\nabla E(\vec{w}^{(t)})$  constantly changes direction (but the error steadily decreases).



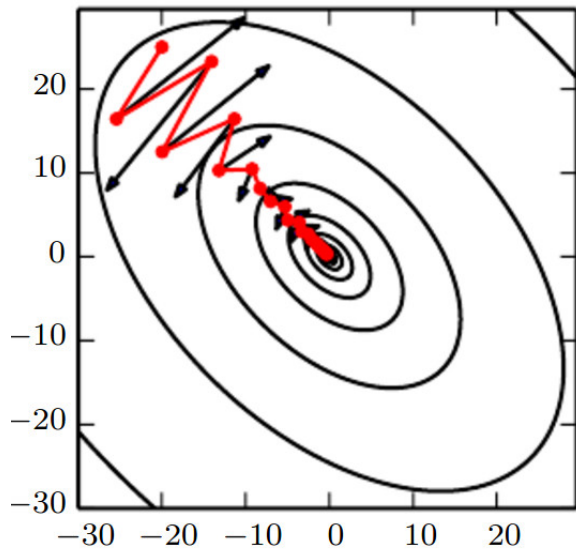
**Solution:** In every step add the change made in the previous step (weighted by a factor  $\alpha$ ):

$$\Delta \vec{w}^{(t)} = -\varepsilon(t) \cdot \sum_{k \in T} \nabla E_k(\vec{w}^{(t)}) + \alpha \cdot \Delta \vec{w}^{(t-1)}$$

where  $0 < \alpha < 1$ .



## Momentum – illustration



# SGD with momentum

- ▶ weights in  $\vec{w}^{(0)}$  are randomly initialized to values close to 0
- ▶ in the step  $t + 1$  (here  $t = 0, 1, 2 \dots$ ), weights  $\vec{w}^{(t+1)}$  are computed as follows:
  - ▶ Choose (randomly) a set of training examples  $T \subseteq \{1, \dots, p\}$
  - ▶ Compute

$$\vec{w}^{(t+1)} = \vec{w}^{(t)} + \Delta \vec{w}^{(t)}$$

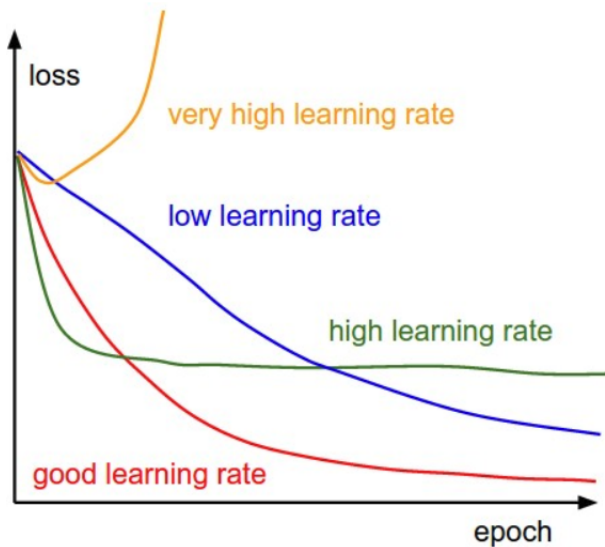
where

$$\Delta \vec{w}^{(t)} = -\varepsilon(t) \cdot \sum_{k \in T} \nabla E_k(\vec{w}^{(t)}) + \alpha \Delta \vec{w}^{(t-1)}$$

- ▶  $0 < \varepsilon(t) \leq 1$  is a *learning rate* in step  $t + 1$
- ▶  $0 < \alpha < 1$  measures the "influence" of the momentum
- ▶  $\nabla E_k(\vec{w}^{(t)})$  is the gradient of the error of the example  $k$

Note that the random choice of the minibatch is typically implemented by randomly shuffling all data and then choosing minibatches sequentially.

# Learning rate



# Search for the learning rate

- ▶ Use settings from a successful solution of a similar problem as a baseline.
- ▶ Search for the learning rate using the learning monitoring:
  - ▶ Search through values from small (e.g. 0.001) to (0.1), possibly multiplying by 2.
  - ▶ Train for several epochs, observe the learning curves (see cross-validation later).

# Adaptive learning rate

- ▶ Power scheduling: Set  $\epsilon(t) = \epsilon_0 / (1 + t/s)$  where  $\epsilon_0$  is an initial learning rate and  $s$  is a number  
(after  $s$  steps the learning rate is  $\epsilon_0/2$ , after  $2s$  it is  $\epsilon_0/3$  etc.)

# Adaptive learning rate

- ▶ Power scheduling: Set  $\epsilon(t) = \epsilon_0 / (1 + t/s)$  where  $\epsilon_0$  is an initial learning rate and  $s$  is a number  
(after  $s$  steps the learning rate is  $\epsilon_0/2$ , after  $2s$  it is  $\epsilon_0/3$  etc.)
- ▶ Exponential scheduling: Set  $\epsilon(t) = \epsilon_0 \cdot 0.1^{t/s}$ .  
(the learning rate decays faster than in the power scheduling)

# Adaptive learning rate

- ▶ Power scheduling: Set  $\epsilon(t) = \epsilon_0 / (1 + t/s)$  where  $\epsilon_0$  is an initial learning rate and  $s$  is a number  
(after  $s$  steps the learning rate is  $\epsilon_0/2$ , after  $2s$  it is  $\epsilon_0/3$  etc.)
- ▶ Exponential scheduling: Set  $\epsilon(t) = \epsilon_0 \cdot 0.1^{t/s}$ .  
(the learning rate decays faster than in the power scheduling)
- ▶ Piecewise constant scheduling: A constant learning rate for a number of steps/epochs, then a smaller learning rate, and so on.

# Adaptive learning rate

- ▶ Power scheduling: Set  $\epsilon(t) = \epsilon_0 / (1 + t/s)$  where  $\epsilon_0$  is an initial learning rate and  $s$  is a number  
(after  $s$  steps the learning rate is  $\epsilon_0/2$ , after  $2s$  it is  $\epsilon_0/3$  etc.)
- ▶ Exponential scheduling: Set  $\epsilon(t) = \epsilon_0 \cdot 0.1^{t/s}$ .  
(the learning rate decays faster than in the power scheduling)
- ▶ Piecewise constant scheduling: A constant learning rate for a number of steps/epochs, then a smaller learning rate, and so on.
- ▶ 1cycle scheduling: Start by increasing the initial learning rate from  $\epsilon_0$  linearly to  $\epsilon_1$  (approx.  $\epsilon_1 = 10\epsilon_0$ ) halfway through training. Then decrease from  $\epsilon_1$  linearly to  $\epsilon_0$ . Finish by dropping the learning rate by several orders of magnitude (still linearly).  
According to a 2018 paper by Leslie Smith this may converge much faster (100 epochs vs 800 epochs on CIFAR10 dataset).

For comparison of some methods see: AN EMPIRICAL STUDY OF LEARNING RATES IN DEEP NEURAL NETWORKS FOR SPEECH RECOGNITION, Senior et al



So far we have considered fixed schedules for learning rates.

It is better to have

- ▶ larger rates for weights with smaller updates,
- ▶ smaller rates for weights with larger updates.

AdaGrad uses individually adapting learning rate for each weight.

# SGD with AdaGrad

- ▶ weights in  $\vec{w}^{(0)}$  are randomly initialized to values close to 0
- ▶ in the step  $t + 1$  (here  $t = 0, 1, 2 \dots$ ), compute  $\vec{w}^{(t+1)}$  :
  - ▶ Choose (randomly) a minibatch  $T \subseteq \{1, \dots, p\}$
  - ▶ Compute

$$w_{ji}^{(t+1)} = w_{ji}^{(t)} + \Delta w_{ji}^{(t)}$$

# SGD with AdaGrad

- ▶ weights in  $\vec{w}^{(0)}$  are randomly initialized to values close to 0
- ▶ in the step  $t + 1$  (here  $t = 0, 1, 2 \dots$ ), compute  $\vec{w}^{(t+1)}$  :
  - ▶ Choose (randomly) a minibatch  $T \subseteq \{1, \dots, p\}$
  - ▶ Compute

$$w_{ji}^{(t+1)} = w_{ji}^{(t)} + \Delta w_{ji}^{(t)}$$

where

$$\Delta w_{ji}^{(t)} = - \frac{\eta}{\sqrt{r_{ji}^{(t)}} + \delta} \cdot \sum_{k \in T} \frac{\partial E_k}{\partial w_{ji}}(\vec{w}^{(t)})$$

and

$$r_{ji}^{(t)} = r_{ji}^{(t-1)} + \left( \sum_{k \in T} \frac{\partial E_k}{\partial w_{ji}}(\vec{w}^{(t)}) \right)^2$$

- ▶  $\eta$  is a constant expressing the influence of the learning rate, typically 0.01.
- ▶  $\delta > 0$  is a smoothing term (typically 1e-8) avoiding division by 0.

The main disadvantage of AdaGrad is the accumulation of the gradient throughout the whole learning process.

In case the learning needs to get over several "hills" before settling in a deep "valley", the weight updates get far too small before getting to it.

RMSProp uses an exponentially decaying average to discard history from the extreme past so that it can converge rapidly after finding a convex bowl, as if it were an instance of the AdaGrad algorithm initialized within that bowl.

# SGD with RMSProp

- ▶ weights in  $\vec{w}^{(0)}$  are randomly initialized to values close to 0
- ▶ in the step  $t + 1$  (here  $t = 0, 1, 2 \dots$ ), compute  $\vec{w}^{(t+1)}$  :
  - ▶ Choose (randomly) a minibatch  $T \subseteq \{1, \dots, p\}$
  - ▶ Compute

$$w_{ji}^{(t+1)} = w_{ji}^{(t)} + \Delta w_{ji}^{(t)}$$

# SGD with RMSProp

- ▶ weights in  $\vec{w}^{(0)}$  are randomly initialized to values close to 0
- ▶ in the step  $t + 1$  (here  $t = 0, 1, 2 \dots$ ), compute  $\vec{w}^{(t+1)}$  :
  - ▶ Choose (randomly) a minibatch  $T \subseteq \{1, \dots, p\}$
  - ▶ Compute

$$w_{ji}^{(t+1)} = w_{ji}^{(t)} + \Delta w_{ji}^{(t)}$$

where

$$\Delta w_{ji}^{(t)} = -\frac{\eta}{\sqrt{r_{ji}^{(t)} + \delta}} \cdot \sum_{k \in T} \frac{\partial E_k}{\partial w_{ji}}(\vec{w}^{(t)})$$

and

$$r_{ji}^{(t)} = \rho r_{ji}^{(t-1)} + (1 - \rho) \left( \sum_{k \in T} \frac{\partial E_k}{\partial w_{ji}}(\vec{w}^{(t)}) \right)^2$$

- ▶  $\eta$  is a constant expressing the influence of the learning rate (Hinton suggests  $\rho = 0.9$  and  $\eta = 0.001$ ).
- ▶  $\delta > 0$  is a smoothing term (typically  $1e-8$ ) avoiding division by 0.

## Other optimization methods

There are more methods such as AdaDelta, Adam (roughly RMSProp combined with momentum), etc.

A natural question: Which algorithm should one choose?

## Other optimization methods

There are more methods such as AdaDelta, Adam (roughly RMSProp combined with momentum), etc.

A natural question: Which algorithm should one choose?

Unfortunately, there is currently no consensus on this point.

According to a recent study, the family of algorithms with adaptive learning rates (represented by RMSProp and AdaDelta) performed fairly robustly, no single best algorithm has emerged.



## Other optimization methods

There are more methods such as AdaDelta, Adam (roughly RMSProp combined with momentum), etc.

A natural question: Which algorithm should one choose?

Unfortunately, there is currently no consensus on this point.

According to a recent study, the family of algorithms with adaptive learning rates (represented by RMSProp and AdaDelta) performed fairly robustly, no single best algorithm has emerged.

Currently, the most popular optimization algorithms actively in use include SGD, SGD with momentum, RMSProp, RMSProp with momentum, AdaDelta and Adam.

The choice of which algorithm to use, at this point, seems to depend largely on the user's familiarity with the algorithm.

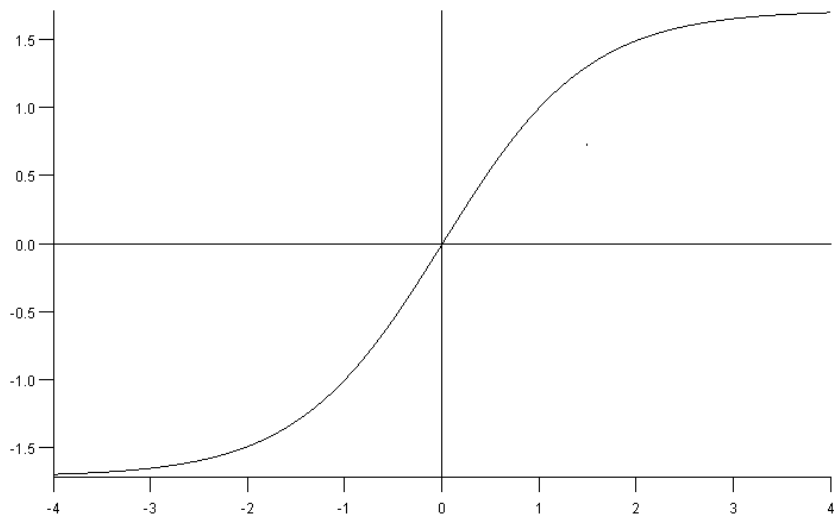
# Choice of (hidden) activations

Generic requirements imposed on activation functions:

1. differentiability  
(to do gradient descent)
2. non-linearity  
(linear multi-layer networks are equivalent to single-layer)
3. monotonicity  
(local extrema of activation functions induce local extrema of the error function)
4. "linearity"  
(i.e. preserve as much linearity as possible; linear models are easiest to fit; find the "minimum" non-linearity needed to solve a given task)

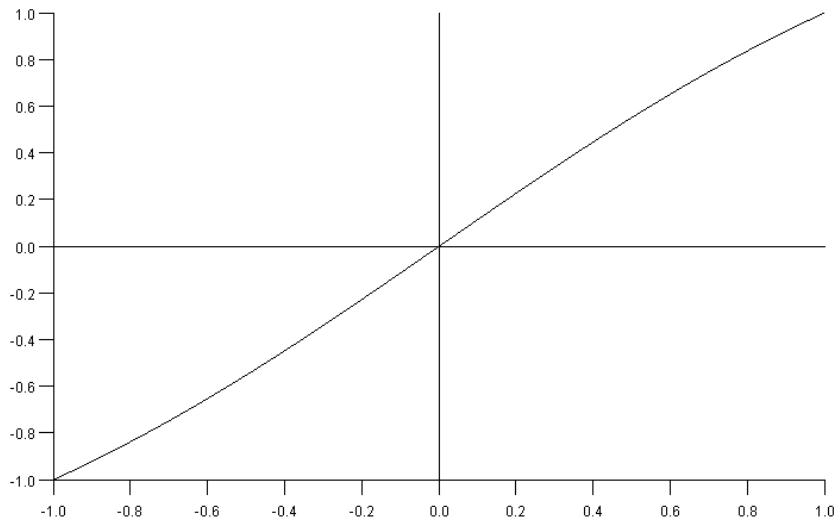
The choice of activation functions is closely related to input preprocessing and the initial choice of weights. I will illustrate the reasoning on sigmoidal functions; say few words about other activation functions later.

## Activation functions – tanh



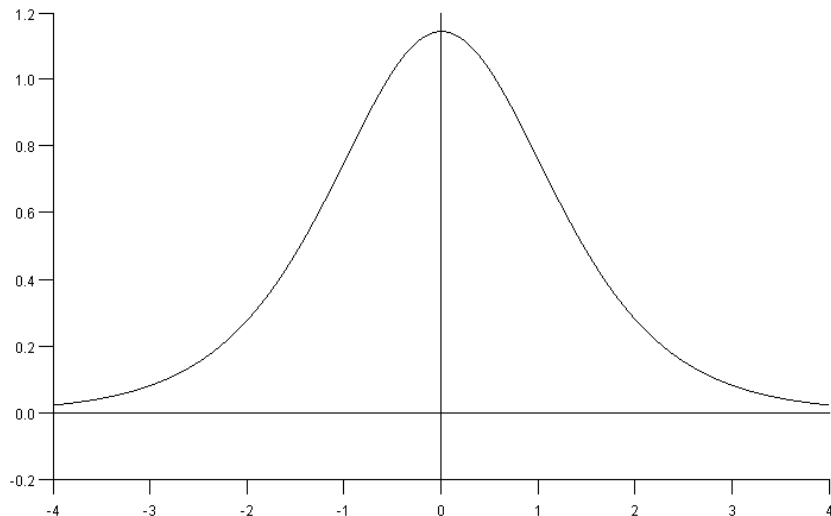
$\sigma(\xi) = 1.7159 \cdot \tanh(\frac{2}{3} \cdot \xi)$ , we have  $\lim_{\xi \rightarrow \infty} \sigma(\xi) = 1.7159$  and  $\lim_{\xi \rightarrow -\infty} \sigma(\xi) = -1.7159$

## Activation functions – tanh



$\sigma(\xi) = 1.7159 \cdot \tanh(\frac{2}{3} \cdot \xi)$  is almost linear on  $[-1, 1]$

## Activation functions – tanh



first derivative:  $\sigma(\xi) = 1.7159 \cdot \tanh(\frac{2}{3} \cdot \xi)$

# Input preprocessing

- ▶ Some inputs may be much larger than others.

E.g.: Height vs weight of a person, maximum speed of a car (in km/h) vs its price (in CZK), etc.

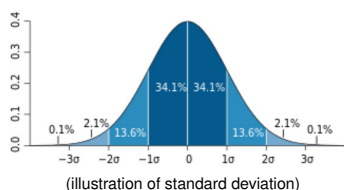
# Input preprocessing

- ▶ Some inputs may be much larger than others.  
E.g.: Height vs weight of a person, maximum speed of a car (in km/h) vs its price (in CZK), etc.
- ▶ Large inputs have greater influence on the training than the small ones. In addition, too large inputs may slow down learning (saturation of activation functions).

# Input preprocessing

- ▶ Some inputs may be much larger than others.  
E.g.: Height vs weight of a person, maximum speed of a car (in km/h) vs its price (in CZK), etc.
- ▶ Large inputs have greater influence on the training than the small ones. In addition, too large inputs may slow down learning (saturation of activation functions).
- ▶ Typical standardization:
  - ▶ average = 0 (subtract the mean)
  - ▶ variance = 1 (divide by the standard deviation)

Here the mean and standard deviation may be estimated from data (the training set).





## Initial weights (for tanh)

- ▶ Assume weights chosen in random. What distribution?

## Initial weights (for tanh)

- ▶ Assume weights chosen in random. What distribution?
- ▶ Consider the activation function  $\sigma(\xi) = 1.7159 \cdot \tanh(\frac{2}{3} \cdot \xi)$  for all neurons.
  - ▶  $\sigma$  is almost linear on  $[-1, 1]$ ,
  - ▶  $\sigma$  saturates out of the interval  $[-4, 4]$  (i.e. it is close to its limit values and its derivative is close to 0).

## Initial weights (for tanh)

- ▶ Assume weights chosen in random. What distribution?
- ▶ Consider the activation function  $\sigma(\xi) = 1.7159 \cdot \tanh(\frac{2}{3} \cdot \xi)$  for all neurons.
  - ▶  $\sigma$  is almost linear on  $[-1, 1]$ ,
  - ▶  $\sigma$  saturates out of the interval  $[-4, 4]$  (i.e. it is close to its limit values and its derivative is close to 0).

Thus

- ▶ for too small weights we may get (almost) linear model.
- ▶ for too large weights the activations may get saturated and the learning will be very slow.

Hence, we want to choose weights so that the inner potentials of neurons will be roughly in the interval  $[-1, 1]$ .

# Normal LeCun initialization

- ▶ Assume the input data have the mean = 0 and the variance = 1. Consider a neuron  $j$  from the first layer with  $n$  inputs. Assume its weights chosen randomly by the normal distribution  $\mathcal{N}(0, w^2)$ . Assume that all random choices are independent of each other.
- ▶ **The rule:** Choose the standard deviation of weights  $w$  so that the *standard deviation* of  $\xi_j$  (denote by  $\sigma_j$ ) satisfies  $\sigma_j \approx 1$ .

# Normal LeCun initialization

- ▶ Assume the input data have the mean = 0 and the variance = 1. Consider a neuron  $j$  from the first layer with  $n$  inputs. Assume its weights chosen randomly by the normal distribution  $\mathcal{N}(0, w^2)$ . Assume that all random choices are independent of each other.
- ▶ **The rule:** Choose the standard deviation of weights  $w$  so that the *standard deviation* of  $\xi_j$  (denote by  $\sigma_j$ ) satisfies  $\sigma_j \approx 1$ .
- ▶ Basic properties of the variance of independent variables give  $\sigma_j = \sqrt{n} \cdot w$ .

Thus by putting  $w = \sqrt{\frac{1}{n}}$  we obtain  $\sigma_j = 1$ .

# Normal LeCun initialization

- ▶ Assume the input data have the mean = 0 and the variance = 1. Consider a neuron  $j$  from the first layer with  $n$  inputs. Assume its weights chosen randomly by the normal distribution  $\mathcal{N}(0, w^2)$ . Assume that all random choices are independent of each other.
- ▶ **The rule:** Choose the standard deviation of weights  $w$  so that the *standard deviation* of  $\xi_j$  (denote by  $o_j$ ) satisfies  $o_j \approx 1$ .
- ▶ Basic properties of the variance of independent variables give  $o_j = \sqrt{n} \cdot w$ .

Thus by putting  $w = \sqrt{\frac{1}{n}}$  we obtain  $o_j = 1$ .

- ▶ The same works for higher layers,  $n$  corresponds to the number of neurons in the layer one level lower.

This gives *normal LeCun initialization*:

$$w_i \sim \mathcal{N}\left(0, \frac{1}{n}\right)$$

# Derivation of the LeCun initialization

Consider a single neuron without bias with the inner potential

$$\xi = \sum_{i=1}^n w_i x_i$$

# Derivation of the LeCun initialization

Consider a single neuron without bias with the inner potential

$$\xi = \sum_{i=1}^n w_i x_i$$

Consider all  $w_i$  and  $x_i$  as **independent** random variables (hence also  $\xi$  is a random variable) where

- ▶  $w_i \in \mathcal{N}(0, w^2)$  for  $i = 1, \dots, n$  where  $w$  is a constant,
- ▶  $\mathbb{E}x_i = 0$  and  $\text{Var}[x_i] = \mathbb{E}[(x_i - \mathbb{E}x_i)^2] = 1$  for  $i = 1, \dots, n$



# Derivation of the LeCun initialization

Consider a single neuron without bias with the inner potential

$$\xi = \sum_{i=1}^n w_i x_i$$

Consider all  $w_i$  and  $x_i$  as **independent** random variables (hence also  $\xi$  is a random variable) where

- ▶  $w_i \in \mathcal{N}(0, w^2)$  for  $i = 1, \dots, n$  where  $w$  is a constant,
- ▶  $\mathbb{E}x_i = 0$  and  $\text{Var}[x_i] = \mathbb{E}[(x_i - \mathbb{E}x_i)^2] = 1$  for  $i = 1, \dots, n$

We prove that  $\mathbb{E}\xi = n \cdot w^2$  as follows:

# Derivation of the LeCun initialization

Consider a single neuron without bias with the inner potential

$$\xi = \sum_{i=1}^n w_i x_i$$

Consider all  $w_i$  and  $x_i$  as **independent** random variables (hence also  $\xi$  is a random variable) where

- ▶  $w_i \in \mathcal{N}(0, w^2)$  for  $i = 1, \dots, n$  where  $w$  is a constant,
- ▶  $\mathbb{E}x_i = 0$  and  $\text{Var}[x_i] = \mathbb{E}[(x_i - \mathbb{E}x_i)^2] = 1$  for  $i = 1, \dots, n$

We prove that  $\mathbb{E}\xi = n \cdot w^2$  as follows:

$$\mathbb{E}\xi = \mathbb{E} \sum_{i=1}^n w_i x_i = \sum_{i=1}^n \mathbb{E} w_i x_i \stackrel{\text{ind.}}{=} \sum_{i=1}^n \mathbb{E} w_i \mathbb{E} x_i = 0$$

# Derivation of the LeCun initialization

Consider a single neuron without bias with the inner potential

$$\xi = \sum_{i=1}^n w_i x_i$$

Consider all  $w_i$  and  $x_i$  as **independent** random variables (hence also  $\xi$  is a random variable) where

- ▶  $w_i \in \mathcal{N}(0, w^2)$  for  $i = 1, \dots, n$  where  $w$  is a constant,
- ▶  $\mathbb{E}x_i = 0$  and  $\text{Var}[x_i] = \mathbb{E}[(x_i - \mathbb{E}x_i)^2] = 1$  for  $i = 1, \dots, n$

We prove that  $\mathbb{E}\xi = n \cdot w^2$  as follows:

$$\mathbb{E}\xi = \mathbb{E} \sum_{i=1}^n w_i x_i = \sum_{i=1}^n \mathbb{E} w_i x_i \stackrel{\text{ind.}}{=} \sum_{i=1}^n \mathbb{E} w_i \mathbb{E} x_i = 0$$

$$\text{and } \text{Var}[w_i x_i] = \mathbb{E}[w_i^2 x_i^2] - \mathbb{E}[w_i x_i]^2 \stackrel{\text{ind.}}{=} \mathbb{E}[w_i^2] \mathbb{E}[x_i^2] - 0 = w^2$$

# Derivation of the LeCun initialization

Consider a single neuron without bias with the inner potential

$$\xi = \sum_{i=1}^n w_i x_i$$

Consider all  $w_i$  and  $x_i$  as **independent** random variables (hence also  $\xi$  is a random variable) where

- ▶  $w_i \in \mathcal{N}(0, w^2)$  for  $i = 1, \dots, n$  where  $w$  is a constant,
- ▶  $\mathbb{E}x_i = 0$  and  $\text{Var}[x_i] = \mathbb{E}[(x_i - \mathbb{E}x_i)^2] = 1$  for  $i = 1, \dots, n$

We prove that  $\mathbb{E}\xi = n \cdot w^2$  as follows:

$$\mathbb{E}\xi = \mathbb{E} \sum_{i=1}^n w_i x_i = \sum_{i=1}^n \mathbb{E} w_i x_i \stackrel{\text{ind.}}{=} \sum_{i=1}^n \mathbb{E} w_i \mathbb{E} x_i = 0$$

and  $\text{Var}[w_i x_i] = \mathbb{E}[w_i^2 x_i^2] - \mathbb{E}[w_i x_i]^2 \stackrel{\text{ind.}}{=} \mathbb{E}[w_i^2] \mathbb{E}[x_i^2] - 0 = w^2$   
implies

$$\text{Var}[\xi] = \text{Var}\left[\sum_{i=1}^n w_i x_i\right] \stackrel{\text{ind.}}{=} \sum_{i=1}^n \text{Var}[w_i x_i] = \sum_{i=1}^n w^2 = n \cdot w^2$$

# Normal Glorot initialization

The previous heuristics for weight initialization ignores variance of the gradient (i.e. it is concerned only with the "size" of activations in the forward pass).

# Normal Glorot initialization

The previous heuristics for weight initialization ignores variance of the gradient (i.e. it is concerned only with the "size" of activations in the forward pass).

Glorot & Bengio (2010) presented a **normalized initialization** by choosing weights randomly from the following normal distribution:

$$N\left(0, \frac{2}{m+n}\right) = N\left(0, \frac{1}{(m+n)/2}\right)$$

Here  $n$  is the number of inputs to the layer,  $m$  is the number of neurons in the layer above.

# Normal Glorot initialization

The previous heuristics for weight initialization ignores variance of the gradient (i.e. it is concerned only with the "size" of activations in the forward pass).

Glorot & Bengio (2010) presented a **normalized initialization** by choosing weights randomly from the following normal distribution:

$$N\left(0, \frac{2}{m+n}\right) = N\left(0, \frac{1}{(m+n)/2}\right)$$

Here  $n$  is the number of inputs to the layer,  $m$  is the number of neurons in the layer above.

This is designed to compromise between the goal of initializing all layers to have the same activation variance and the goal of initializing all layers to have the same gradient variance.

This gives *normal Glorot initialization* (also called *normal Xavier initialization*):

$$w_i \sim \mathcal{N}\left(0, \frac{2}{m+n}\right)$$

# Uniform LeCun initialization

- ▶ Assume that the input data have mean = 0 and variance = 1.

Consider a neuron  $j$  from the first layer with  $n$  inputs. Assume its weights chosen randomly by the uniform distribution  $U(-w, w)$ .

Assume that all random choices are independent of each other.

- ▶ As before, we want the standard deviation  $\sigma_j$  of the inner potential  $\xi_j$  to be approximately 1.
- ▶ Basic properties of the variance of independent variables give  $\sigma_j = \sqrt{\frac{n}{3}} \cdot w$ .

Thus by putting  $w = \sqrt{\frac{3}{n}}$  we obtain  $\sigma_j = 1$ .

We obtain *uniform LeCun initialization*:

$$w_i \sim U\left(-\sqrt{\frac{3}{n}}, \sqrt{\frac{3}{n}}\right)$$



# Uniform Glorot initialization

Similarly to the normal case, we want to normalize the initialization w.r.t. both forward and backward passes.

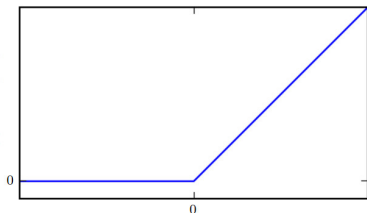
We obtain *uniform Glorot initialization* (aka *uniform Xavier init.*):

$$w_i \sim U\left(-\sqrt{\frac{6}{m+n}}, \sqrt{\frac{6}{m+n}}\right) = U\left(-\sqrt{\frac{3}{(m+n)/2}}, \sqrt{\frac{3}{(m+n)/2}}\right)$$

Here  $n$  is the number of inputs to the layer,  $m$  is the number of neurons in the layer above.

# Modern activation functions

For hidden neurons sigmoidal functions are often substituted with piece-wise linear activations functions. Most prominent is ReLU:



$$\sigma(\xi) = \max\{0, \xi\}$$

- ▶ THE default activation function recommended for use with most feedforward neural networks.
- ▶ As close to linear function as possible; very simple; does not saturate for large potentials.
- ▶ Dead for negative potentials.

## Normal He initialization

- ▶ The ReLU is not as sensitive to the large variance of the inner potential as sigmoidal functions (large variance does not matter as much).

## Normal He initialization

- ▶ The ReLU is not as sensitive to the large variance of the inner potential as sigmoidal functions (large variance does not matter as much).
- ▶ Still the variance is good to be constant (at least due to the output layer).

## Normal He initialization

- ▶ The ReLU is not as sensitive to the large variance of the inner potential as sigmoidal functions (large variance does not matter as much).
- ▶ Still the variance is good to be constant (at least due to the output layer).
- ▶ LeCun initialization cannot be justified for ReLU due to the following reason:

The ReLU is not a symmetric function. So even if the inner potential  $\xi_j$  has mean = 0 and variance = 1, it is not true of the output (the variance is halved).

## Normal He initialization

- ▶ The ReLU is not as sensitive to the large variance of the inner potential as sigmoidal functions (large variance does not matter as much).
- ▶ Still the variance is good to be constant (at least due to the output layer).
- ▶ LeCun initialization cannot be justified for ReLU due to the following reason:

The ReLU is not a symmetric function. So even if the inner potential  $\xi_j$  has mean = 0 and variance = 1, it is not true of the output (the variance is halved).

Modifying the normal LeCun initialization to take the halving variance into account, we obtain *normal He initialization*:

$$w_i \in \mathcal{N}\left(0, \frac{2}{n}\right) \quad \left(\text{LeCun is } w_i \in \mathcal{N}\left(0, \frac{1}{n}\right)\right)$$

# More modern activation functions

- ▶ Leaky ReLU (greenboard):
  - ▶ Generalizes ReLU, not dead for negative potentials.
  - ▶ Experimentally not much better than ReLU.

# More modern activation functions

- ▶ Leaky ReLU (greenboard):
  - ▶ Generalizes ReLU, not dead for negative potentials.
  - ▶ Experimentally not much better than ReLU.
- ▶ ELU: "Smoothed" ReLU:

$$\sigma(\xi) = \begin{cases} \alpha(\exp(\xi) - 1) & \text{for } \xi < 0 \\ \xi & \text{for } \xi \geq 0 \end{cases}$$

Here  $\alpha$  is a parameter, ELU converges to  $-\alpha$  as  $\xi \rightarrow -\infty$ . As opposed to ReLU: Smooth, always non-zero gradient (but saturates), slower to compute.



# More modern activation functions

- ▶ Leaky ReLU (greenboard):
  - ▶ Generalizes ReLU, not dead for negative potentials.
  - ▶ Experimentally not much better than ReLU.
- ▶ ELU: "Smoothed" ReLU:

$$\sigma(\xi) = \begin{cases} \alpha(\exp(\xi) - 1) & \text{for } \xi < 0 \\ \xi & \text{for } \xi \geq 0 \end{cases}$$

Here  $\alpha$  is a parameter, ELU converges to  $-\alpha$  as  $\xi \rightarrow -\infty$ . As opposed to ReLU: Smooth, always non-zero gradient (but saturates), slower to compute.

- ▶ SELU: Scaled variant of ELU: :

$$\sigma(\xi) = \lambda \begin{cases} \alpha(\exp(\xi) - 1) & \text{for } \xi < 0 \\ \xi & \text{for } \xi \geq 0 \end{cases}$$

*Self-normalizing*, i.e. output of each layer will tend to preserve a mean (close to) 0 and a standard deviation (close to) 1 for  $\lambda \approx 1.050$  and  $\alpha \approx 1.673$ , properly initialized weights (see below) and normalized inputs (zero mean, standard deviation 1).

# Initializing with Normal Distribution

Denote by  $n$  the number of inputs to the initialized layer, and  $m$  the number of neurons in the layer.

- ▶ normal Glorot:

$$w_i \sim \mathcal{N}\left(0, \frac{2}{m+n}\right)$$

Suitable for none, tanh, logistic, softmax

# Initializing with Normal Distribution

Denote by  $n$  the number of inputs to the initialized layer, and  $m$  the number of neurons in the layer.

- ▶ normal Glorot:

$$w_i \sim \mathcal{N}\left(0, \frac{2}{m+n}\right)$$

Suitable for none, tanh, logistic, softmax

- ▶ normal He:

$$w_i \in \mathcal{N}\left(0, \frac{2}{n}\right)$$

Suitable for ReLU, leaky ReLU

# Initializing with Normal Distribution

Denote by  $n$  the number of inputs to the initialized layer, and  $m$  the number of neurons in the layer.

- ▶ normal Glorot:

$$w_i \sim \mathcal{N}\left(0, \frac{2}{m+n}\right)$$

Suitable for none, tanh, logistic, softmax

- ▶ normal He:

$$w_i \in \mathcal{N}\left(0, \frac{2}{n}\right)$$

Suitable for ReLU, leaky ReLU

- ▶ normal LeCun:

$$w_i \sim \mathcal{N}\left(0, \frac{1}{n}\right)$$

Suitable for SELU (by the authors)

# How to choose activation of hidden neurons

- ▶ The default is ReLU.
- ▶ According to Aurélien Géron:

*SELU > ELU > leakyReLU > ReLU > tanh > logistic*

For discussion see: Hands-On Machine Learning with Scikit-Learn, Keras, and TensorFlow: Concepts, Tools, and Techniques to Build Intelligent Systems, Aurélien Géron

## Batch normalization (roughly)

**Intuition:** Instead of keeping mean = 0 and variance = 1 implicitly due to a clever weight initialization, we may **renormalize values of neurons** throughout the layers.

## Batch normalization (roughly)

**Intuition:** Instead of keeping mean = 0 and variance = 1 implicitly due to a clever weight initialization, we may **renormalize values of neurons** throughout the layers.

Consider the  $\ell$ -th layer of the network.

Note that the output values of neurons in the  $\ell$ -th layer can be seen as inputs to the sub-network consisting of all layers above the  $\ell$ -th one.

## Batch normalization (roughly)

**Intuition:** Instead of keeping mean = 0 and variance = 1 implicitly due to a clever weight initialization, we may **renormalize values of neurons** throughout the layers.

Consider the  $\ell$ -th layer of the network.

Note that the output values of neurons in the  $\ell$ -th layer can be seen as inputs to the sub-network consisting of all layers above the  $\ell$ -th one.

What if we standardize the values of the  $\ell$ -th layer as we did with the input data?

For this we need to form a "dataset" of values of the  $\ell$ -th layer.



## Batch normalization (roughly)

Let us consider the  $\ell$ -th layer with  $n$  neurons.

Consider a batch of training examples:

$$\{(\vec{x}_k, \vec{d}_k) \mid k = 1, \dots, p\}$$

(This is typically a minibatch.)

# Batch normalization (roughly)

Let us consider the  $\ell$ -th layer with  $n$  neurons.

Consider a batch of training examples:

$$\{(\vec{x}_k, \vec{d}_k) \mid k = 1, \dots, p\}$$

(This is typically a minibatch.)

- ▶ For every  $k = 1, \dots, p$ : Compute the values of neurons in the  $\ell$ -th layer for the input  $\vec{x}_k$  and obtain a vector

$$\vec{z}_k = (\vec{z}_{k1}, \dots, \vec{z}_{kn})$$

## Batch normalization (roughly)

Let us consider the  $\ell$ -th layer with  $n$  neurons.

Consider a batch of training examples:

$$\{(\vec{x}_k, \vec{d}_k) \mid k = 1, \dots, p\}$$

(This is typically a minibatch.)

- ▶ For every  $k = 1, \dots, p$ : Compute the values of neurons in the  $\ell$ -th layer for the input  $\vec{x}_k$  and obtain a vector

$$\vec{z}_k = (\vec{z}_{k1}, \dots, \vec{z}_{kn})$$

- ▶ Set all components of all vectors  $\vec{z}_k$  to the mean  $= 0$  and the variance  $= 1$  and obtain *normalized vectors*:  $\hat{z}_1, \dots, \hat{z}_p$ .

## Batch normalization (roughly)

Let us consider the  $\ell$ -th layer with  $n$  neurons.

Consider a batch of training examples:

$$\{(\vec{x}_k, \vec{d}_k) \mid k = 1, \dots, p\}$$

(This is typically a minibatch.)

- ▶ For every  $k = 1, \dots, p$ : Compute the values of neurons in the  $\ell$ -th layer for the input  $\vec{x}_k$  and obtain a vector

$$\vec{z}_k = (z_{k1}, \dots, z_{kn})$$

- ▶ Set all components of all vectors  $\vec{z}_k$  to the mean  $= 0$  and the variance  $= 1$  and obtain *normalized vectors*:  $\hat{z}_1, \dots, \hat{z}_p$ .
- ▶ For every  $k = 1, \dots, p$  give

$$\vec{\gamma} \cdot \hat{z}_k + \delta$$

as the output of the  $\ell$ -th layer instead of  $\vec{z}_k$ . Here  $\vec{\gamma}$  and  $\delta$  are new trainable weights.

# Generalization

**Intuition:** Generalization = ability to cope with new unseen instances.

Data are mostly noisy, so it is not good idea to fit exactly.

In case of function approximation, the network should not return exact results as in the training set.

# Generalization

**Intuition:** Generalization = ability to cope with new unseen instances.

Data are mostly noisy, so it is not good idea to fit exactly.

In case of function approximation, the network should not return exact results as in the training set.

More formally: It is typically assumed that the training set has been generated as follows:

$$d_{kj} = g_j(\vec{x}_k) + \Theta_{kj}$$

where  $g_j$  is the "underlying" function corresponding to the output neuron  $j \in Y$  and  $\Theta_{kj}$  is random noise.

The network should fit  $g_j$  not the noise.

Methods improving generalization are called **regularization methods**.

# Regularization

Regularization is a big issue in neural networks, as they typically use a huge amount of parameters and thus are very susceptible to overfitting.

Regularization is a big issue in neural networks, as they typically use a huge amount of parameters and thus are very susceptible to overfitting.

von Neumann: **"With four parameters I can fit an elephant, and with five I can make him wiggle his trunk."**

... and I ask you prof. Neumann:

What can you fit with 40GB of parameters??



## Early stopping

Early stopping means that we stop learning before it reaches a minimum of the error  $E$ .

When to stop?

# Early stopping

Early stopping means that we stop learning before it reaches a minimum of the error  $E$ .

When to stop?

In many applications the error function is not the main thing we want to optimize.

E.g. in the case of a trading system, we typically want to maximize our profit not to minimize (strange) error functions designed to be easily differentiable.

Also, as noted before, minimizing  $E$  completely is not good for generalization.

For start: We may employ standard approach of training on one set and stopping on another one.

# Early stopping

Divide your dataset into several subsets:

- ▶ **training set** (e.g. 60%) – train the network here
- ▶ **validation set** (e.g. 20%) – use to stop the training
- ▶ **test set** (e.g. 20%) – use to evaluate the final model

What to use as a stopping rule?

# Early stopping

Divide your dataset into several subsets:

- ▶ **training set** (e.g. 60%) – train the network here
- ▶ **validation set** (e.g. 20%) – use to stop the training
- ▶ **test set** (e.g. 20%) – use to evaluate the final model

What to use as a stopping rule?

You may observe  $E$  (or any other function of interest) on the validation set, if it does not improve for last  $k$  steps, stop.

Alternatively, you may observe the gradient, if it is small for some time, stop.

(recent studies shown that this traditional rule is not too good: it may happen that the gradient is larger close to minimum values; on the other hand,  $E$  does not have to be evaluated which saves time.

To compare models you may use ML techniques such as various types of cross-validation etc.

## Size of the network

Similar problem as in the case of the training duration:

- ▶ Too small network is not able to capture intrinsic properties of the training set.
- ▶ Large networks overfit faster.

**Solution:** Optimal number of neurons :-)

# Size of the network

Similar problem as in the case of the training duration:

- ▶ Too small network is not able to capture intrinsic properties of the training set.
- ▶ Large networks overfit faster.

**Solution:** Optimal number of neurons :-)

- ▶ there are some (useless) theoretical bounds
- ▶ there are algorithms dynamically adding/removing neurons (not much use nowadays)
- ▶ In practice: Start with an existing network solving similar problem.

If you are trully desperate trying to solve a brand new problem, you may try an ancient rule of thumb: the number of neurons  $\approx$  ten times less than the number of training instances.

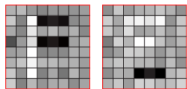
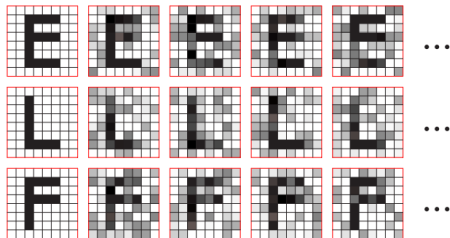
Experiment, experiment, experiment.

# Feature extraction

Consider a two layer network. Hidden neurons are supposed to represent "patterns" in the inputs.

Example: Network 64-2-3 for letter classification:

*sample training patterns*



*learned input-to-hidden weights*

# Ensemble methods

Techniques for reducing generalization error by combining several models.

The reason that ensemble methods work is that different models will usually not make all the same errors on the test set.

**Idea:** Train several different models separately, then have all of the models vote on the output for test examples.



# Ensemble methods

Techniques for reducing generalization error by combining several models.

The reason that ensemble methods work is that different models will usually not make all the same errors on the test set.

**Idea:** Train several different models separately, then have all of the models vote on the output for test examples.

## Bagging:

- ▶ Generate  $k$  training sets  $T_1, \dots, T_k$  by *sampling from  $\mathcal{T}$  uniformly with replacement*.

If the number of samples is  $|\mathcal{T}|$ , then on average  $|T_i| = (1 - 1/e)|\mathcal{T}|$ .

- ▶ For each  $i$ , train a model  $M_i$  on  $T_i$ .
- ▶ Combine outputs of the models: for regression by averaging, for classification by (majority) voting.

# Dropout

**The algorithm:** In every step of the gradient descent

- ▶ choose randomly a set  $N$  of neurons, each neuron is included independently with probability  $1/2$ ,  
(in practice, different probabilities are used as well).
- ▶ do forward and backward propagations only using the selected neurons  
(i.e. leave weights of the other neurons unchanged)

# Dropout

**The algorithm:** In every step of the gradient descent

- ▶ choose randomly a set  $N$  of neurons, each neuron is included independently with probability  $1/2$ ,  
(in practice, different probabilities are used as well).
- ▶ do forward and backward propagations only using the selected neurons  
(i.e. leave weights of the other neurons unchanged)

Dropout resembles bagging: Large ensemble of neural networks is trained "at once" on parts of the data.

Dropout is not exactly the same as bagging: The models share parameters, with each model inheriting a different subset of parameters from the parent neural network. This parameter sharing makes it possible to represent an exponential number of models with a tractable amount of memory.

In the case of bagging, each model is trained to convergence on its respective training set. This would be infeasible for large networks/training sets.

# Dropout – details

- ▶ The inner potential of a neuron  $j$  **without dropout**:

$$\xi_j = \sum_{i \in j_{\leftarrow}} w_{ji} y_i$$

- ▶ The inner potential of a neuron  $j$  **with dropout**:

$$r_i \sim \text{Bernoulli}(1/2) \quad \text{for all } i \in j_{\leftarrow} \setminus \{0\}$$

$$\xi_j = \sum_{i \in j_{\leftarrow}} w_{ji} (r_i y_i)$$

(Intuitively, randomly chosen neurons are masked out.)

- ▶ During inference do not drop out neurons and multiply values of neurons with 1/2.

This compensates for the fact that without the drop out there are twice as many neurons.

# Weight decay and L2 regularization

Generalization can be improved by removing "unimportant" weights.

Penalising large weights gives stronger indication about their importance.

# Weight decay and L2 regularization

Generalization can be improved by removing "unimportant" weights.

Penalising large weights gives stronger indication about their importance.

In every step we decrease weights (multiplicatively) as follows:

$$w_{ji}^{(t+1)} = (1 - \zeta) w_{ji}^{(t)} - \varepsilon \cdot \frac{\partial E}{\partial w_{ji}}(\vec{w}^{(t)})$$

**Intuition:** Unimportant weights will be pushed to 0, important weights will survive the decay.

# Weight decay and L2 regularization

Generalization can be improved by removing "unimportant" weights.

Penalising large weights gives stronger indication about their importance.

In every step we decrease weights (multiplicatively) as follows:

$$w_{ji}^{(t+1)} = (1 - \zeta) w_{ji}^{(t)} - \varepsilon \cdot \frac{\partial E}{\partial w_{ji}}(\vec{w}^{(t)})$$

**Intuition:** Unimportant weights will be pushed to 0, important weights will survive the decay.

Weight decay is equivalent to the gradient descent with a constant learning rate  $\varepsilon$  and the following error function:

$$E'(\vec{w}) = E(\vec{w}) + \frac{\zeta}{2\varepsilon}(\vec{w} \cdot \vec{w})$$

Here  $\frac{\zeta}{2\varepsilon}(\vec{w} \cdot \vec{w})$  is the L2 regularization that penalizes large weights.

We use the gradient descent with a constant learning rate to illustrate the equivalence between L2 regularization and the weight decay. Both methods can be combined with other learning algorithms (AdaGrad, etc.).

## More optimization, regularization ...

There are many more practical tips, optimization methods, regularization methods, etc.

For a very nice survey see

<http://www.deeplearningbook.org/>

... and also all other infinitely many urls concerned with deep learning.