

# **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(Extremely well-written online textbook (a little outdated))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("Classical" overview of the theory of neural networks (a little outdated))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>  
(Greatly 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 the 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 the 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 can we not use specialized libraries in projects?

**Q:** Why can we not 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 can we not 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 resources 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

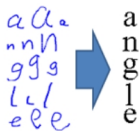
- ▶ Machine learning = construction of systems that learn their functionality from data  
(... and thus do not need to be programmed.)

# 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 "labeled" emails
    - ▶ consequently can distinguish spam from ham

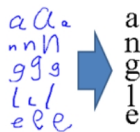
# Machine learning in general

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



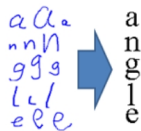
# Machine learning in general

- ▶ 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 "labeled" emails
    - ▶ consequently can distinguish spam from ham
  - ▶ handwritten text reader
    - ▶ learns from a database of handwritten letters (or text) labeled 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 "labeled" emails
    - ▶ consequently can distinguish spam from ham
  - ▶ handwritten text reader
    - ▶ learns from a database of handwritten letters (or text) labeled 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.

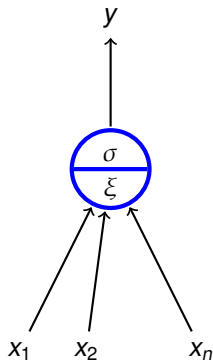
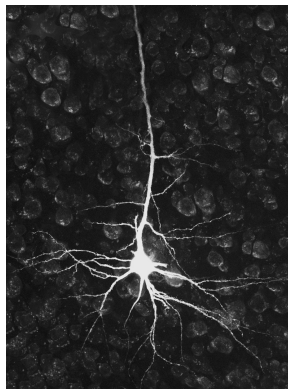
There are many types of models:

- ▶ 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.

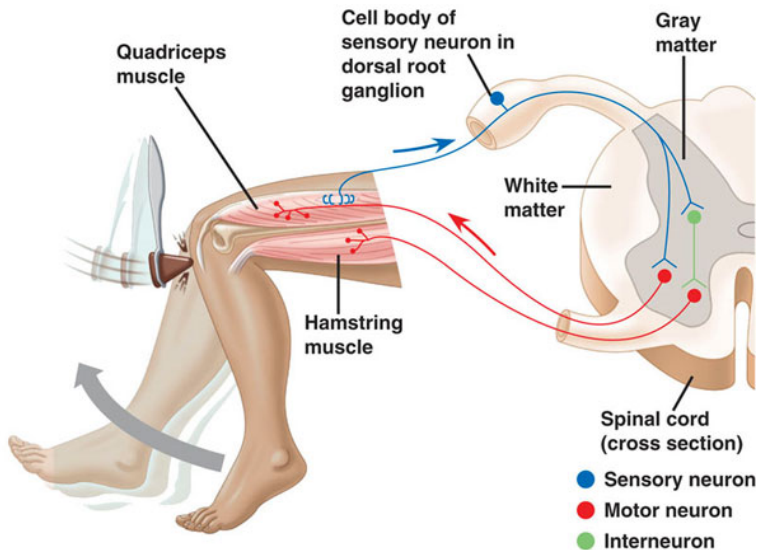
# Biological neural network

- ▶ 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

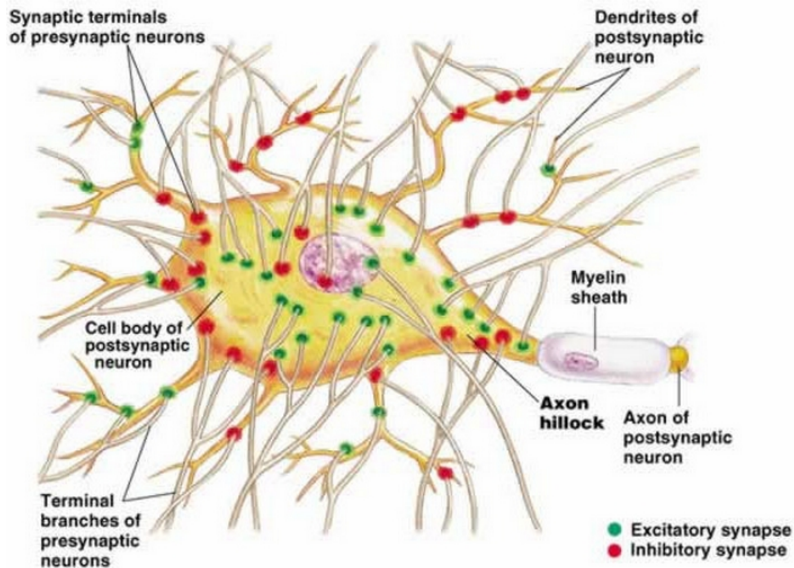
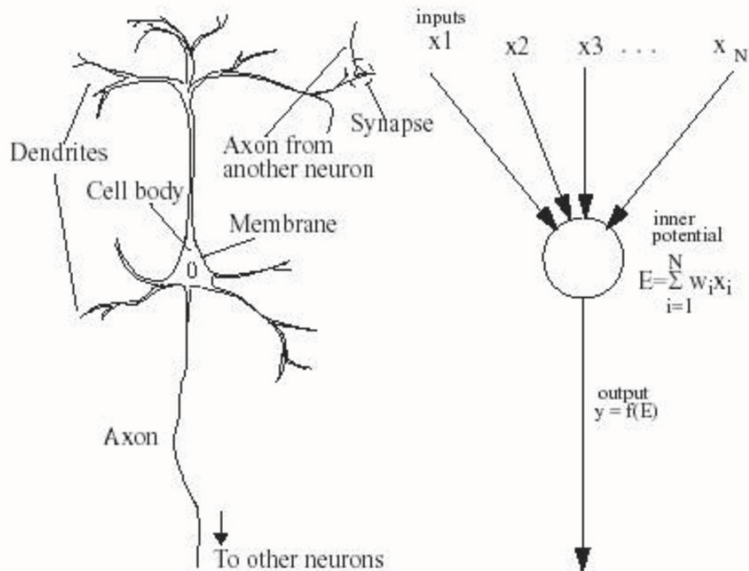


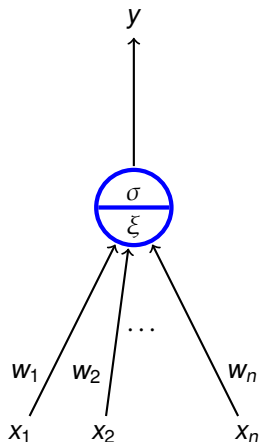
Figure 48.11(a), page 972, Campbell's *Biology*, 5th Edition

# Biological and Mathematical neurons



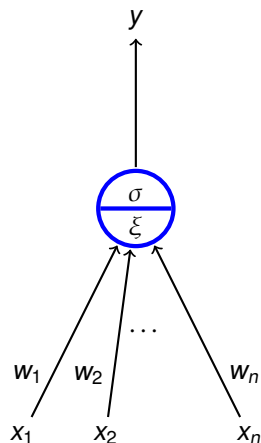
# Formal neuron (without bias)

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



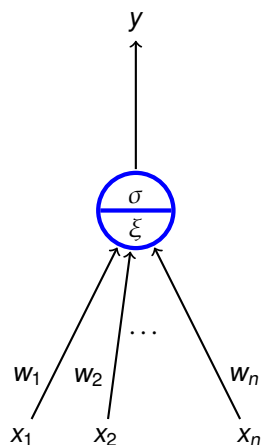


## Formal neuron (without bias)



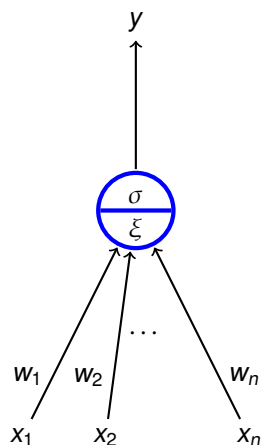
- ▶  $x_1, \dots, x_n \in \mathbb{R}$  are **inputs**
- ▶  $w_1, \dots, w_n \in \mathbb{R}$  are **weights**

## 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)



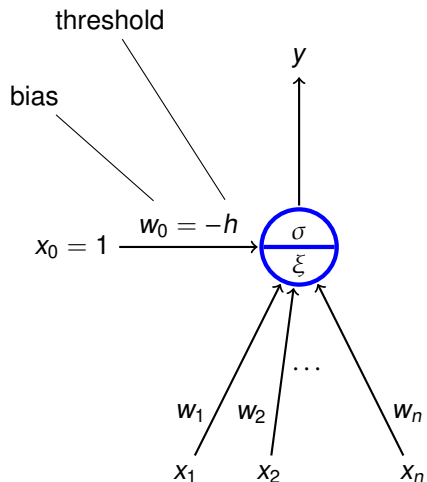
- ▶  $x_1, \dots, x_n \in \mathbb{R}$  are **inputs**
- ▶  $w_1, \dots, w_n \in \mathbb{R}$  are **weights**
- ▶  $\xi$  is an **inner potential**;  
almost always  $\xi = \sum_{i=1}^n w_i x_i$
- ▶  $y$  is an **output** given by  $y = \sigma(\xi)$   
where  $\sigma$  is an **activation function**; 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

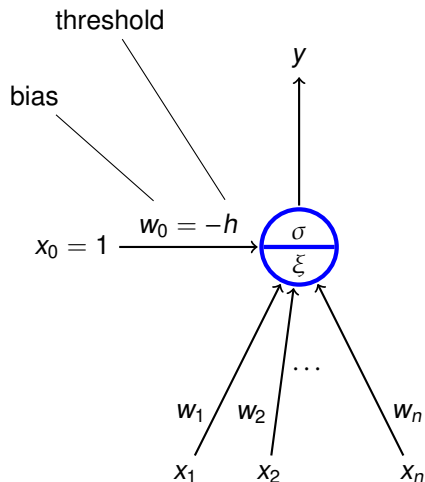
# Formal neuron (with bias)

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

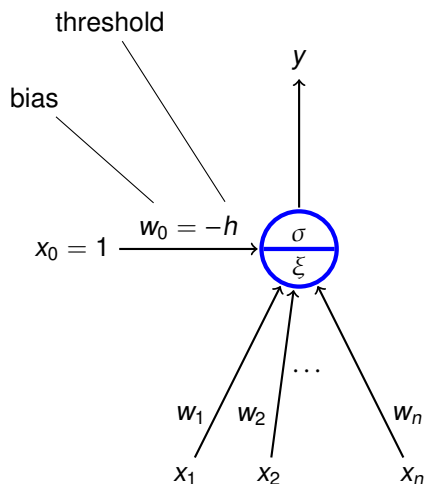


# Formal neuron (with bias)

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

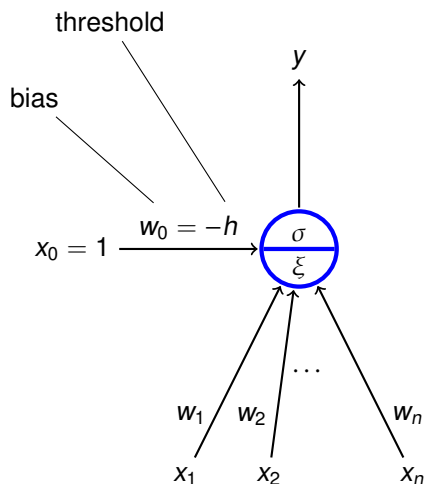


# Formal neuron (with bias)



- ▶  $x_0 = 1, x_1, \dots, x_n \in \mathbb{R}$  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)

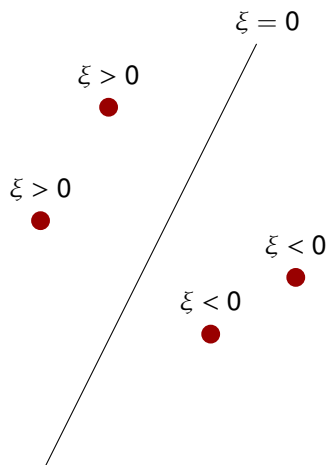


- ▶  $x_0 = 1, x_1, \dots, x_n \in \mathbb{R}$  are **inputs**
- ▶  $w_0, w_1, \dots, w_n \in \mathbb{R}$  are **weights**
- ▶  $\xi$  is an **inner potential**;  
almost always  $\xi = w_0 + \sum_{i=1}^n w_i x_i$
- ▶  $y$  is an **output** given by  $y = \sigma(\xi)$   
where  $\sigma$  is an **activation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

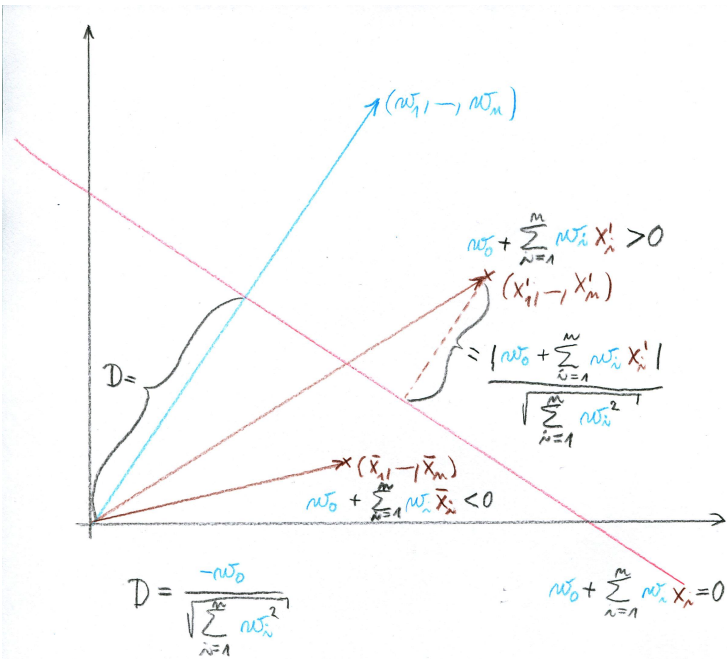
$$\xi = w_0 + \sum_{i=1}^n w_i x_i$$

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

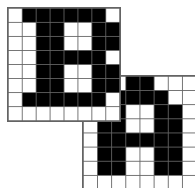
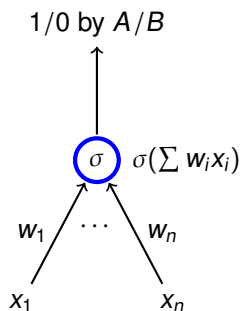
- ▶ in 2d line
- ▶ in 3d plane
- ▶ ...



# Neuron geometry

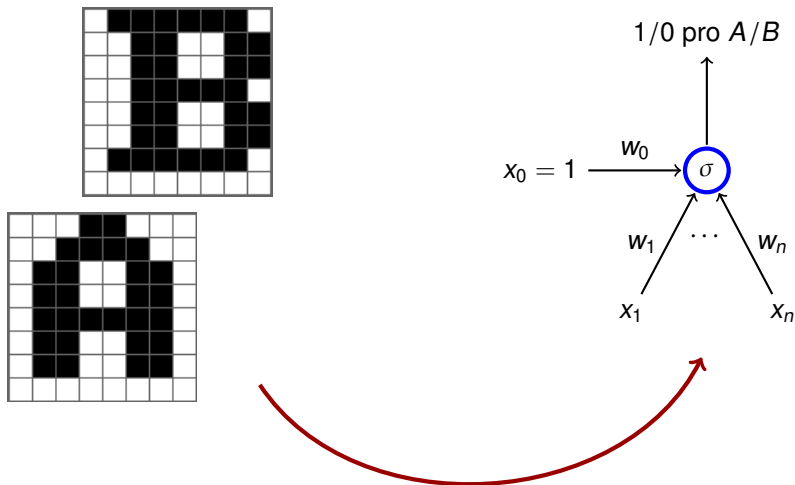


# Neuron and linear separation



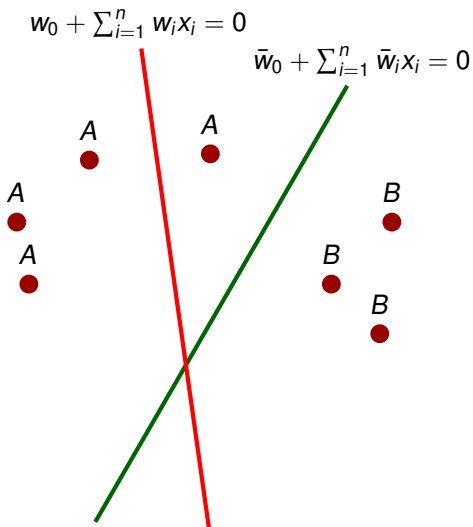
$n = 8 \cdot 8$ , i.e. the number of pixels in the images. Inputs are binary vectors of dimension  $n$  (black pixel  $\approx 1$ , white pixel  $\approx 0$ ).

# Neuron and linear separation



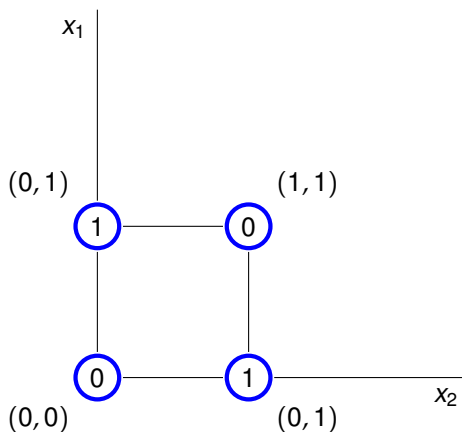
$n = 8 \cdot 8$ , i.e. the number of pixels in the images. Inputs are binary vectors of dimension  $n$  (black pixel  $\approx 1$ , white pixel  $\approx 0$ ).

# 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.

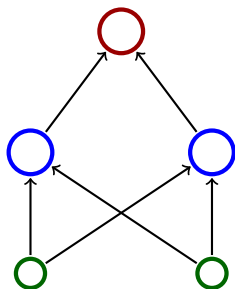
# Architecture

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

We distinguish several categories of neurons:

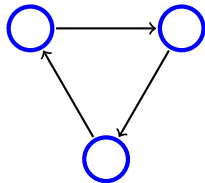
- ▶ **Output neurons**
- ▶ **Hidden neurons**
- ▶ **Input neurons**

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



# Architecture – Cycles

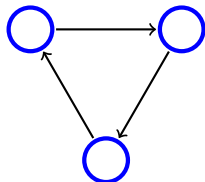
- ▶ A network is **cyclic** (recurrent) if its architecture contains a directed cycle.



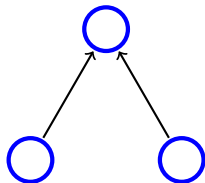


# Architecture – Cycles

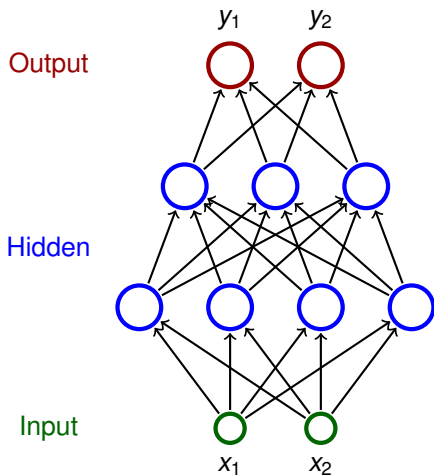
- ▶ A network is **cyclic** (recurrent) if its architecture contains a directed cycle.



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

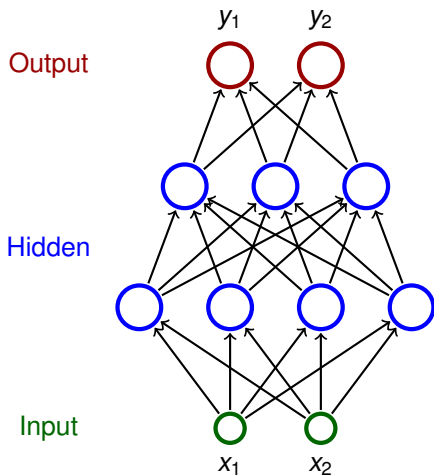


# Architecture – Multilayer Perceptron (MLP)



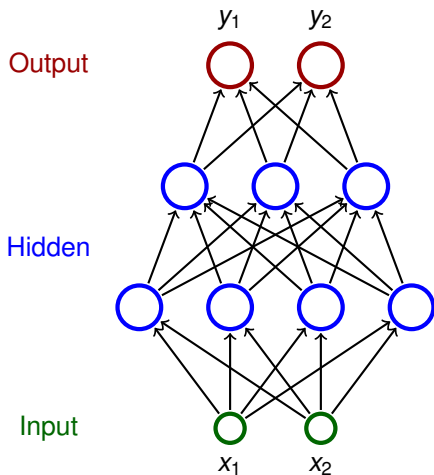
- ▶ Neurons partitioned into **layers**; one input layer, one output layer, possibly several hidden layers

# Architecture – Multilayer Perceptron (MLP)



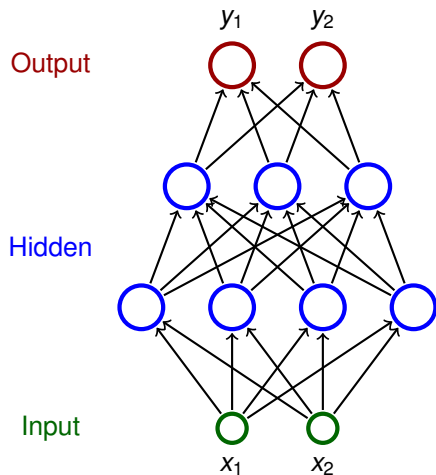
- ▶ Neurons partitioned into **layers**; one input layer, one output layer, possibly several hidden layers
- ▶ layers numbered from 0; the input layer has number 0
  - ▶ E.g. three-layer network has two hidden layers and one output layer

# 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

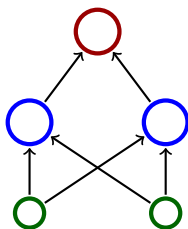
Consider a network with  $n$  neurons,  $k$  input,  $\ell$  output.

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

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

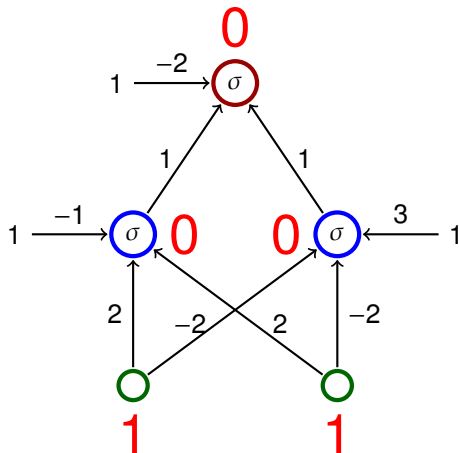
- ▶ Hyperbolic tangens

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

- ▶ ReLU

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

# Activity – XOR



- ▶ Activation function is a unit step function

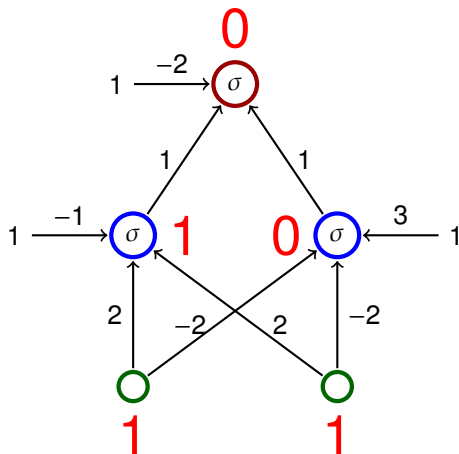
$$\sigma(\xi) = \begin{cases} 1 & \xi \geq 0; \\ 0 & \xi < 0. \end{cases}$$

- ▶ 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



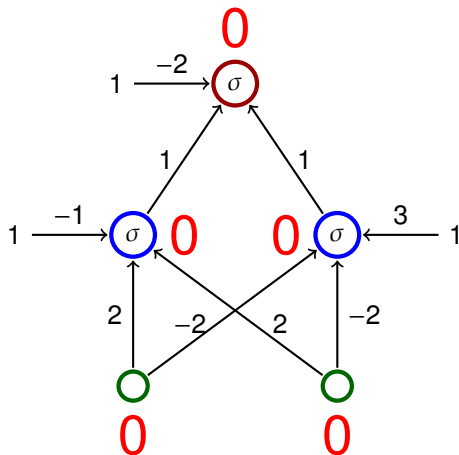
- ▶ Activation function is a unit step function

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

- ▶ 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



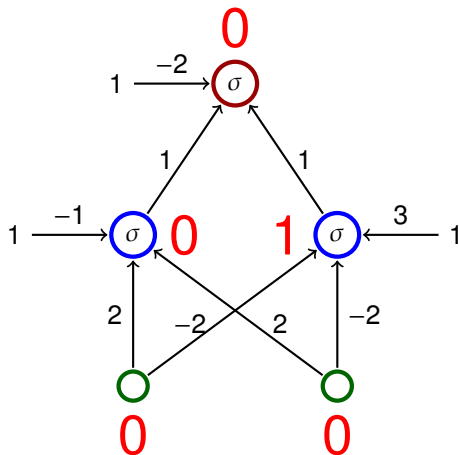
- ▶ Activation function is a unit step function

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

- ▶ 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



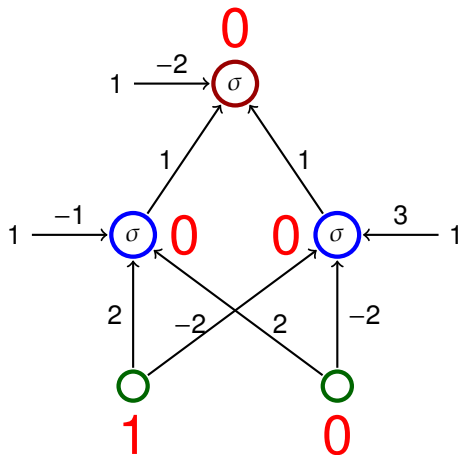
- ▶ Activation function is a unit step function

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

- ▶ 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



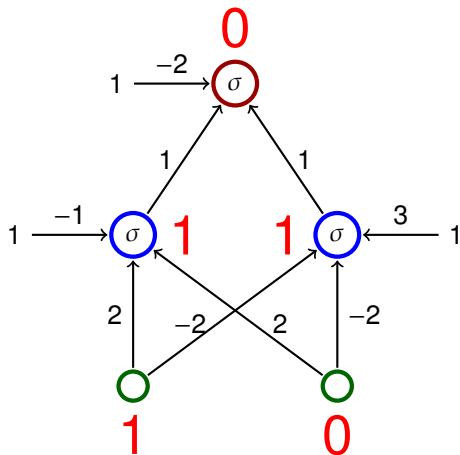
- ▶ Activation function is a unit step function

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

- ▶ 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



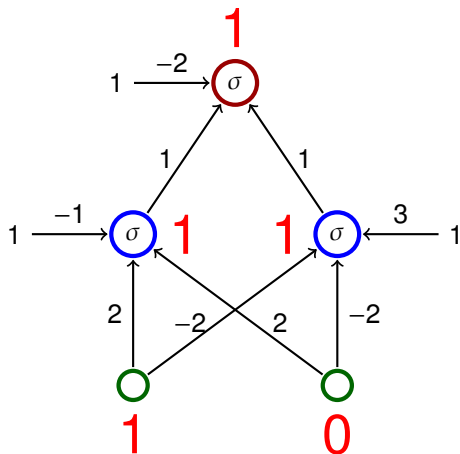
- ▶ Activation function is a unit step function

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

- ▶ 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



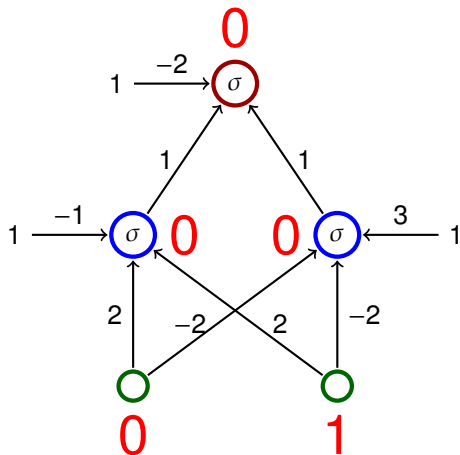
- ▶ Activation function is a unit step function

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

- ▶ 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



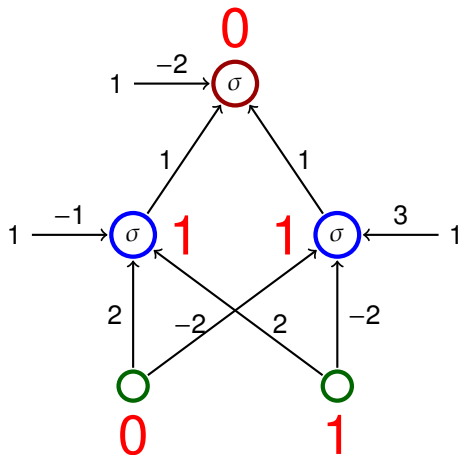
- ▶ Activation function is a unit step function

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

- ▶ 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

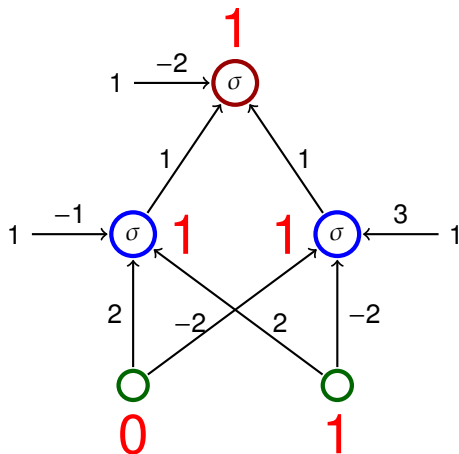
$$\sigma(\xi) = \begin{cases} 1 & \xi \geq 0; \\ 0 & \xi < 0. \end{cases}$$

- ▶ 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



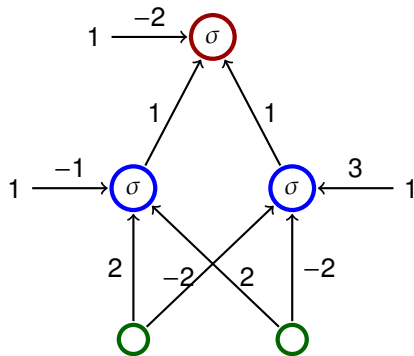
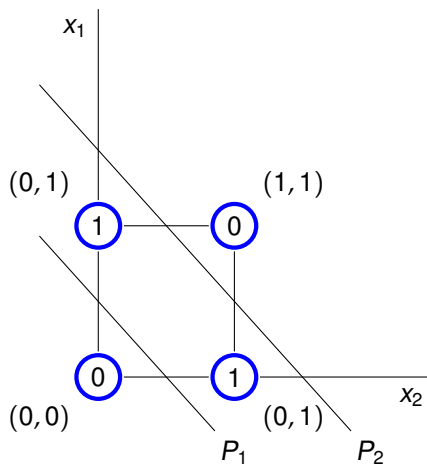
- ▶ Activation function is a unit step function

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

- ▶ 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



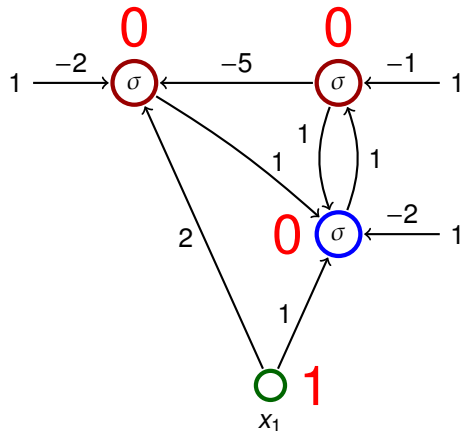
- ▶ The line  $P_1$  is given by  $-1 + 2x_1 + 2x_2 = 0$
- ▶ The line  $P_2$  is given by  $3 - 2x_1 - 2x_2 = 0$

## Activity – example

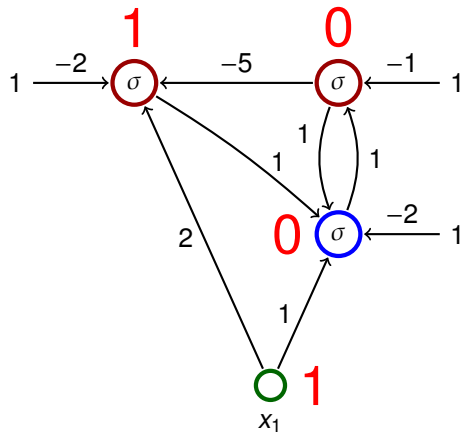
The activation function is the unit step function

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

The input is equal to 1



## Activity – example



The activation function is the unit step function

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

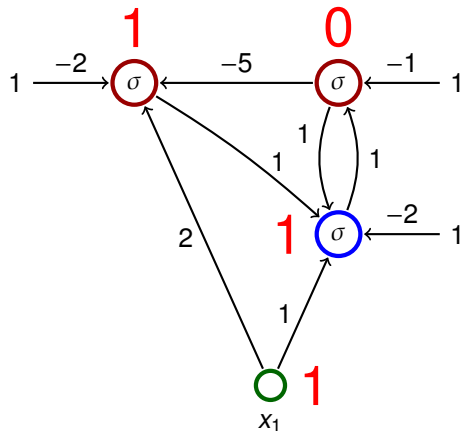
The input is equal to 1

## 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

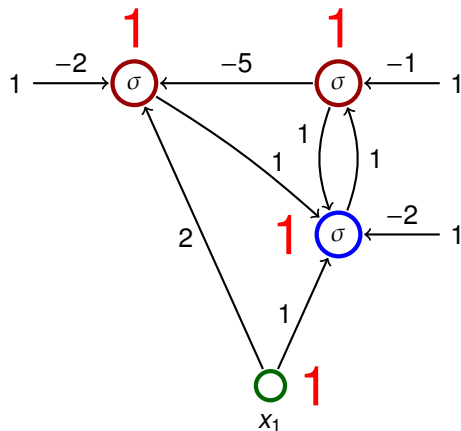


## Activity – example

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

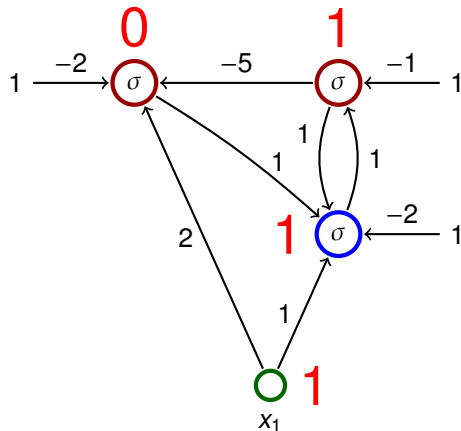


## Activity – example

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.

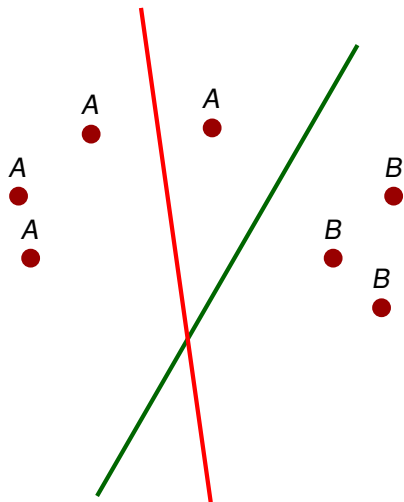
**Learning rule** for weight adaptation.

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

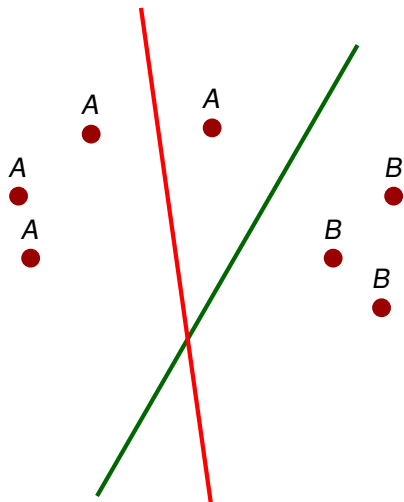
- ▶ Supervised learning
  - ▶ The desired function is described using *training examples* that are pairs of the form (input, output).
  - ▶ Learning algorithm searches for a configuration which "corresponds" to the training examples, typically by minimizing an error function.
- ▶ Unsupervised learning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

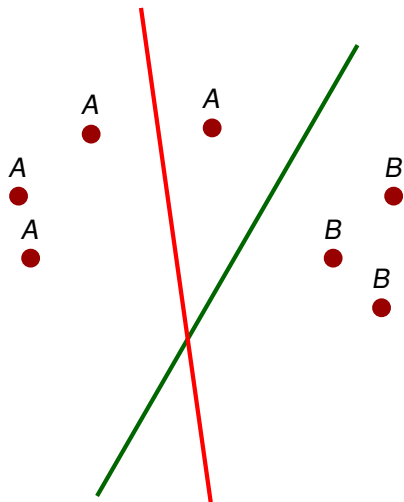


# Supervised learning – illustration



- ▶ 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

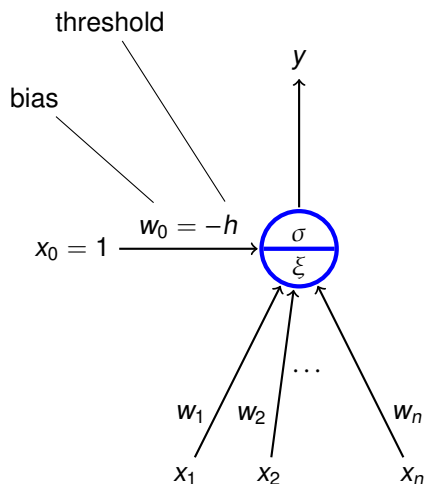
- ▶ Massive parallelism
  - ▶ neurons can be evaluated in parallel
- ▶ Learning
  - ▶ many sophisticated learning algorithms used to "program" neural networks
- ▶ generalization and robustness
  - ▶ information is encoded in a distributed manner in weights
  - ▶ "close" inputs typically get similar values

# 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*

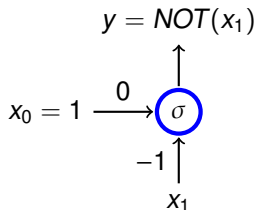
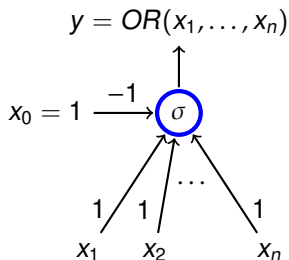
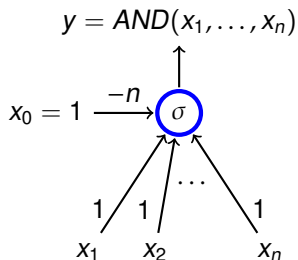
$$\sigma(\xi) = \begin{cases} 1 & \xi \geq 0; \\ 0 & \xi < 0. \end{cases}$$

# 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  $F : \{0, 1\}^n \rightarrow \{0, 1\}$ .*

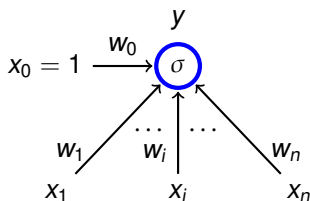
# Boolean functions

## 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  $\vec{v}$ :

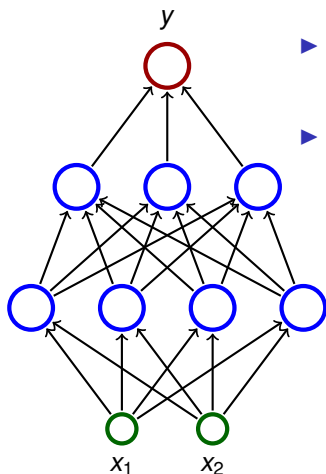


$$w_0 = -\sum_{i=1}^n v_i$$

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

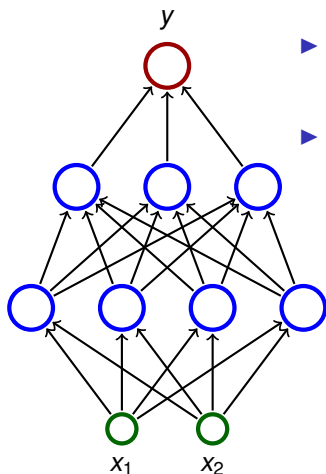
- ▶ Now let us connect all outputs of all neurons  $N_{\vec{v}}$  satisfying  $F(\vec{v}) = 1$  using a neuron implementing OR. □

# Non-linear separation



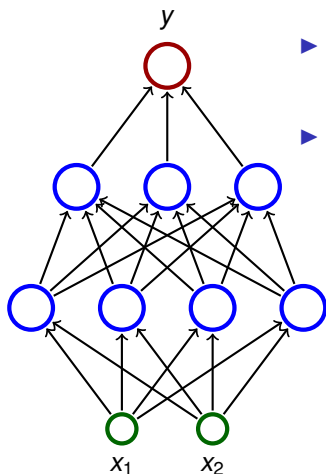
- ▶ Consider a three layer network; each neuron has the unit step activation function.
- ▶ The network divides the input space in two subspaces according to the output (0 or 1).

# Non-linear separation



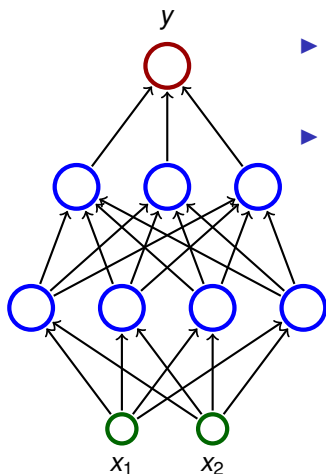
- ▶ Consider a three layer network; each neuron has the unit step activation function.
- ▶ The network divides the input space in two subspaces according to the output (0 or 1).
  - ▶ The first (hidden) layer divides the input space into half-spaces.

# Non-linear separation



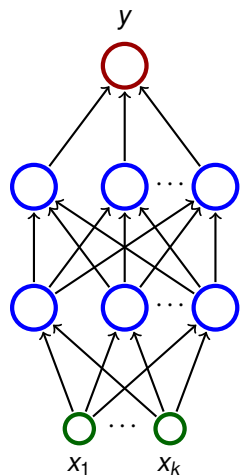
- ▶ Consider a three layer network; each neuron has the unit step activation function.
- ▶ The network divides the input space in two subspaces according to the output (0 or 1).
  - ▶ The first (hidden) layer divides the input space into half-spaces.
  - ▶ The second layer may e.g. make intersections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

# Non-linear separation



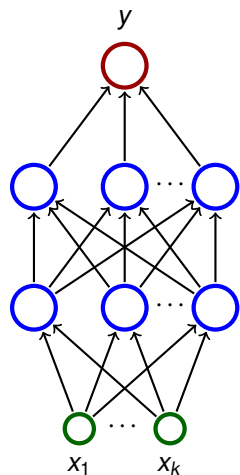
- ▶ Consider a three layer network; each neuron has the unit step activation function.
- ▶ The network divides the input space in two subspaces according to the output (0 or 1).
  - ▶ The first (hidden) layer divides the input space into half-spaces.
  - ▶ The second layer may e.g. make intersections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

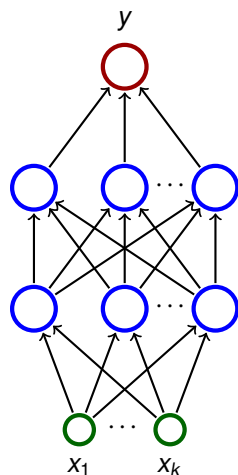
## Non-linear separation – illustration



- ▶ 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, ...)

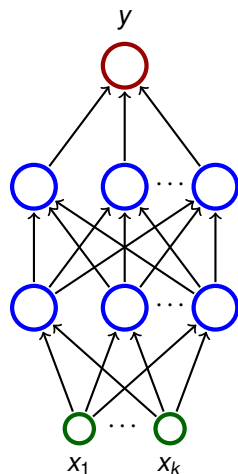


# Non-linear separation – illustration



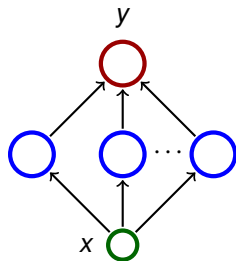
- ▶ Consider three layer networks; each neuron has the unit step activation function.
- ▶ Three layer nets are capable of "approximating" any "reasonable" subset  $A$  of the input space  $\mathbb{R}^k$ .
  - ▶ Cover  $A$  with hypercubes (in 2D squares, in 3D cubes, ...)
  - ▶ Each hypercube  $K$  can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

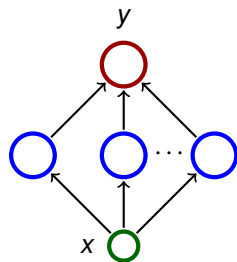
# Power of ReLU



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

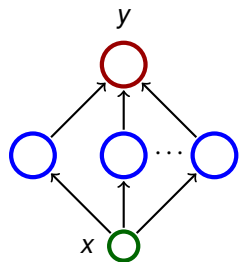


Consider a two layer network

- ▶ 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



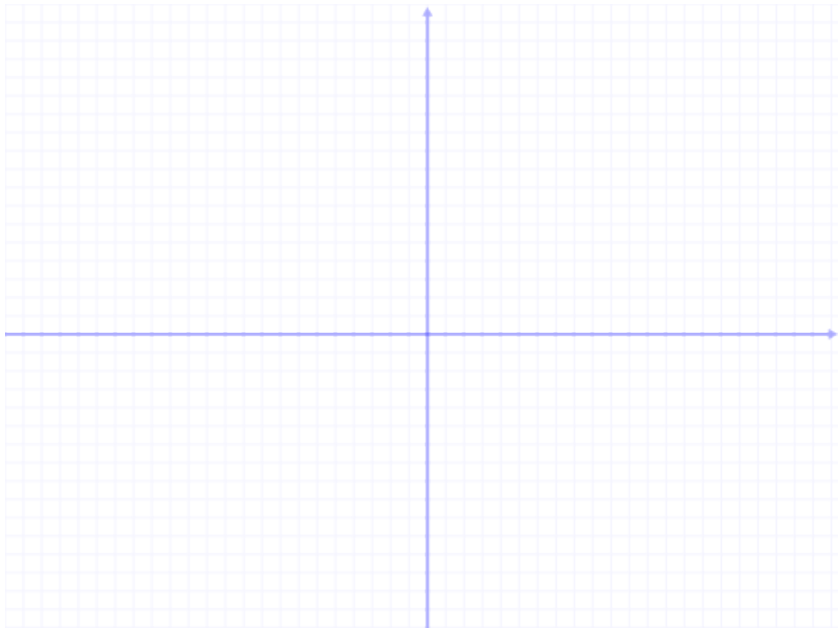
Consider a two layer network

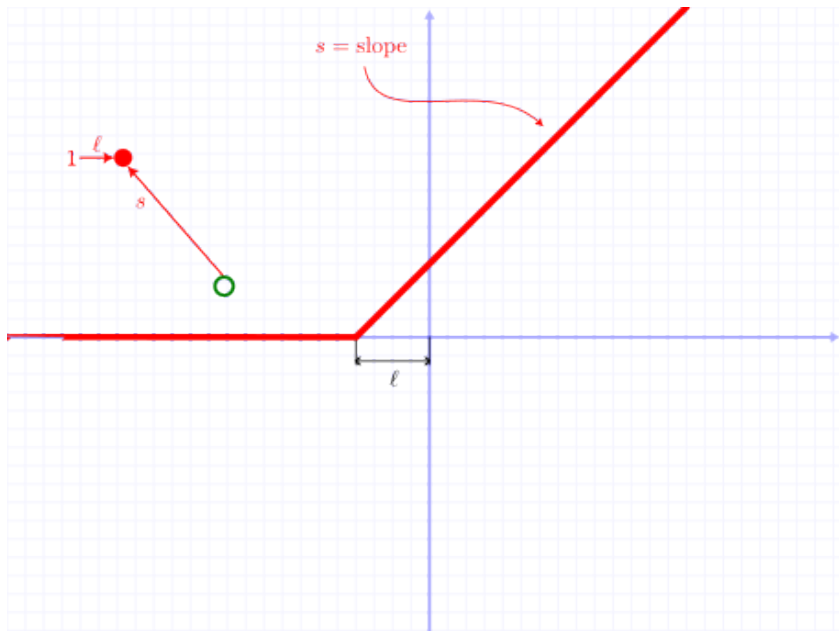
- ▶ with a single input and single output;
- ▶ hidden neurons with the ReLU activation:  
 $\sigma(\xi) = \max(\xi, 0)$ ;
- ▶ the output neuron with identity activation:  
 $\sigma(\xi) = \xi$  (linear model)

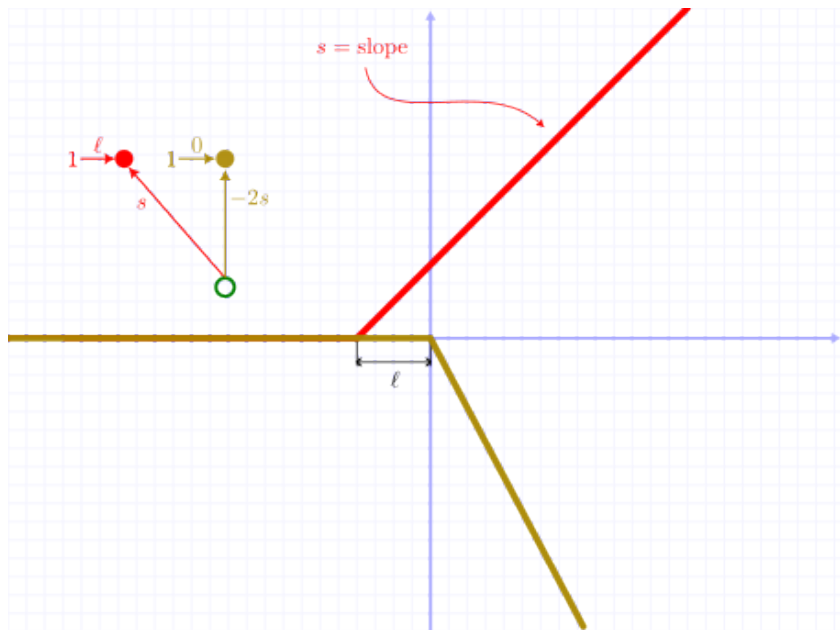
For every continuous function  $f : [0, 1] \rightarrow [0, 1]$  and  $\varepsilon > 0$  there is a network of the above type computing a function  $F : [0, 1] \rightarrow \mathbb{R}$  such that  $|f(x) - F(x)| \leq \varepsilon$  for all  $x \in [0, 1]$ .

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  $x$ .

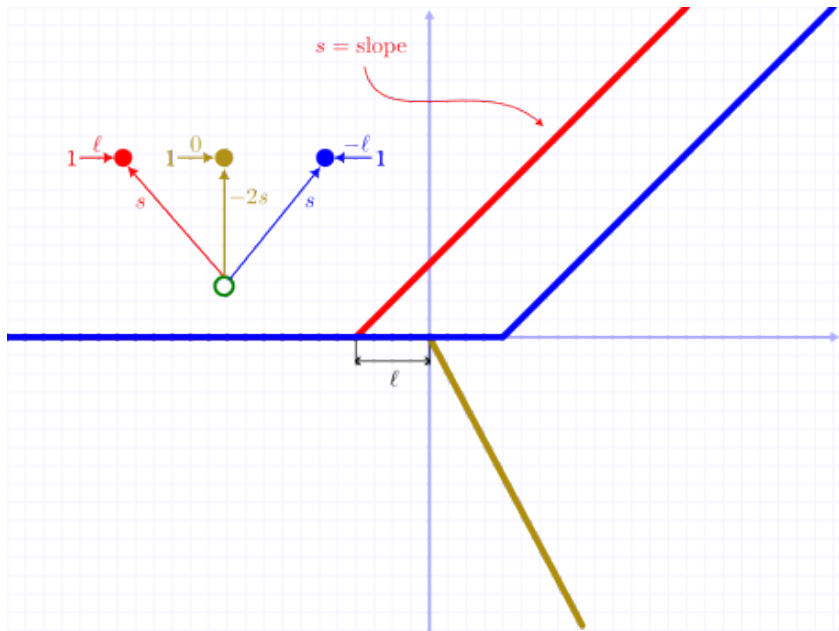
Just consider a continuous function  $f$  where  $f(x)$  is the minimum difference between  $x$  and a point on the boundary of  $A$ . Then uniformly approximate  $f$  using the networks.

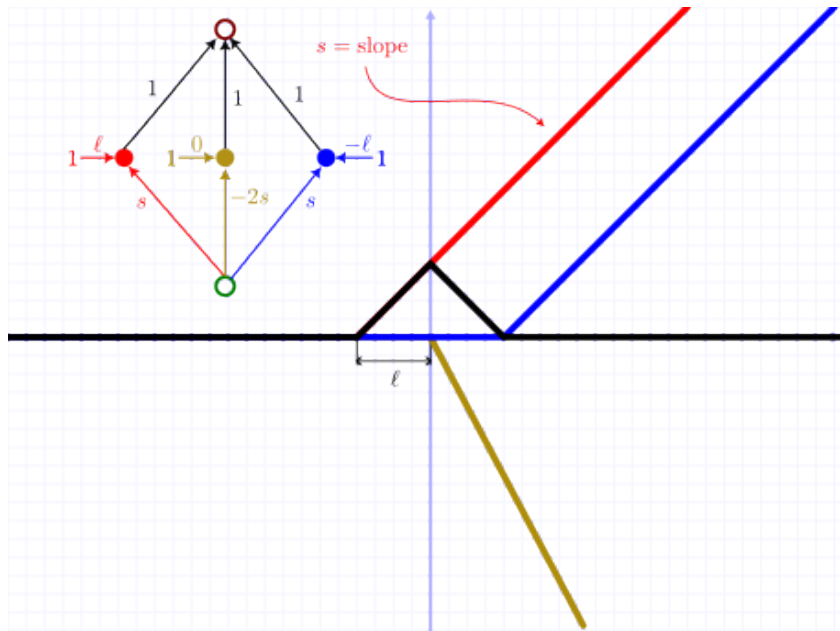


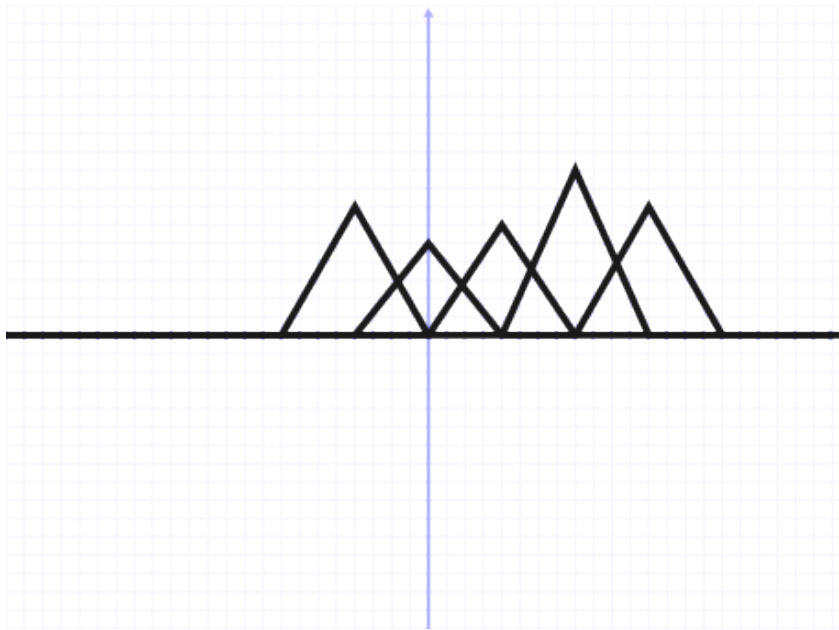


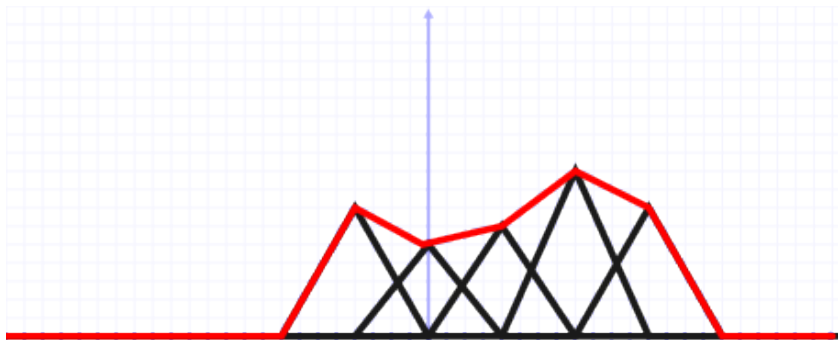












**Red** = sum of black

# Non-linear separation - sigmoid

## Theorem (Cybenko 1989 - informal version)

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

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

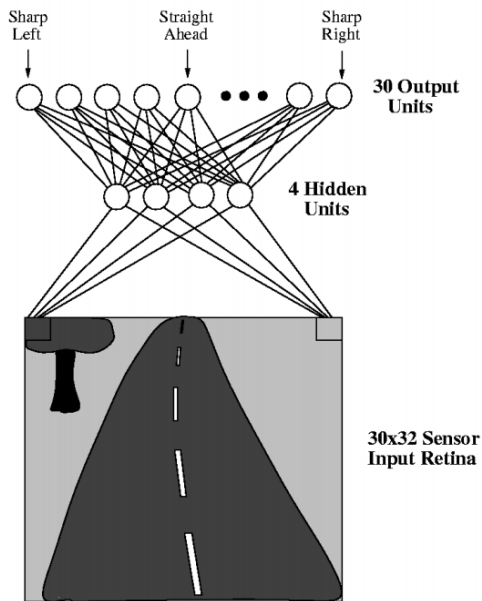
For every "reasonable" set  $A \subseteq [0, 1]^n$ , there is a **two layer network** where each hidden neuron has the activation function  $\sigma$  (output neurons are linear), that satisfies the following:

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:

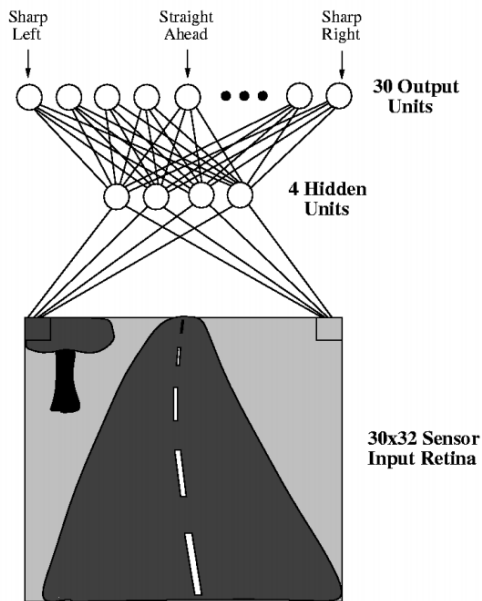
- ▶ "reasonable" means Lebesgue measurable
- ▶ "most" means that the set of incorrectly classified vectors has the Lebesgue measure smaller than a given  $\varepsilon > 0$

# Non-linear separation - practical illustration



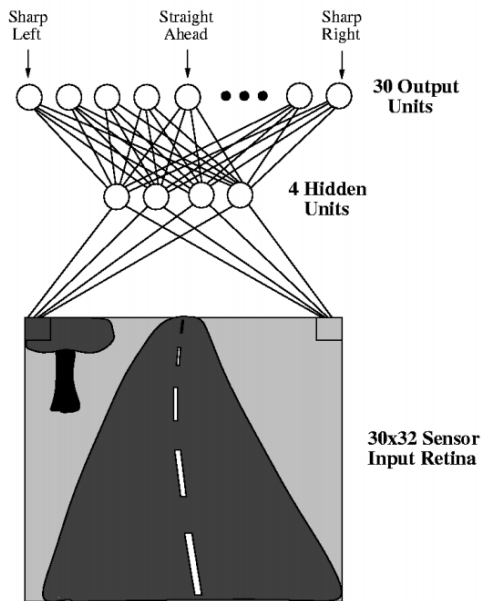
► ALVINN drives a car

# Non-linear separation - practical illustration



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

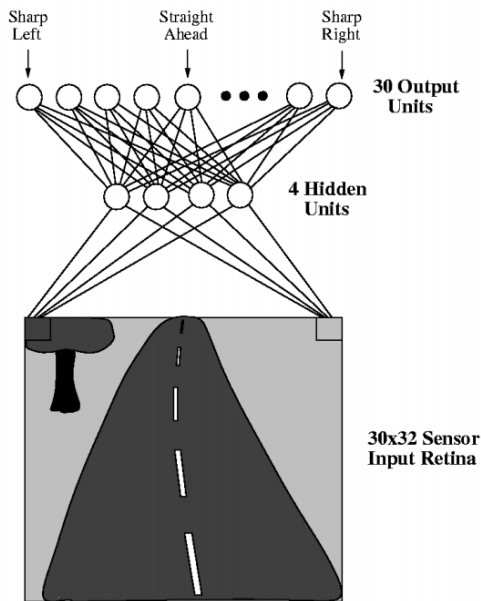
# Non-linear separation - practical illustration



- ▶ 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

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

# 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

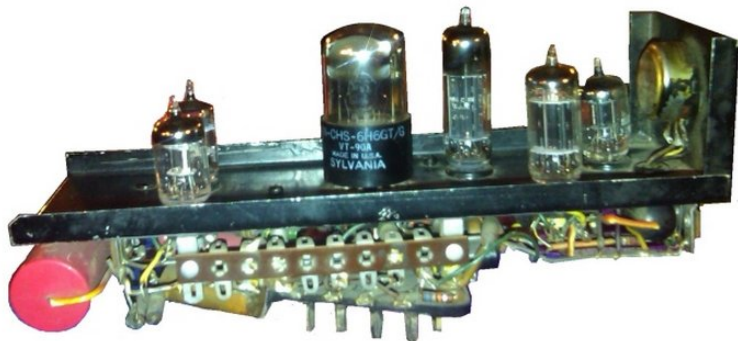
	Neural networks	"Classical" computers
Data	implicitly in weights	explicitly
Computation	naturally parallel	sequential, localized
Robustness	robust w.r.t. input corruption & damage	changing one bit may completely crash the computation
Precision	imprecise, network recalls a training example "similar" to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

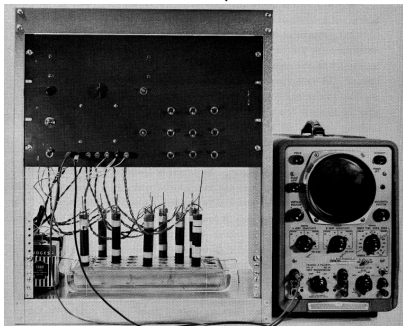
- ▶ 1957: Mark I Perceptron (Rosenblatt et al) - the first successful network for image recognition



- ▶ 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.

## Example

This slide was automatically generated by asking GPT-4 "Give me a beamer slide with complexity of Steepest descent, Neton's method and BFGS".

### Computational Complexity

Algorithm	Computational Complexity
Steepest Descent	$O(n)$ per iteration
Newton's Method	$O(n^3)$ to compute Hessian and solve system
BFGS	$O(n^2)$ to update Hessian approximation

Table: Summary of the computational complexity for each optimization algorithm.

- ▶ Steepest Descent: Simple but often slow, requiring many iterations.
- ▶ Newton's Method: Fast convergence but expensive per iteration.
- ▶ BFGS: Quasi-Newton, no Hessian needed, good speed and iteration count balance.

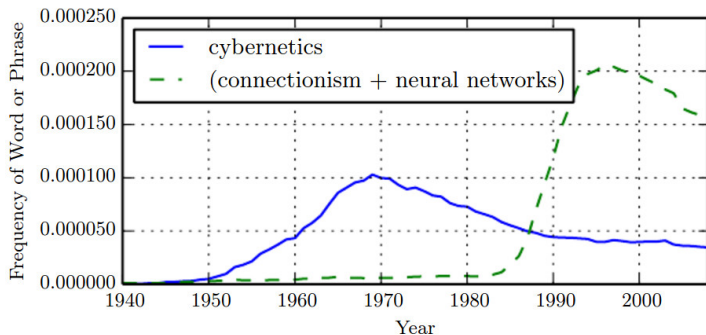
# Example Source

```
\begin{frame}{Computational Complexity}
\begin{table}
\begin{tabular}{l c}
\hline
\textbf{Algorithm} & \textbf{Computational Complexity} \\
\hline
Steepest Descent &  $O(n)$  per iteration \\
Newton's Method &  $O(n^3)$  to compute Hessian and solve system \\
BFGS &  $O(n^2)$  to update Hessian approximation \\
\hline
\end{tabular}
\caption{Summary of the computational complexity for each optimization algorithm.}
\end{table}

\begin{itemize}
\item Steepest Descent: Simple but often slow, requiring many iterations.
\item Newton's Method: Fast convergence but expensive per iteration.
\item BFGS: Quasi-Newton, no Hessian needed, good speed and iteration count balance.
\end{itemize}
\end{frame}
```



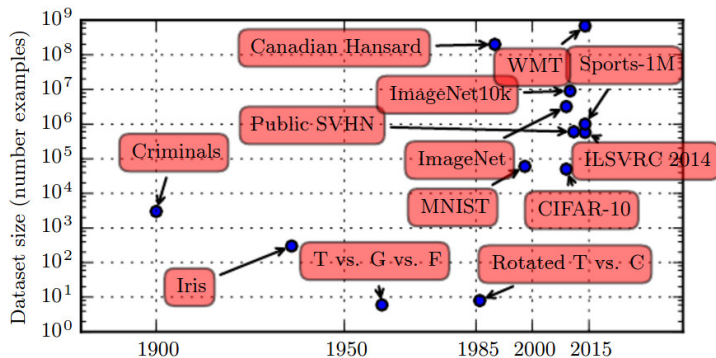
## History in waves ...



**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  $3.6 * 10^9$  images.

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

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

# GPT-3 Training Dataset

45 TB text data from multiple sources

<b><i>Dataset</i></b>	<b>Quantity (tokens)</b>	<b>Weight in training mix</b>	<b>Epochs elapsed when training for 300B tokens</b>
<i>Common Crawl (filtered)</i>	410 billion	60%	0.44
<i>WebText2</i>	19 billion	22%	2.9
<i>Books1</i>	12 billion	8%	1.9
<i>Books2</i>	55 billion	8%	0.43
<i>Wikipedia</i>	3 billion	3%	3.4

***Common Crawl*** corpus contains petabytes of data collected over 8 years of web crawling. The corpus contains raw web page data, metadata extracts and text extracts with light filtering.

***WebText2*** is the text of web pages from all outbound Reddit links from posts with 3+ upvotes.

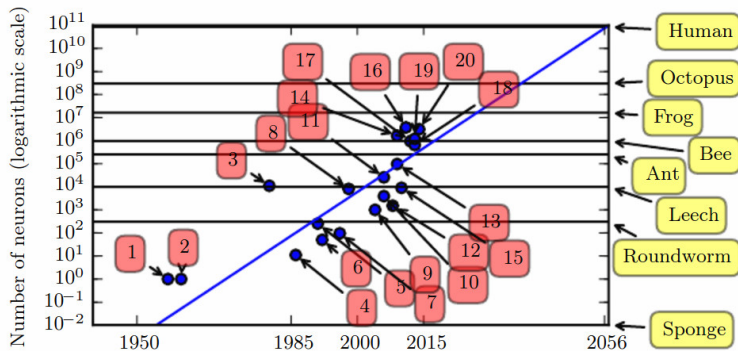
***Books1 & Books2*** are two internet-based books corpora.

***Wikipedia*** pages in the English language are also part of the training corpus.

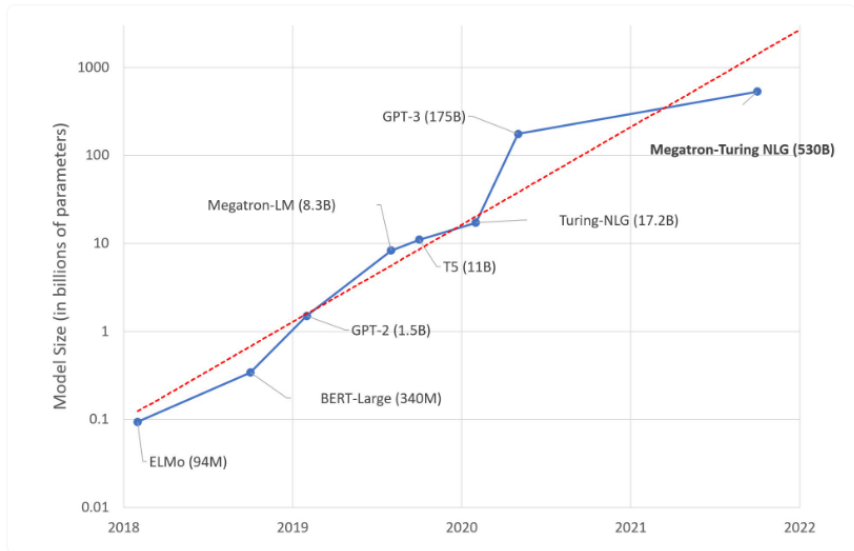
**Source:** Kindra Cooper. OpenAI GPT-3: Everything You Need to Know. Springboard. 2023

# Current hardware – What do we face?

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



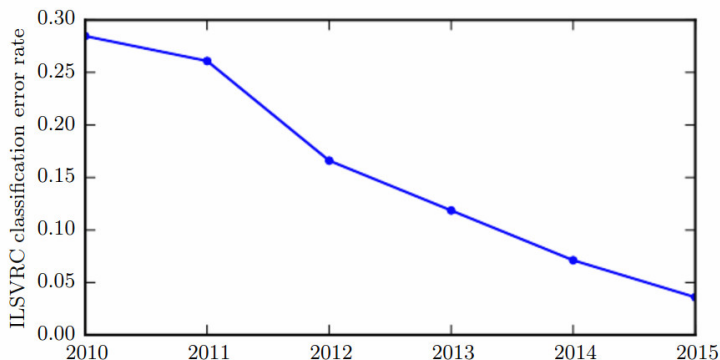
2. ADALINE
4. Early back-propagation network (Rumelhart et al., 1986b)
8. Image recognition: LeNet-5 (LeCun et al., 1998b)
10. Dimensionality reduction: Deep belief network (Hinton et al., 2006)  
... here the third "wave" of neural networks started
15. Digit recognition: GPU-accelerated multilayer perceptron (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 ...



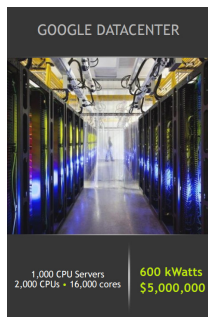
**Figure:** Since deep networks reached the scale necessary to compete in the ImageNet Large Scale Visual Recognition Challenge, they have consistently won the competition every year, and yielded lower and lower error rates each time. Data from Russakovsky et al. (2014b) and He et al. (2015).

# 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

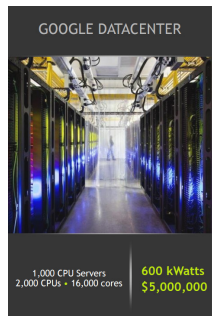
The task was image recognition (10 million youtube video frames)

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.



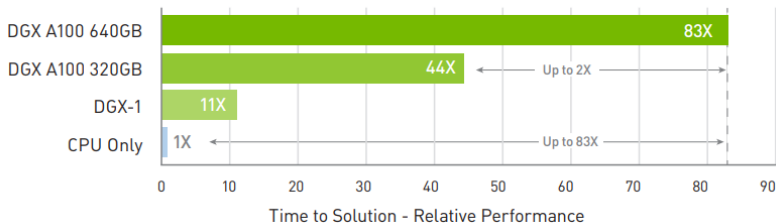


# Current hardware – NVIDIA DGX Station

- ▶ 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

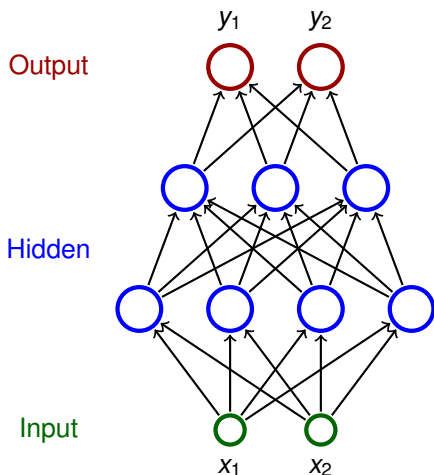
Most "mathematical" software packages contain some support of neural networks:

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

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

# 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., a 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 the numbers of neurons in individual layers (e.g., 2-4-3-2)

## 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*

## 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)

## 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

- ▶ inner potential of neuron  $j$ :

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

## MLP – 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)

# MLP – 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

- ▶ 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  $\mathbb{R}^{|X|} \rightarrow \mathbb{R}^{|Y|}$ . 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.

# MLP – 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  $(d_{kj})_{j \in Y}$ ).

# MLP – 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

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

This is just an example of an error function; we shall see other error functions later.

# 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

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

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

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

## Derivation of backprop.

Consider  $k = 1, \dots, p$  and a weight  $w_{ji}$ . By the chain rule:

$$\frac{\partial E_k}{\partial w_{ji}} =$$

## Derivation of backprop.

Consider  $k = 1, \dots, p$  and a weight  $w_{ji}$ . By the chain rule:

$$\frac{\partial E_k}{\partial w_{ji}} = \frac{\partial E_k}{\partial y_j} \cdot \frac{\partial y_j}{\partial w_{ji}} =$$



## Derivation of backprop.

Consider  $k = 1, \dots, p$  and a weight  $w_{ji}$ . By the chain rule:

$$\frac{\partial E_k}{\partial w_{ji}} = \frac{\partial E_k}{\partial y_j} \cdot \frac{\partial y_j}{\partial w_{ji}} = \frac{\partial E_k}{\partial y_j} \cdot \frac{\partial y_j}{\partial \xi_j} \cdot \frac{\partial \xi_j}{\partial w_{ji}} =$$

## Derivation of backprop.

Consider  $k = 1, \dots, p$  and a weight  $w_{ji}$ . By the chain rule:

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

since

$$\frac{\partial y_j}{\partial \xi_j} = \frac{\partial(\sigma_j(\xi_j))}{\partial \xi_j} = \sigma'_j(\xi_j)$$

$$\frac{\partial \xi_j}{\partial w_{ji}} = \frac{\partial(\sum_{r \in j_{\leftarrow}} w_{jr} y_r)}{\partial w_{ji}} = y_i$$

## Derivation of backdrop. (cont.)

$$\text{For } j \in Y: \frac{\partial E_k}{\partial y_j} = \frac{\partial \left( \frac{1}{2} \sum_{r \in Y} (y_r - d_{kr})^2 \right)}{\partial y_j} = y_j - d_{kj}$$

## Derivation of backdrop. (cont.)

$$\text{For } j \in Y: \frac{\partial E_k}{\partial y_j} = \frac{\partial \left( \frac{1}{2} \sum_{r \in Y} (y_r - d_{kr})^2 \right)}{\partial y_j} = y_j - d_{kj}$$

... and another application of the chain rule:

$$\text{For } j \notin Y: \frac{\partial E_k}{\partial y_j} =$$

## Derivation of backdrop. (cont.)

$$\text{For } j \in Y: \frac{\partial E_k}{\partial y_j} = \frac{\partial \left( \frac{1}{2} \sum_{r \in Y} (y_r - d_{kr})^2 \right)}{\partial y_j} = y_j - d_{kj}$$

... and another application of the chain rule:

$$\text{For } j \notin Y: \frac{\partial E_k}{\partial y_j} = \sum_{r \in j^*} \frac{\partial E_k}{\partial y_r} \cdot \frac{\partial y_r}{\partial y_j} =$$

## Derivation of backdrop. (cont.)

$$\text{For } j \in Y: \frac{\partial E_k}{\partial y_j} = \frac{\partial \left( \frac{1}{2} \sum_{r \in Y} (y_r - d_{kr})^2 \right)}{\partial y_j} = y_j - d_{kj}$$

... and another application of the chain rule:

$$\begin{aligned} \text{For } j \notin Y: \frac{\partial E_k}{\partial y_j} &= \sum_{r \in j \rightarrow} \frac{\partial E_k}{\partial y_r} \cdot \frac{\partial y_r}{\partial y_j} = \sum_{r \in j \rightarrow} \frac{\partial E_k}{\partial y_r} \cdot \frac{\partial y_r}{\partial \xi_r} \cdot \frac{\partial \xi_r}{\partial y_j} \\ &= \end{aligned}$$

## Derivation of backdrop. (cont.)

$$\text{For } j \in Y: \frac{\partial E_k}{\partial y_j} = \frac{\partial \left( \frac{1}{2} \sum_{r \in Y} (y_r - d_{kr})^2 \right)}{\partial y_j} = y_j - d_{kj}$$

... and another application of the chain rule:

$$\begin{aligned} \text{For } j \notin Y: \frac{\partial E_k}{\partial y_j} &= \sum_{r \in j \rightarrow} \frac{\partial E_k}{\partial y_r} \cdot \frac{\partial y_r}{\partial y_j} = \sum_{r \in j \rightarrow} \frac{\partial E_k}{\partial y_r} \cdot \frac{\partial y_r}{\partial \xi_r} \cdot \frac{\partial \xi_r}{\partial y_j} \\ &= \sum_{r \in j \rightarrow} \frac{\partial E_k}{\partial y_r} \cdot \sigma'_r(\xi_r) \cdot w_{rj} \end{aligned}$$

since

$$\frac{\partial y_r}{\partial \xi_r} = \frac{\partial(\sigma_r(\xi_r))}{\partial \xi_r} = \sigma'_r(\xi_r)$$

$$\frac{\partial \xi_r}{\partial y_j} = \frac{\partial \left( \sum_{s \in r \leftarrow} w_{rs} y_s \right)}{\partial y_j} = w_{rj}$$

## 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  $j \in Z \setminus X$ :

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

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

# 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  $\frac{\partial E}{\partial w_{ji}}(\vec{w}^{(t-1)})$  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  $\xi_r$ )

## 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 w.r.t. the number of network weights.

# 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  $\xi_r$ )

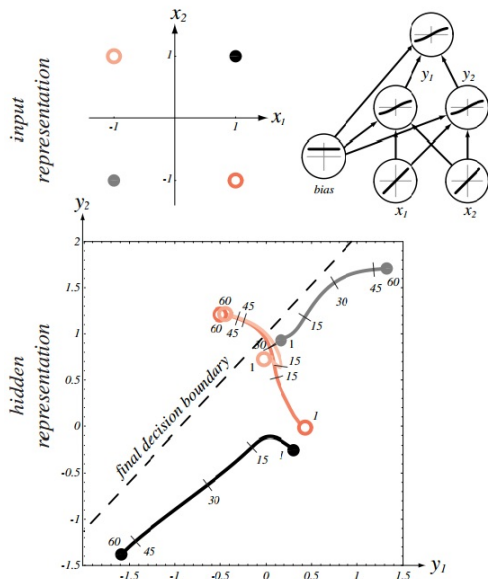
**Proof sketch:** The algorithm does the following  $p$  times:

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 w.r.t. the number of network weights.

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 the 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.

## Regression: Output and Error

- ▶ For **regression**, the output activation is typically the identity, i.e.,  $y_i = \sigma(\xi_i) = \xi_i$  for  $i \in Y$ .

## Regression: Output and Error

- ▶ For **regression**, the output activation is typically the identity, i.e.,  $y_i = \sigma(\xi_i) = \xi_i$  for  $i \in Y$ .
- ▶ A training dataset

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

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  $(d_{ki})_{i \in Y}$ ).

## Regression: Output and Error

- ▶ For **regression**, the output activation is typically the identity, i.e.,  $y_i = \sigma(\xi_i) = \xi_i$  for  $i \in Y$ .
- ▶ A training dataset

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

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

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



# Maximum Likelihood vs Least Squares

**Fix a training set**  $D = \{(x_1, d_1), (x_2, d_2), \dots, (x_p, d_p)\}$ ,  $d_k \in \mathbb{R}$ .

Consider a single output neuron  $o$ .

# Maximum Likelihood vs Least Squares

**Fix a training set**  $D = \{(x_1, d_1), (x_2, d_2), \dots, (x_p, d_p)\}$ ,  $d_k \in \mathbb{R}$ .

Consider a single output neuron  $o$ .

Assume that each  $d_k$  was generated randomly as follows

$$d_k = y_o(\vec{w}, \vec{x}_k) + \epsilon_k$$

- ▶  $\vec{w}$  are **unknown constants**
- ▶  $\epsilon_k$  are normally distributed with mean 0 and an unknown variance  $\sigma^2$

# Maximum Likelihood vs Least Squares

Fix a training set  $D = \{(x_1, d_1), (x_2, d_2), \dots, (x_p, d_p)\}$ ,  $d_k \in \mathbb{R}$ .

Consider a single output neuron  $o$ .

Assume that each  $d_k$  was generated randomly as follows

$$d_k = y_o(\vec{w}, \vec{x}_k) + \epsilon_k$$

- ▶  $\vec{w}$  are **unknown constants**
- ▶  $\epsilon_k$  are normally distributed with mean 0 and an unknown variance  $\sigma^2$

Assume that  $\epsilon_1, \dots, \epsilon_p$  have been generated **independently**.

# Maximum Likelihood vs Least Squares

Fix a training set  $D = \{(x_1, d_1), (x_2, d_2), \dots, (x_p, d_p)\}$ ,  $d_k \in \mathbb{R}$ .

Consider a single output neuron  $o$ .

Assume that each  $d_k$  was generated randomly as follows

$$d_k = y_o(\vec{w}, \vec{x}_k) + \epsilon_k$$

- ▶  $\vec{w}$  are **unknown constants**
- ▶  $\epsilon_k$  are normally distributed with mean 0 and an unknown variance  $\sigma^2$

Assume that  $\epsilon_1, \dots, \epsilon_p$  have been generated **independently**.

Denote by  $p(d_1, \dots, d_p \mid \vec{w}, \sigma^2)$  the probability density of the values  $d_1, \dots, d_p$  assuming fixed  $x_1, \dots, x_p, \vec{w}, \sigma^2$ .

(For the interested: The independence and definition of  $d_k$ 's imply

$$p(d_1, \dots, d_p \mid \vec{w}, \sigma^2) = \prod_{k=1}^p N[y_o(\vec{w}, \vec{x}_k), \sigma^2](d_k)$$

$N[y_o(\vec{w}, \vec{x}_k), \sigma^2](d_k)$  is a normal dist. with the mean  $y_o(\vec{w}, \vec{x}_k)$  and var.  $\sigma^2$ .)

# Maximum Likelihood vs Least Squares

Our goal is to find the weights  $\vec{w}$  that maximize the likelihood

$$L(\vec{w}, \sigma^2) := p(d_1, \dots, d_p \mid \vec{w}, \sigma^2)$$

But now with the **fixed** values  $d_1, \dots, d_n$  from the training set!

# Maximum Likelihood vs Least Squares

Our goal is to find the weights  $\vec{w}$  that maximize the likelihood

$$L(\vec{w}, \sigma^2) := p(d_1, \dots, d_p \mid \vec{w}, \sigma^2)$$

But now with the **fixed** values  $d_1, \dots, d_n$  from the training set!

## Theorem

*The unique  $\vec{w}$  that minimize the least squares error  $E[\vec{w}]$  maximize  $L(\vec{w}, \sigma^2)$  for an arbitrary variance  $\sigma^2$ .*

## Classification: Output and Error

- ▶ The output activation function *softmax*:

$$y_i = \sigma_i(\xi_{j_1}, \dots, \xi_{j_k}) = \frac{e^{\xi_i}}{\sum_{j \in Y} e^{\xi_j}} \quad \text{Here } Y = \{j_1, \dots, j_k\}$$

# Classification: Output and Error

- ▶ The output activation function *softmax*:

$$y_i = \sigma_i(\xi_{j_1}, \dots, \xi_{j_k}) = \frac{e^{\xi_i}}{\sum_{j \in Y} e^{\xi_j}} \quad \text{Here } Y = \{j_1, \dots, j_k\}$$

- ▶ A training dataset

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

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  $(d_{ki})_{i \in Y}$ ).



# Classification: Output and Error

- ▶ The output activation function *softmax*:

$$y_i = \sigma_i(\xi_{j_1}, \dots, \xi_{j_k}) = \frac{e^{\xi_i}}{\sum_{j \in Y} e^{\xi_j}} \quad \text{Here } Y = \{j_1, \dots, j_k\}$$

- ▶ A training dataset

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

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  $\ell \in V$ .

## Binary Classification: Output and Error

Assume a single output neuron  $o \in Y = \{0\}$ .

- ▶ The output activation function *logistic sigmoid*:

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

# Binary Classification: Output and Error

Assume a single output neuron  $o \in Y = \{0\}$ .

- ▶ 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.

# Binary Classification: Output and Error

Assume a single output neuron  $o \in Y = \{0\}$ .

- ▶ 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*:

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

## Cross-entropy vs max likelihood

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

## Cross-entropy vs max likelihood

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.

## Cross-entropy vs max likelihood

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*:

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

$\log(L) =$

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

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

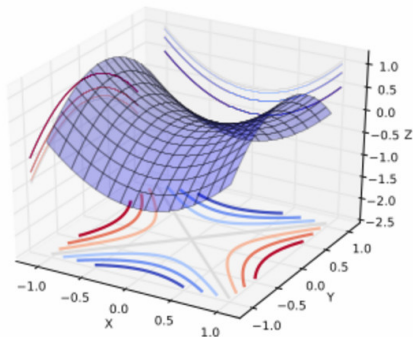
- ▶ Small networks: Lots of local minima where the descent gets stuck.
- ▶ The model identifiability problem: Swapping incoming weights of neurons  $i$  and  $j$  leaves the same network topology – **weight space symmetry**.
- ▶ Recent studies show that for sufficiently large networks, all local minima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

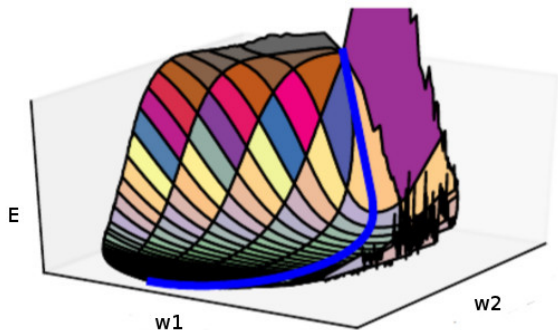
## Saddle points

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



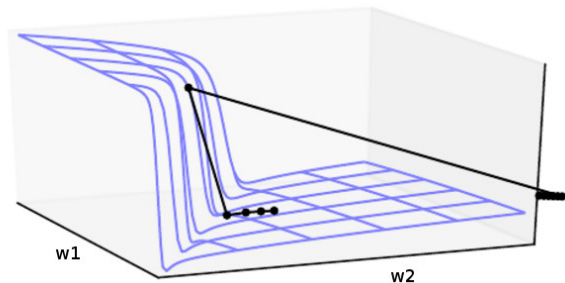
# Issues in gradient descent – too slow descent

- ▶ flat regions

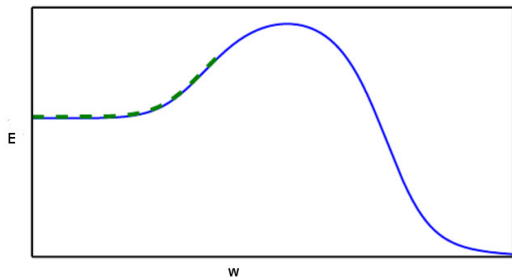


# Issues in gradient descent – too fast descent

- ▶ 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 a linear rate.*

Later, we get to a special type of network 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. The typical power of 2 batch sizes ranges 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. The typical power of 2 batch sizes ranges 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

The issue in the gradient descent:

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



# Momentum

The issue in the gradient descent:

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



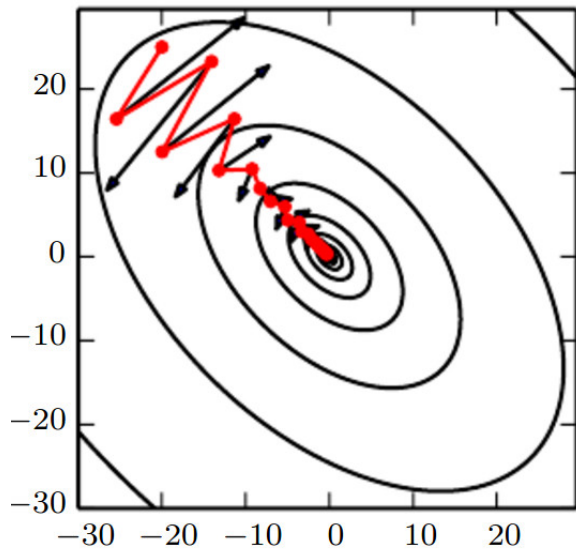
**Solution:** In every step, add the change made in the previous step (weighted by a factor  $\alpha$ ):

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

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

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

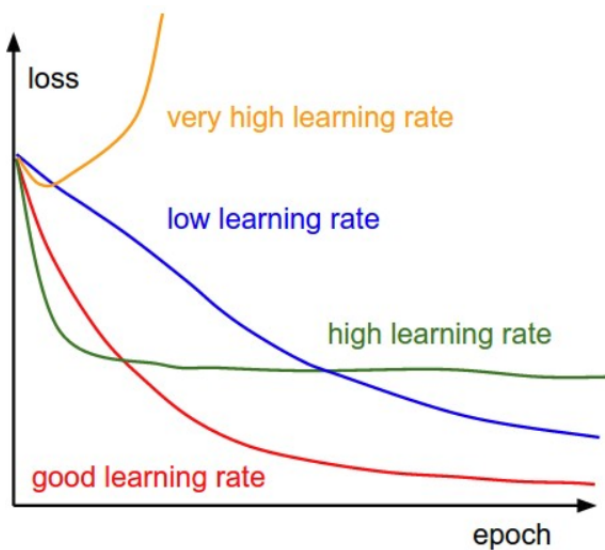
where

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

- ▶  $0 < \varepsilon(t) \leq 1$  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 constant 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 constant 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 constant 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 constant 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 the CIFAR10 dataset).

For a 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 rates 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 gradients throughout the 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 and Adam (RMSProp combined with momentum).

A natural question: Which algorithm should one choose?

## Other optimization methods

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

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 and Adam (RMSProp combined with momentum).

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.

# Input preprocessing

- ▶ Some inputs may be much larger than others.

For example, the height vs. weight of a person, the max. speed of a car (in km/h) vs. its price (in CZK), etc.

# Input preprocessing

- ▶ Some inputs may be much larger than others.

For example, the height vs. weight of a person, the max. speed of a car (in km/h) vs. its price (in CZK), etc.

- ▶ Large inputs have a greater influence on the training than the small ones. Also, too large inputs may slow down learning (saturation of some activation functions).

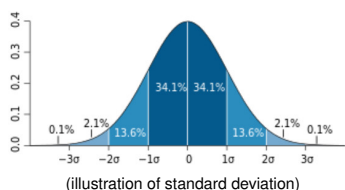
# Input preprocessing

- ▶ Some inputs may be much larger than others.

For example, the height vs. weight of a person, the max. speed of a car (in km/h) vs. its price (in CZK), etc.

- ▶ Large inputs have a greater influence on the training than the small ones. Also, too large inputs may slow down learning (saturation of some 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 the data (the training set).



# Initial weights - intuition

- ▶ Assume weights are chosen randomly. What distribution?

# Initial weights - intuition

- ▶ Assume weights are chosen randomly. What distribution?

Consider the behavior of a deep network:

- ▶ Small weights make the values of inner potentials vanish.
- ▶ Large weights make the values of inner potentials explode.

Hence, we want to choose weights so that the inner potentials of neurons are stable (similar in all layers of the network).

## 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 are 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 are 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 are 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  $\sigma_j = 1$ .

- ▶ 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  $\text{Var}[\xi] = n \cdot w^2$  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  $\text{Var}[\xi] = n \cdot w^2$  as follows:

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

## 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  $\text{Var}[\xi] = n \cdot w^2$  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  $\text{Var}[\xi] = n \cdot w^2$  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 heuristic for weight initialization ignores the variance of the gradient (i.e., it is concerned only with the "size" of activations in the forward pass).

## Normal Glorot initialization

The previous heuristic for weight initialization ignores the 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 heuristic for weight initialization ignores the 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 are 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.

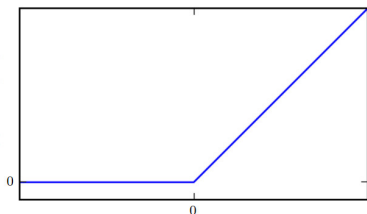
We obtain *uniform Glorot initialization* (aka *uniform Xavier init.*):

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

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 activation functions. Most prominent is ReLU:



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

- ▶ THE default activation function recommended for 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 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 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 (green board):
  - ▶ Generalizes ReLU, not dead for negative potentials.
  - ▶ Experimentally not much better than ReLU.

## More modern activation functions

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

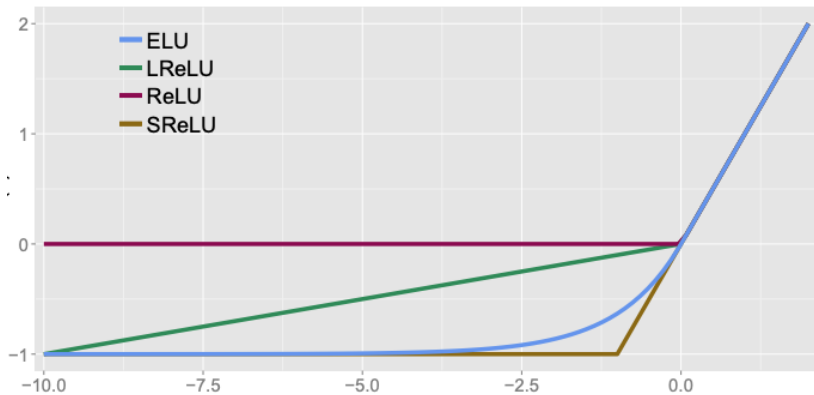
$$\sigma(\xi) = \begin{cases} \alpha(\exp(\xi) - 1) & \text{for } \xi < 0 \\ \xi & \text{for } \xi \geq 0 \end{cases}$$

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

$$\vec{z}_k = (z_{k1}, \dots, z_{kn})$$

## 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

$$\vec{z}_k = (z_{k1}, \dots, z_{kn})$$

- ▶ Set all components of all vectors  $\vec{z}_k$  to the mean = 0 and the variance = 1 and obtain *normalized vectors*:  $\hat{z}_1, \dots, \hat{z}_p$ .



## 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.

# Normalization

During the **training**, the normalized vectors  $\hat{z}_1, \dots, \hat{z}_p$  are computed as follows:

$$\hat{z}_{ki} = \frac{z_{ki} - \mu_i}{s_i}$$

Here

$$\mu_i = \frac{1}{p} \sum_{k=1}^p z_{ki}$$

$$s_i = \sqrt{\frac{1}{p} \sum_{k=1}^p (z_{ki} - \mu_i)^2}$$

During **inference**, where we have just a single value  $\vec{z}$  of the layer  $\ell$  for an input  $\vec{x}$ , we use  $\mu_i$  and  $s_i$  estimated on a population (e.g., a larger sample of the training set).

# Generalization

# 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

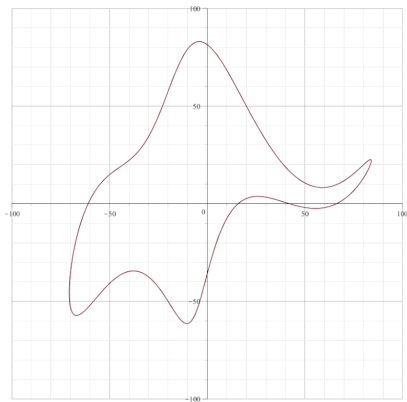
Regularization is a big issue in neural networks, as they typically use a huge amount of parameters and thus are very susceptible to overfitting.

# Regularization

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."**

# Elephant



$$x(t) = -60 \cos(t) + 30 \sin(t) - 8 \sin(2t) + 10 \sin(3t)$$

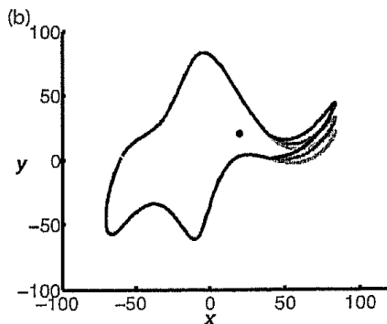
$$y(t) = 50 \sin(t) + 18 \sin(2t) - 12 \cos(3t) + 14 \cos(5t)$$

The four parameters are complex numbers (e.g.,  $-60 + 50i$ ).

Mayer, Jurgen; Khairy, Khaled; Howard, Jonathon (May 12, 2010). "Drawing an elephant with four complex parameters". American Journal of Physics. 78 (6)



# Fifth Elephant



Parameter	Real part	Imaginary part
$p_1 = 50 - 30i$	$B_1^x = 50$	$B_1^y = -30$
$p_2 = 18 + 8i$	$B_2^x = 18$	$B_2^y = 8$
$p_3 = 12 - 10i$	$A_3^x = 12$	$B_3^y = -10$
$p_4 = -14 - 60i$	$A_3^x = -14$	$A_1^y = -60$
$p_5 = 40 + 20i$	Wiggle coeff. = 40	$x_{cyc} = y_{cyc} = 20$

# 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.

(some 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.

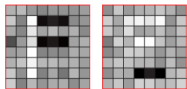
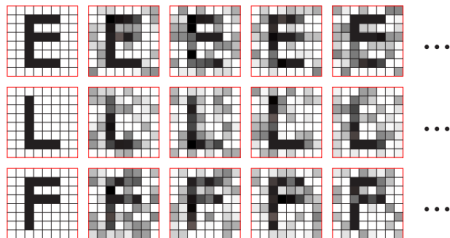
Experiment, experiment, experiment.

# 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.).

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.

Some application(s)

## MLP applications

- ▶ MLP is the basic network used when it is unclear what topology is appropriate.

# MLP applications

- ▶ MLP is the basic network used when it is unclear what topology is appropriate.
- ▶ MLP is often used as a component of larger deep models (especially at the head of the network, interpreting the features extracted by the lower, more specialized layers.)

# MLP applications

- ▶ MLP is the basic network used when it is unclear what topology is appropriate.
- ▶ MLP is often used as a component of larger deep models (especially at the head of the network, interpreting the features extracted by the lower, more specialized layers.)
- ▶ The most prominent NN architecture in the 80s and 90s.

# MLP applications

- ▶ MLP is the basic network used when it is unclear what topology is appropriate.
- ▶ MLP is often used as a component of larger deep models (especially at the head of the network, interpreting the features extracted by the lower, more specialized layers.)
- ▶ The most prominent NN architecture in the 80s and 90s.
- ▶ Various applications:
  - ▶ ALVINN - autonomous driving car
  - ▶ Characters/digits recognition
  - ▶ Table data processing
  - ▶ ...



# MLP applications

- ▶ MLP is the basic network used when it is unclear what topology is appropriate.
- ▶ MLP is often used as a component of larger deep models (especially at the head of the network, interpreting the features extracted by the lower, more specialized layers.)
- ▶ The most prominent NN architecture in the 80s and 90s.
- ▶ Various applications:
  - ▶ ALVINN - autonomous driving car
  - ▶ Characters/digits recognition
  - ▶ Table data processing
  - ▶ ...
- ▶ Two-layer MLPs are usually not much better than non-neural models.

# MNIST – handwritten digits recognition

- ▶ Database of labeled images of handwritten digits: 60 000 training examples, 10 000 testing.
- ▶ Dimensions: 28 x 28, digits are centered to the "center of gravity" of pixel values and normalized to a fixed size.
- ▶ More at <http://yann.lecun.com/exdb/mnist/>

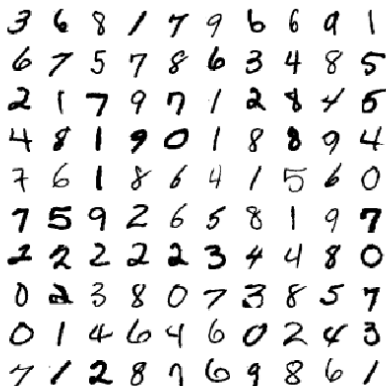


Fig. 4. Size-normalized examples from the MNIST database.

The database is used as a standard benchmark in lots of publications.

# MNIST – handwritten digits recognition

- ▶ Database of labeled images of handwritten digits: 60 000 training examples, 10 000 testing.
- ▶ Dimensions: 28 x 28, digits are centered to the "center of gravity" of pixel values and normalized to a fixed size.
- ▶ More at <http://yann.lecun.com/exdb/mnist/>

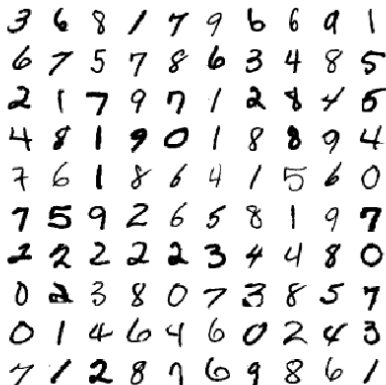


Fig. 4. Size-normalized examples from the MNIST database.

The database is used as a standard benchmark in lots of publications.

Allows comparison of various methods.

One of the best "old" results is the following:

6-layer NN 784-2500-2000-1500-1000-500-10 (on GPU)  
(Ciresan et al. 2010)

**Abstract:** Good old on-line back-propagation for plain multi-layer perceptrons yields a very low 0.35 error rate on the famous MNIST handwritten digits benchmark. All we need to achieve this best result so far are many hidden layers, many neurons per layer, numerous deformed training images, and graphics cards to greatly speed up learning.

A famous application of a learning **convolutional** network LeNet-1 in 1998.

Y. LeCun, L. Bottou, Y. Bengio, and P. Haffner. Gradient-based learning applied to document recognition. Proceedings of the IEEE, 1998

# MNIST – LeNet1

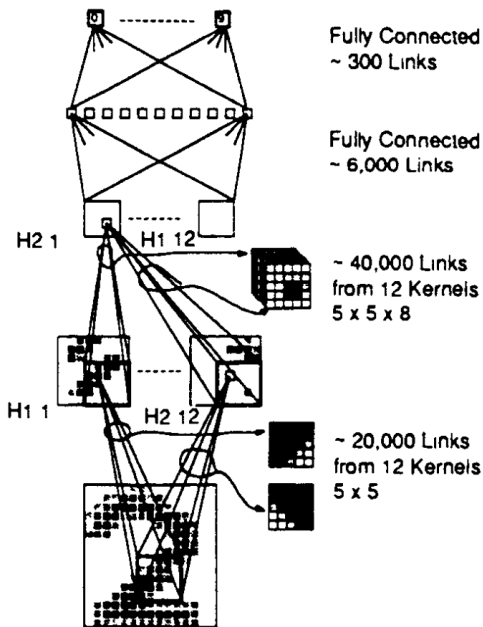
10 Output Units

Layer H3  
30 Hidden Units

Layer H2  
 $12 \times 16 = 192$   
Hidden Units

Layer H1  
 $12 \times 64 = 768$   
Hidden Units

256 Input Units



Interpretation of output:

- ▶ the output neuron with the highest value identifies the digit.
- ▶ the same, but if the two largest neuron values are too close together, the input is rejected (i.e. no answer).

## **Learning:**

Inputs:

- ▶ training on 7291 samples, tested on 2007 samples

## **Results:**

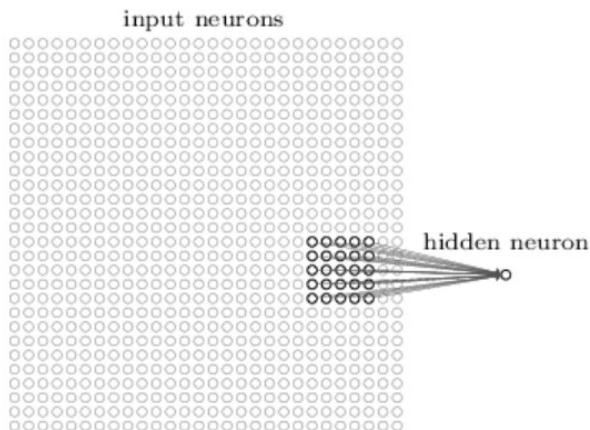
- ▶ error on test set without rejection: 5%
- ▶ error on test set with rejection: 1% (12% rejected)
  
- ▶ compare with dense MLP with 40 hidden neurons: error 1% (19.4% rejected)

# Convolutional Networks

Some parts of the lecture are based on the online book Neural Networks and Deep Learning by Michael Nielsen.

<http://neuralnetworksanddeeplearning.com/index.html>

# Convolutional networks - local receptive fields



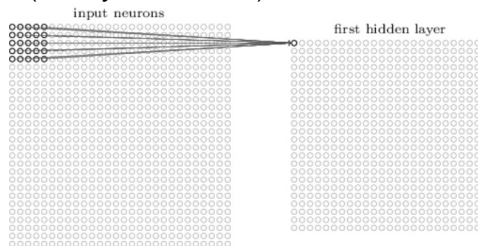
Every neuron is connected with a field of  $k \times k$  (in this case  $5 \times 5$ ) neurons in the lower layer (this field is *receptive field*).

The neuron is "standard": Computes a weighted sum of its inputs and applies an activation function.

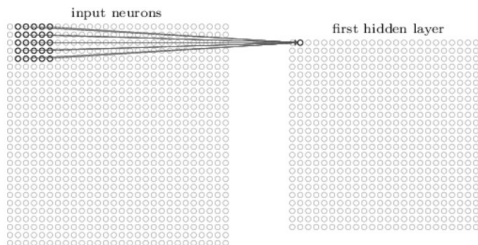


# Convolutional networks - stride length

Then we slide the local receptive field over by one pixel to the right (i.e., by one neuron) to connect to a second hidden neuron:

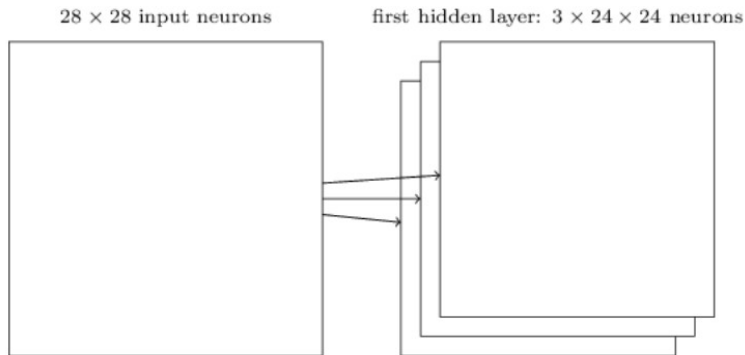


The "size" of the slide is called *stride length*.



The group of all such neurons is *feature map*. all these neurons *share weights and biases!*

# Feature maps

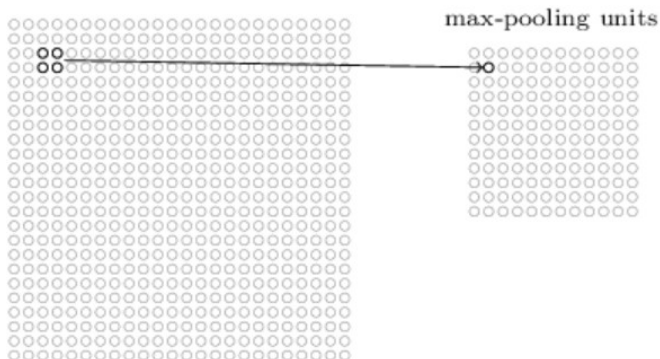


Each feature map represents a property of the input that is supposed to be spatially invariant.

Typically, we consider several feature maps in a single layer.

# Pooling

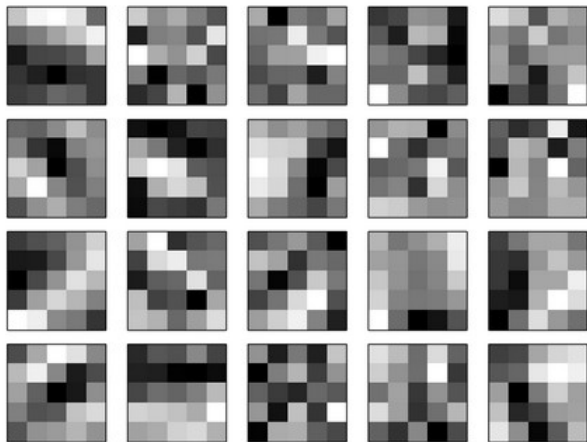
hidden neurons (output from feature map)



Neurons in the pooling layer compute functions of their receptive fields:

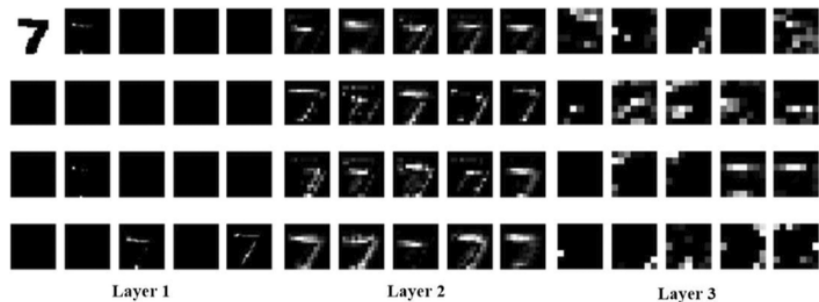
- ▶ **Max-pooling** : maximum of inputs
- ▶ **L2-pooling** : square root of the sum of squares
- ▶ **Average-pooling** : mean
- ▶ ...

## Trained receptive fields

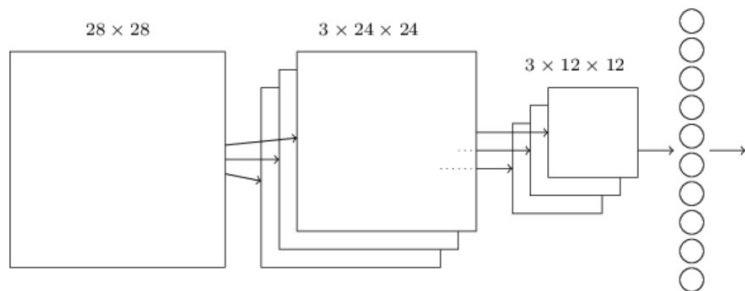


(20 feature maps, receptive fields  $5 \times 5$ )

# Trained feature maps



# Simple convolutional network

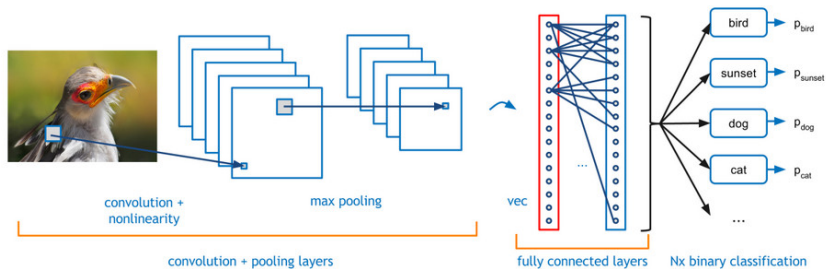


$28 \times 28$  input image, 3 feature maps, each feature map has its own max-pooling (field  $5 \times 5$ , stride = 1), 10 output neurons.

Each neuron in the output layer gets input from each neuron in the pooling layer.

Trained using backprop, which can be easily adapted to convolutional networks.

# Convolutional network



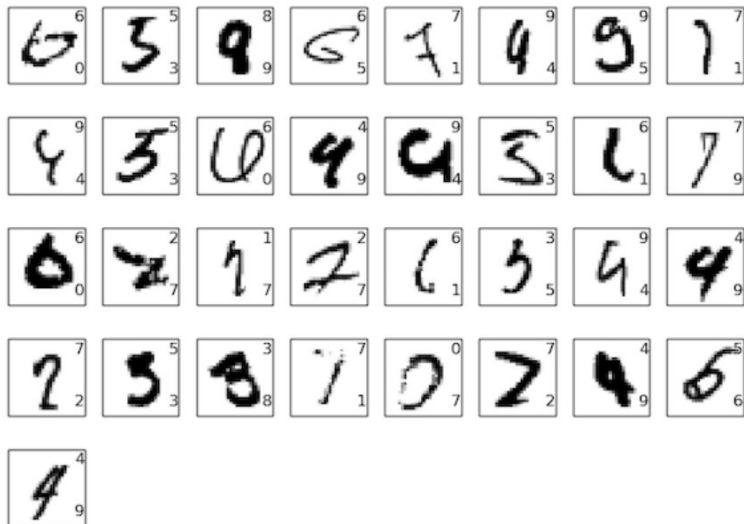
## Simple convolutional network vs MNIST

two convolutional-pooling layers, one 20, second 40 feature maps, two dense (MLP) layers (1000-1000), outputs (10)

- ▶ Activation functions of the feature maps and dense layers: ReLU
- ▶ max-pooling
- ▶ output layer: soft-max
- ▶ Error function: negative log-likelihood (= cross-entropy)
- ▶ Training: SGD, mini-batch size 10
- ▶ learning rate 0.03
- ▶ L2 regularization with "weight"  $\lambda = 0.1$  + dropout with prob. 1/2
- ▶ training for 40 epochs (i.e., every training example is considered 40 times)
- ▶ Expanded dataset: displacement by one pixel to an arbitrary direction.
- ▶ Committee voting of 5 networks.



Out of 10,000 images in the test set, only these 33 have been incorrectly classified:





# ImageNet Large-Scale Visual Recognition Challenge (ILSVRC)

Competition in classification over a subset of images from ImageNet.

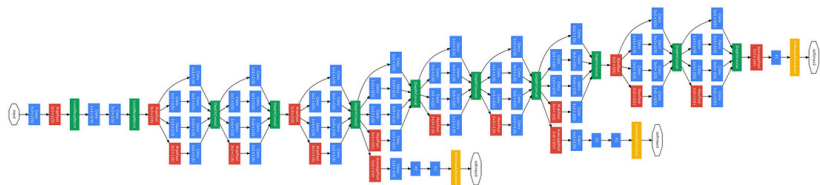
Started in 2010, assisted in a breakthrough in image recognition.

The training set 1.2 million images, 1000 classes. Validation set: 50 000, test set: 150 000.

Many images contain more than one object  $\Rightarrow$  model is allowed to choose five classes, the correct label must be among the five. (top-5 criterion).

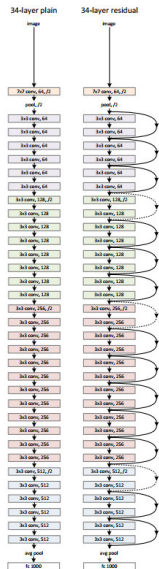
The same set as in 2012, top-5 criterion.

GoogLeNet: deep convolutional network, 22 layers

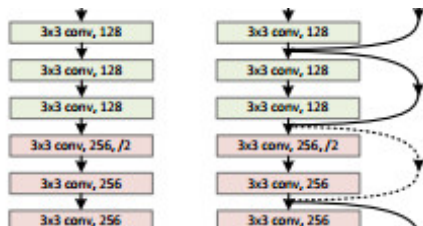


Results:

- ▶ Accuracy 93.33% top-5



- ▶ Deep convolutional network - ResNet
- ▶ Various numbers of layers, the winner has 152 layers
- ▶ Skip connections implementing residual learning
- ▶ Error **3.57%** in top-5.



## Further development of CNN architectures

Convolutional networks made a breakthrough in image recognition.

See IB031 for discussion of ILSVRC challenge.

# Further development of CNN architectures

Convolutional networks made a breakthrough in image recognition.

See IB031 for discussion of ILSVRC challenge.

There are myriads of CNNs for

- ▶ Image classification  
Binary, or many classes, 2D, 3D, ... nD "images" and movies. Various architectures.
- ▶ Image segmentation  
I.e., partitioning pixels of images, a HUGE amount of CNN variants (U-Net, fully convolutional, etc.)
- ▶ Object detection  
Bounding box prediction, R-CNN architectures
- ▶ ...

# Further development of CNN architectures

Convolutional networks made a breakthrough in image recognition.

See IB031 for discussion of ILSVRC challenge.

There are myriads of CNNs for

- ▶ Image classification  
Binary, or many classes, 2D, 3D, ... nD "images" and movies. Various architectures.
- ▶ Image segmentation  
I.e., partitioning pixels of images, a HUGE amount of CNN variants (U-Net, fully convolutional, etc.)
- ▶ Object detection  
Bounding box prediction, R-CNN architectures
- ▶ ...

For more details, see

PA228 Machine Learning in Image Processing.



# Convolutional networks – learning theory

## Convolutional networks – architecture

Neurons organized in layers,  $L_0, L_1, \dots, L_n$ , connections (typically) only from  $L_m$  to  $L_{m+1}$ .

# Convolutional networks – architecture

Neurons organized in layers,  $L_0, L_1, \dots, L_n$ , connections (typically) only from  $L_m$  to  $L_{m+1}$ .

Several types of layers:

- ▶ **input** layer  $L_0$

# Convolutional networks – architecture

Neurons organized in layers,  $L_0, L_1, \dots, L_n$ , connections (typically) only from  $L_m$  to  $L_{m+1}$ .

Several types of layers:

- ▶ **input** layer  $L_0$
- ▶ **dense** layer  $L_m$ : Each neuron of  $L_m$  connected with each neuron of  $L_{m-1}$ .

# Convolutional networks – architecture

Neurons organized in layers,  $L_0, L_1, \dots, L_n$ , connections (typically) only from  $L_m$  to  $L_{m+1}$ .

Several types of layers:

- ▶ **input** layer  $L_0$
- ▶ **dense** layer  $L_m$ : Each neuron of  $L_m$  connected with each neuron of  $L_{m-1}$ .
- ▶ **convolutional** layer  $L_m$ : Neurons organized into disjoint **feature maps**, all neurons of a given feature map *share weights* (but have different inputs)

# Convolutional networks – architecture

Neurons organized in layers,  $L_0, L_1, \dots, L_n$ , connections (typically) only from  $L_m$  to  $L_{m+1}$ .

Several types of layers:

- ▶ **input** layer  $L_0$
- ▶ **dense** layer  $L_m$ : Each neuron of  $L_m$  connected with each neuron of  $L_{m-1}$ .
- ▶ **convolutional** layer  $L_m$ : Neurons organized into disjoint **feature maps**, all neurons of a given feature map *share weights* (but have different inputs)
- ▶ **pooling** layer: "Neurons" organized into **pooling maps**, all neurons
  - ▶ compute a simple aggregate function (such as max),
  - ▶ have *disjoint inputs*.

Pooling after convolution is usually applied to each feature map separately. I.e., a single pooling map after each feature map.

# Convolutional networks – architecture

- ▶ 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  $i$ )
- ▶  $[ji]$  is a set of all connections (i.e. pairs of neurons) sharing the weight  $w_{ji}$ .

# Convolutional networks – activity

- ▶ neurons of dense and convolutional layers:

- ▶ 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):

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



# Convolutional networks – activity

- ▶ neurons of dense and convolutional layers:

- ▶ 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):

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

- ▶ Neurons of pooling layers: Apply the "pooling" function:

- ▶ max-pooling:

$$y_j = \max_{i \in j_{\leftarrow}} y_i$$

- ▶ avg-pooling:

$$y_j = \frac{\sum_{i \in j_{\leftarrow}} y_i}{|j_{\leftarrow}|}$$

A convolutional network is evaluated layer-wise (as MLP), for each  $j \in Y$  we have that  $y_j(\vec{w}, \vec{x})$  is the value of the output neuron  $j$  after evaluating the network with weights  $\vec{w}$  and input  $\vec{x}$ .

# Convolutional networks – learning

## Learning:

- ▶ Given a **training set**  $\mathcal{T}$  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 – squared error (for example):**

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

where

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

# Convolutional networks – SGD

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:
  - ▶ 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

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

Here  $T$  is a *minibatch* (of a fixed size),

- ▶  $0 < \varepsilon(t) \leq 1$  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.

**Epoch** consists of one round through all data.

# Backprop

Recall that  $\nabla E_k(\vec{w}^{(t)})$  is a vector of all partial derivatives of the form  $\frac{\partial E_k}{\partial w_{ji}}$ .

How to compute  $\frac{\partial E_k}{\partial w_{ji}}$  ?

# Backprop

Recall that  $\nabla E_k(\vec{w}^{(t)})$  is a vector of all partial derivatives of the form  $\frac{\partial E_k}{\partial w_{ji}}$ .

How to compute  $\frac{\partial E_k}{\partial w_{ji}}$  ?

First, switch from derivatives w.r.t.  $w_{ji}$  to derivatives w.r.t.  $y_j$ :

- ▶ Recall that for every  $w_{ji}$  where  $j$  is in a dense layer, i.e. does not share weights:

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

# Backprop

Recall that  $\nabla E_k(\vec{w}^{(t)})$  is a vector of all partial derivatives of the form  $\frac{\partial E_k}{\partial w_{ji}}$ .

How to compute  $\frac{\partial E_k}{\partial w_{ji}}$  ?

First, switch from derivatives w.r.t.  $w_{ji}$  to derivatives w.r.t.  $y_j$ :

- ▶ Recall that for every  $w_{ji}$  where  $j$  is in a dense layer, i.e. does not share weights:

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

- ▶ Now for every  $w_{ji}$  where  $j$  is in a convolutional layer:

$$\frac{\partial E_k}{\partial w_{ji}} = \sum_{r \in [ji]} \frac{\partial E_k}{\partial y_r} \cdot \sigma'_r(\xi_r) \cdot y_\ell$$

- ▶ Neurons of pooling layers do not have weights.

# Backprop

Now compute derivatives w.r.t.  $y_j$ :

- ▶ for every  $j \in Y$ :

$$\frac{\partial E_k}{\partial y_j} = y_j - d_{kj}$$

This holds for the squared error, for other error functions the derivative w.r.t. outputs will be different.

# Backprop

Now compute derivatives w.r.t.  $y_j$ :

- ▶ for every  $j \in Y$ :

$$\frac{\partial E_k}{\partial y_j} = y_j - d_{kj}$$

This holds for the squared error, for other error functions the derivative w.r.t. outputs will be different.

- ▶ for every  $j \in Z \setminus Y$  such that  $j^\rightarrow$  is either a dense layer, or a convolutional layer:

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



# Backprop

Now compute derivatives w.r.t.  $y_j$ :

- ▶ for every  $j \in Y$ :

$$\frac{\partial E_k}{\partial y_j} = y_j - d_{kj}$$

This holds for the squared error, for other error functions the derivative w.r.t. outputs will be different.

- ▶ for every  $j \in Z \setminus Y$  such that  $j \rightarrow$  is either a dense layer, or a convolutional layer:

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

- ▶ for every  $j \in Z \setminus Y$  such that  $j \rightarrow$  is max-pooling: Then  $j \rightarrow = \{i\}$  for a single "max" neuron and we have

$$\frac{\partial E_k}{\partial y_j} = \begin{cases} \frac{\partial E_k}{\partial y_i} & \text{if } j = \arg \max_{r \in i \leftarrow} y_r \\ 0 & \text{otherwise} \end{cases}$$

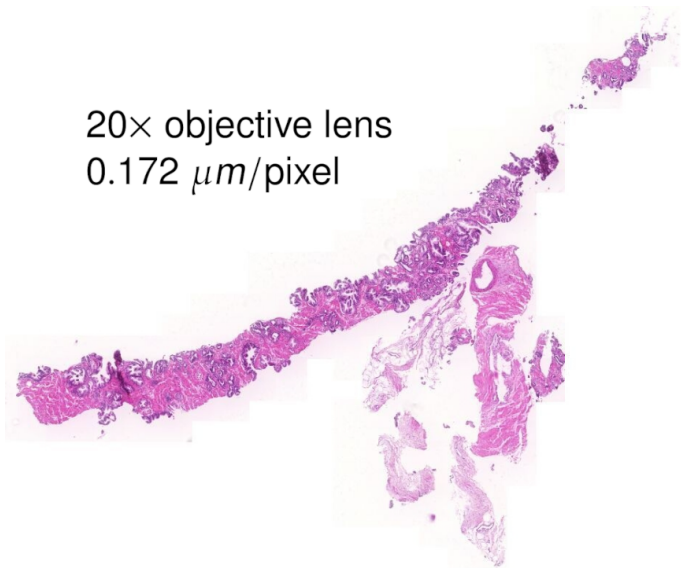
The gradient can be propagated from the output layer downwards as in MLP.

# Convolutional networks – summary

- ▶ Conv. nets. are nowadays the most used networks in image processing (and also in other areas where input has some local, "spatially" invariant properties)  
*2024 update: Vision transformers are gradually taking over.*
- ▶ Typically trained using the gradient descent and its modifications (such as Adam).
- ▶ Due to the weight sharing allow (very) deep architectures.
- ▶ Typically extended with more adjustments and tricks in their topologies.

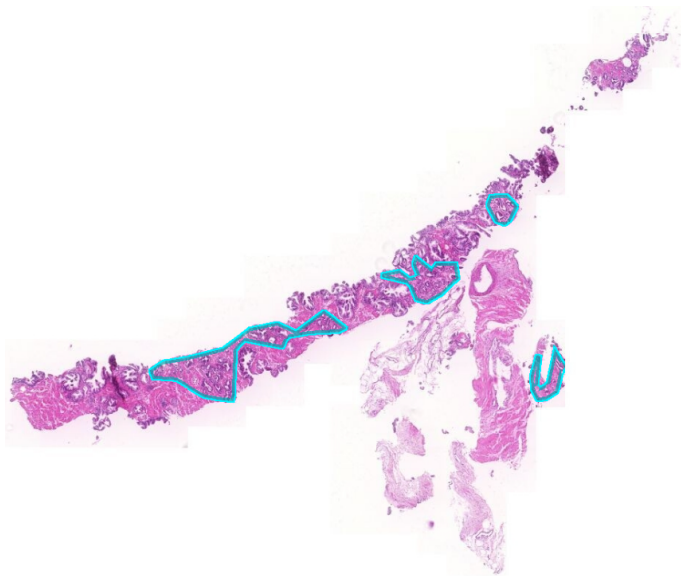
# The problem of cancer detection in WSI

20× objective lens  
0.172  $\mu\text{m}/\text{pixel}$



**The problem:** Detect cancer in this image.

# The problem of cancer detection in WSI

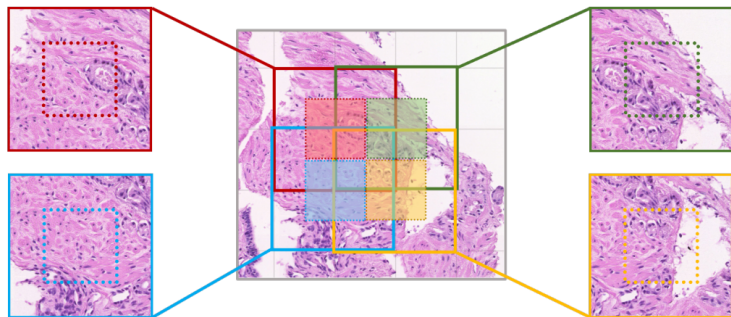


- ▶ WSI annotated by pathologists, **not** pixel level precise!

# Input data

WSI too large, 105,185 px  $\times$  221,772 px

Cut into patches of size 512 px  $\times$  512 px



Patch positive **iff** the inner square intersects the annotation

# Training on WSI

Our dataset from Masaryk Memorial Cancer Institute:

- ▶ 785 WSI from 166 patients  
(698 WSI for training, 87 WSI for testing)
- ▶ Cut into 7,878,675 patches for training, 193,235 patches for testing.

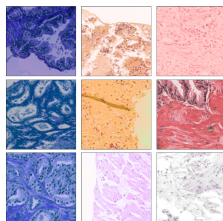
# Training on WSI

Our dataset from Masaryk Memorial Cancer Institute:

- ▶ 785 WSI from 166 patients  
(698 WSI for training, 87 WSI for testing)
- ▶ Cut into 7,878,675 patches for training, 193,235 patches for testing.

Dataset augmentation:

- ▶ random vertical and horizontal flips
- ▶ random color perturbations



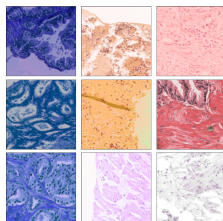
# Training on WSI

Our dataset from Masaryk Memorial Cancer Institute:

- ▶ 785 WSI from 166 patients (698 WSI for training, 87 WSI for testing)
- ▶ Cut into 7,878,675 patches for training, 193,235 patches for testing.

Dataset augmentation:

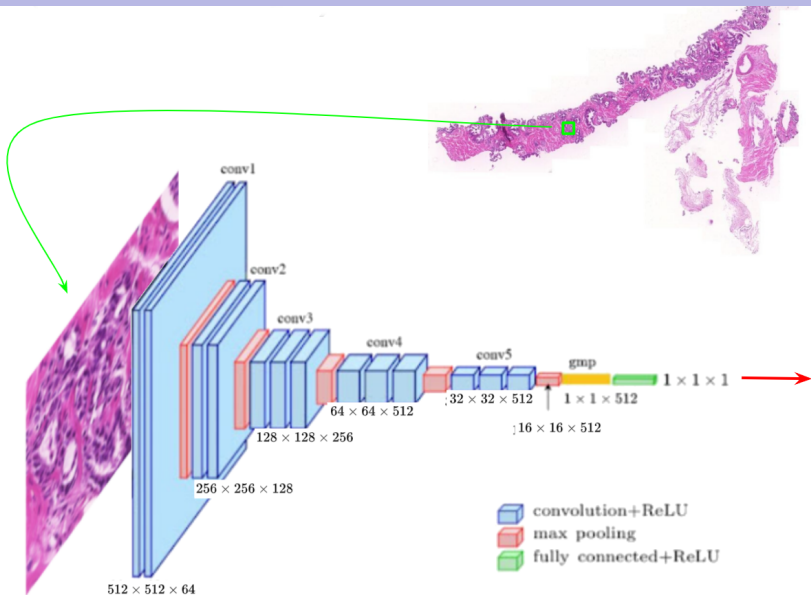
- ▶ random vertical and horizontal flips
- ▶ random color perturbations



- ▶ Training data three step sampling:
  1. randomly select a label
  2. randomly select a slide containing at least a single patch with the label
  3. randomly select a patch with the label from the slide



# VGG16



$3 \times 3$  convolutions, stride 1, padding 1. Max pooling  $2 \times 2$ , stride 2.

- ▶ VGG16 pretrained on the ImageNet (of-the-shelf solution). Top fully connected parts removed, substituted with global max-pooling and a single dense layer.

## Training VGG16 on WSI

- ▶ VGG16 pretrained on the ImageNet (of-the-shelf solution). Top fully connected parts removed, substituted with global max-pooling and a single dense layer.
- ▶ The network has single logistic output - the probability of cancer in the patch

## Training VGG16 on WSI

- ▶ VGG16 pretrained on the ImageNet (of-the-shelf solution). Top fully connected parts removed, substituted with global max-pooling and a single dense layer.
- ▶ The network has single logistic output - the probability of cancer in the patch
- ▶ The error  $E = \text{cross-entropy}$

# Training VGG16 on WSI

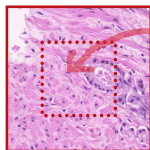
- ▶ VGG16 pretrained on the ImageNet (of-the-shelf solution). Top fully connected parts removed, substituted with global max-pooling and a single dense layer.
- ▶ The network has single logistic output - the probability of cancer in the patch
- ▶ The error  $E$  = cross-entropy
- ▶ Training:
  - ▶ RMSprop optimizer
  - ▶ The "forgetting" hyperparameter:  $\rho = 0.9$
  - ▶ The initial learning rate  $5 \times 10^{-5}$ 
    - ▶ If no improvement in  $E$  on validation data for 3 consecutive epochs  $\Rightarrow$  half the learning rate
    - ▶ If no improvement in ROCAUC on validation data for 5 consecutive epochs  $\Rightarrow$  terminate
  - ▶ Momentum with the weight  $\alpha = 0.9$

# Prediction



## Model evaluation - attempt 2

Can we detect cancer in patches?



Any cancer here?

Predict  $I$  positive iff  $F(I) \geq 0.75$

Single WSI:

		PREDICTED	
		Pos	Neg
TRUE	Pos	805	18
	Neg	48	614

All WSIs:

		PREDICTED	
		Pos	Neg
TRUE	Pos	24796	4340
	Neg	5345	158754

Ok, does it detect cancer?

# Model evaluation – attempt 3 – FROC

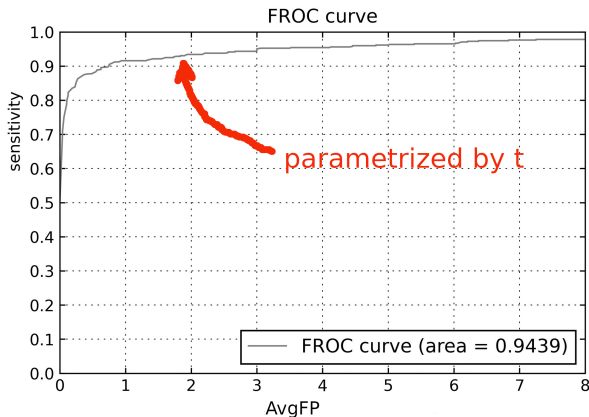
Detect *particular tumors* ?



How to evaluate the quality of tumor detection?



# Model evaluation – attempt 3 – FROC



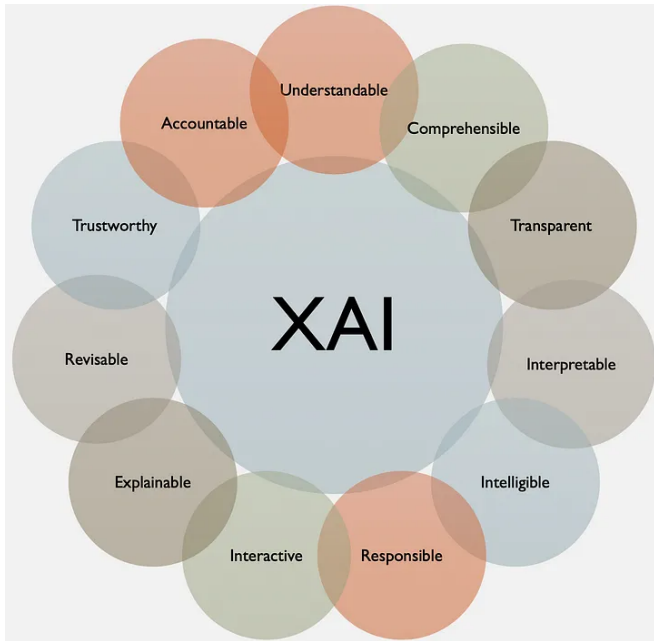
sensitivity  $\approx$  the proportion of tumors containing at least one patch  $I$  with  $F(I) \geq t$  w.r.t. all tumors in all slides

AvgFP  $\approx$  average number of patches  $I$  with  $F(I) \geq t$  in each non-cancerous slide

## Is it good?

- ▶ Tile-based metric does not tell us whether cancer will be detected
- ▶ FROC curve captures tumor detection - however, WSI contains only cuts across tumors
- ▶ WSI level evaluation needs huge amounts of WSIs  
Our later results shown that the system has approx. 0.98 AUC for WSI level tumor detection on thousands of WSI from MMCI
- ▶ The system broke down when we used WSIs from a completely different scanner and different hospital (they used slightly different colors).  
We have corrected this by appropriate color normalization. But more data from more hospitals is needed!

# Explainable methods (XAI)



# What is XAI?

IBM: "Explainable artificial intelligence (XAI) is a set of processes and methods that allows human users to comprehend and trust the results and output created by machine learning algorithms."

<https://www.ibm.com/topics/explainable-ai>

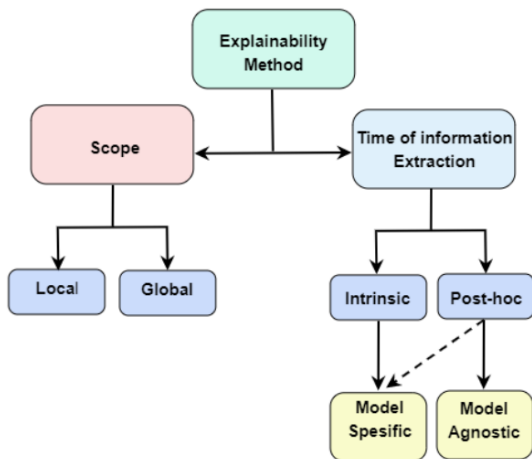
ISO/IEC TR 29119-11:2020(en):

- ▶ Interpretability: Level of understanding how the underlying (AI) technology works.
- ▶ Explainability: Level of understanding how the AI-based system came up with a given result.

# XAI methods

The goal is to understand how and why the network does what it does.

We will consider classification models only.



# A Comprehensive Taxonomy for Explainable Artificial Intelligence: A Systematic Survey of Surveys on Methods and Concepts

Gesina Schwalbe<sup>1,2\*†</sup> and Bettina Finzel<sup>2†</sup>

<sup>1</sup>\*Continental AG, Regensburg, Germany.

<sup>2</sup>Cognitive Systems Group, University of Bamberg, Bamberg, Germany.

\*Corresponding author(s). E-mail(s):

[gesina.schwalbe@continental-corporation.com](mailto:gesina.schwalbe@continental-corporation.com);

Contributing authors: [bettina.finzel@uni-bamberg.de](mailto:bettina.finzel@uni-bamberg.de);

<sup>†</sup>These authors contributed equally to this work.

## Abstract

In the meantime, a wide variety of terminologies, motivations, approaches, and evaluation criteria have been developed within the research field of explainable artificial intelligence (XAI). With the amount of XAI methods vastly growing, a taxonomy of methods is needed by researchers as well as practitioners: To grasp the breadth of the topic, compare methods, and to select the right XAI method based on traits required by a

Many methods are available.

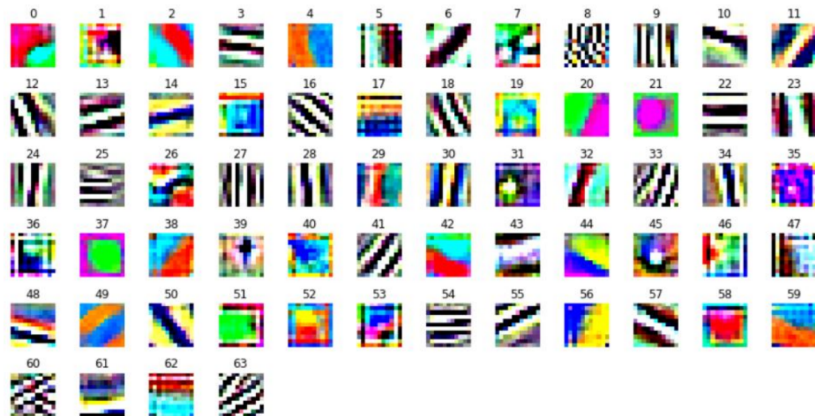
We consider only a few representative methods for interpreting learning models.

Methods based on various principles:

- ▶ Visualize weights and feature maps
- ▶ Visualize the most important inputs for a given class
- ▶ Visualize the effect of input perturbations on the output
- ▶ Construct an interpretable surrogate model

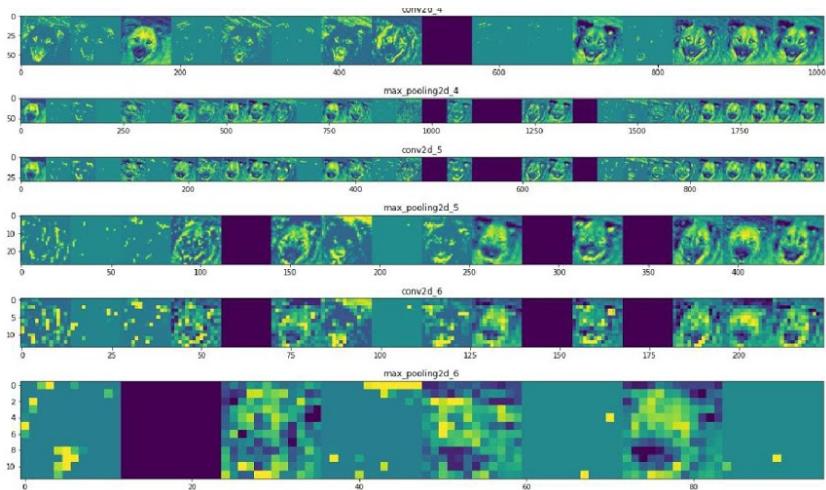


# Alex-net - filters of the first convolutional layer

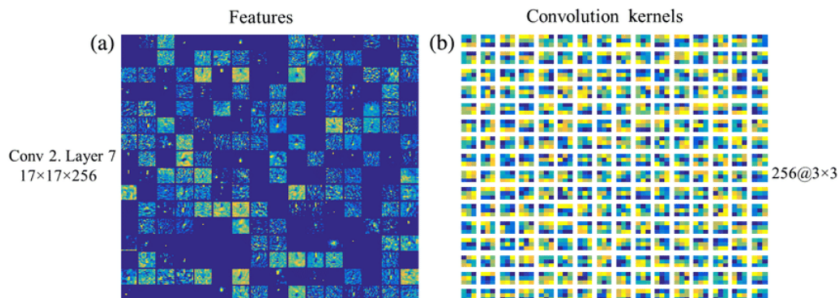


- ▶ 64 filters of depth 3 (RGB)
- ▶ Combined each filter RGB channels into one RGB image of size 11x11x3.

# CNN - feature maps



# CNN - feature maps - radar target classification



Synthetic-aperture radar (SAR) – used to create two-dimensional images or three-dimensional reconstructions of objects, such as landscapes.

## Maximizing input

Now what if we try to find the most "representative" input vector for a given class?

## Maximizing input

Now what if we try to find the most "representative" input vector for a given class?

Assume a trained model giving a score for each class given an input vector.

## Maximizing input

Now what if we try to find the most "representative" input vector for a given class?

Assume a trained model giving a score for each class given an input vector.

- ▶ Denote by  $\xi_i(\vec{x})$  the inner potential of the *output* neuron  $i \in Y$  given a network input vector  $\vec{x}$ .

# Maximizing input

Now what if we try to find the most "representative" input vector for a given class?

Assume a trained model giving a score for each class given an input vector.

- ▶ Denote by  $\xi_i(\vec{x})$  the inner potential of the *output* neuron  $i \in Y$  given a network input vector  $\vec{x}$ .
- ▶ Maximize

$$\xi_i(\vec{x}) - \lambda \|\vec{x}\|_2^2$$

over all input vectors  $\vec{x}$ .

# Maximizing input

Now what if we try to find the most "representative" input vector for a given class?

Assume a trained model giving a score for each class given an input vector.

- ▶ Denote by  $\xi_i(\vec{x})$  the inner potential of the *output* neuron  $i \in Y$  given a network input vector  $\vec{x}$ .
- ▶ Maximize

$$\xi_i(\vec{x}) - \lambda \|\vec{x}\|_2^2$$

over all input vectors  $\vec{x}$ .

- ▶ A maximizing input vector computed using the gradient ascent.
- ▶ Gives the most "representative" input vector of the class represented by the neuron  $i$ .



## Maximizing input - example



**dumbbell**



**cup**



**dalmatian**

## Input specific saliency maps

**The goal:** Label features in a given input that are "most important" for the output of the network.

**The goal:** Label features in a given input that are "most important" for the output of the network.

Various approaches:

- ▶ gradient based
  - ▶ Gradient saliency maps
  - ▶ GradCAM
  - ▶ SmoothGrad
  - ▶ ...
- ▶ occlusion based
  - ▶ Simple occlusion maps
  - ▶ LIME
  - ▶ ...

## Gradient based saliency

- ▶ Let us fix an output neuron  $i$  and an input vector  $\vec{x}$ .

## Gradient based saliency

- ▶ Let us fix an output neuron  $i$  and an input vector  $\vec{x}$ .
- ▶ **Idea:** Rank every input neuron  $k \in X$  based on its influence on the value  $\xi_i(\vec{x})$ .

Note that the vector of input values is *fixed*.

## Gradient based saliency

- ▶ Let us fix an output neuron  $i$  and an input vector  $\vec{x}$ .
- ▶ **Idea:** Rank every input neuron  $k \in X$  based on its influence on the value  $\xi_i(\vec{x})$ .

Note that the vector of input values is *fixed*.

For every input neuron  $k \in X$  we consider

$$\left| \frac{\partial \xi_i}{\partial y_k}(\vec{x}) \right|$$

to measure the importance of the input  $y_k$  for the output potential  $\xi_i$  with respect to the particular input vector  $\vec{x}$ .

# Gradient based saliency

- ▶ Let us fix an output neuron  $i$  and an input vector  $\vec{x}$ .
- ▶ **Idea:** Rank every input neuron  $k \in X$  based on its influence on the value  $\xi_i(\vec{x})$ .

Note that the vector of input values is *fixed*.

For every input neuron  $k \in X$  we consider

$$\left| \frac{\partial \xi_i}{\partial y_k}(\vec{x}) \right|$$

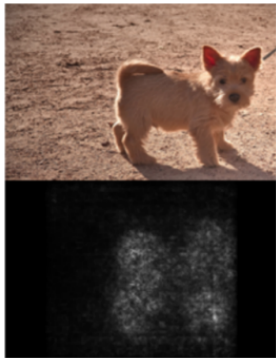
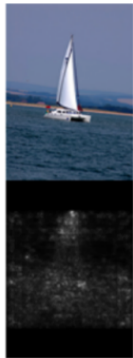
to measure the importance of the input  $y_k$  for the output potential  $\xi_i$  with respect to the particular input vector  $\vec{x}$ .

- ▶ Note that saliency comes from a surrogate local linear model given by the first-order Taylor approximation:

$$\xi_i(\vec{x}') \approx \xi_i(\vec{x}) + \left( \frac{\partial \xi_i}{\partial \mathbf{X}}(\vec{x}) \right) (\vec{x}' - \vec{x})$$

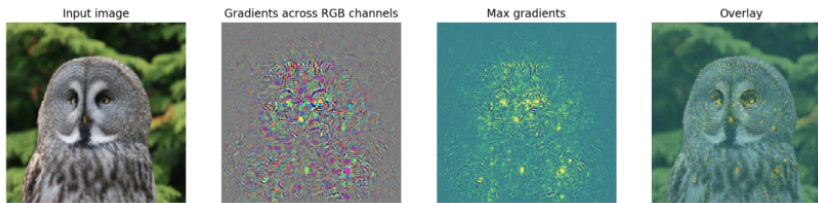
Here  $\frac{\partial \xi_i}{\partial \mathbf{X}}$  is the vector of all partial derivatives  $\frac{\partial \xi_i}{\partial y_k}$  where  $k \in X$ .

## Saliency maps - example



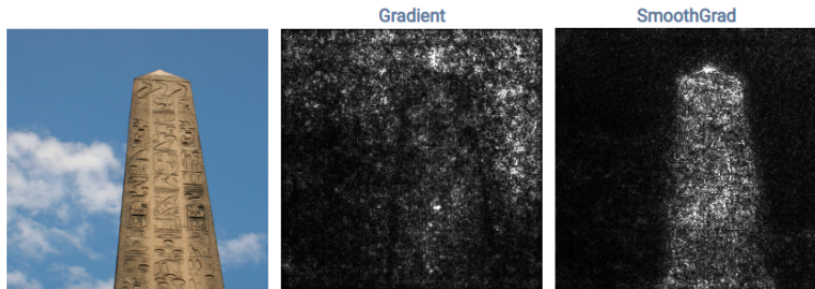


# Saliency maps - example



Quite noisy, the signal is spread and does not say much about the perception of the owl.

# Saliency maps - example



## SmoothGrad:

- ▶ Do the following several times:
  - ▶ Add noise to the input image
  - ▶ Compute a saliency map
- ▶ Average the resulting saliency maps.

- ▶ Consider a convolutional network and fix an input image  $I$  of the network.

ALL values of all neurons  $y_j$  are computed on the input  $I$ .

- ▶ Consider a convolutional network and fix an input image  $I$  of the network.

ALL values of all neurons  $y_j$  are computed on the input  $I$ .

- ▶ Fix a convolutional layer  $L$  consisting of convolutional feature maps  $F^1, \dots, F^k$ .

Each  $F^\ell$  is a set of neurons that belong to the feature map  $F^\ell$ .

Slightly abusing notation, we write  $F^\ell(I)$  to denote the tensor of all values of all neurons in  $F^\ell(I)$ .

- ▶ Consider a convolutional network and fix an input image  $I$  of the network.

ALL values of all neurons  $y_j$  are computed on the input  $I$ .

- ▶ Fix a convolutional layer  $L$  consisting of convolutional feature maps  $F^1, \dots, F^k$ .

Each  $F^\ell$  is a set of neurons that belong to the feature map  $F^\ell$ .

Slightly abusing notation, we write  $F^\ell(I)$  to denote the tensor of all values of all neurons in  $F^\ell(I)$ .

- ▶ Fix an output neuron  $i \in Y$  with the inner potential  $\xi_i$ . Compute the *average importance* of  $F^\ell(I)$ :

$$\alpha_i^\ell = \frac{1}{|F^\ell|} \sum_{j \in F^\ell} \frac{\partial \xi_i}{\partial y_j} (F^\ell(I))$$

- ▶ Consider a convolutional network and fix an input image  $I$  of the network.

ALL values of all neurons  $y_j$  are computed on the input  $I$ .

- ▶ Fix a convolutional layer  $L$  consisting of convolutional feature maps  $F^1, \dots, F^k$ .

Each  $F^\ell$  is a set of neurons that belong to the feature map  $F^\ell$ .

Slightly abusing notation, we write  $F^\ell(I)$  to denote the tensor of all values of all neurons in  $F^\ell(I)$ .

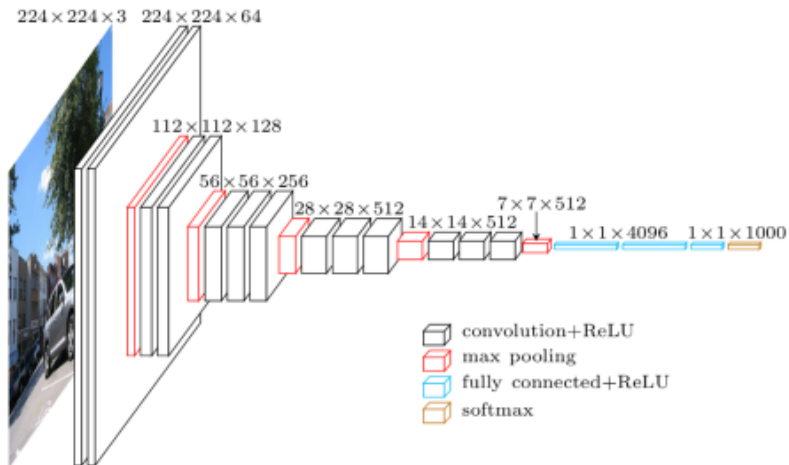
- ▶ Fix an output neuron  $i \in Y$  with the inner potential  $\xi_i$ . Compute the *average importance* of  $F^\ell(I)$ :

$$\alpha_i^\ell = \frac{1}{|F^\ell|} \sum_{j \in F^\ell} \frac{\partial \xi_i}{\partial y_j} (F^\ell(I))$$

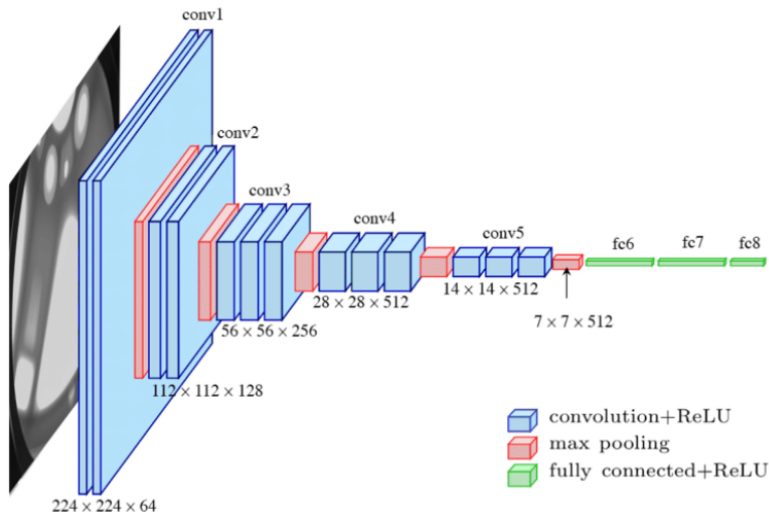
and the final *gradCAM heat map* for  $L$  is obtained using

$$M_i^L = \text{ReLU} \left( \sum_{\ell=1}^k \alpha_i^\ell F^\ell(I) \right)$$

# GradCAM on VGG16



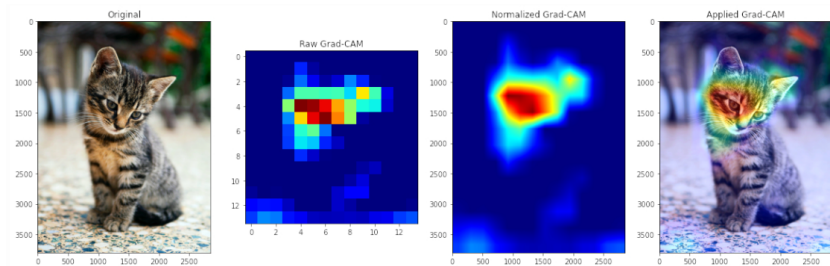
# GradCAM on VGG16



Consider the last convolutional layer of the VGG16 (Block5,



# GradCAM on VGG16

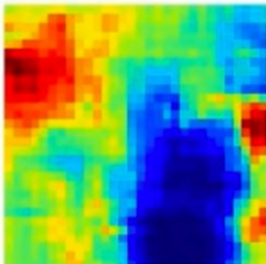
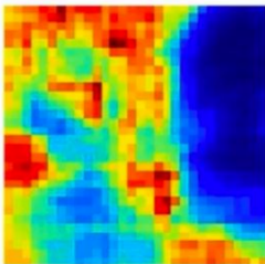
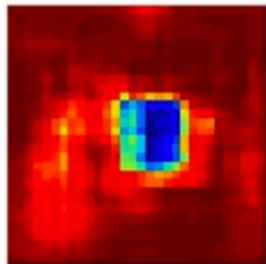
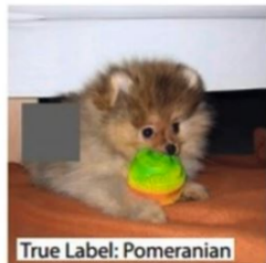


From left to right:

- ▶ An image of a cat (has to be resized to  $224 \times 224$  to fit VGG16)
- ▶ The gradCAM heat map for the last convolutional layer and the class "cat"
- ▶ Rescaled and smoothed gradCAM heat map.
- ▶ The gradCAM overlay.

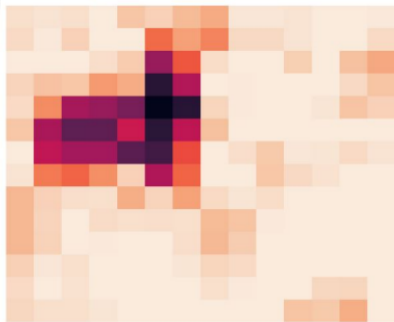
- ▶ Systematically cover parts of the input image.
- ▶ Observe the effect on the output value.
- ▶ Find regions with the largest effect.

# Occlusion - example



# Occlusion - example

['harmonica, mouth organ, harp, mouth harp']



## LIME - for images

Let us fix an image  $I$  to be explained.

# LIME - for images

Let us fix an image  $I$  to be explained.

Outline:

- ▶ Consider superpixels of  $I$  as interpretable components.

# LIME - for images

Let us fix an image  $I$  to be explained.

Outline:

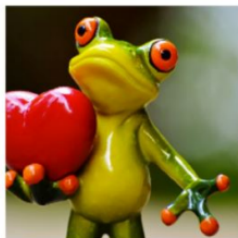
- ▶ Consider superpixels of  $I$  as interpretable components.
- ▶ Construct a linear model approximating the network around the image  $I$  with weights corresponding to the superpixels.

# LIME - for images

Let us fix an image  $I$  to be explained.

Outline:

- ▶ Consider superpixels of  $I$  as interpretable components.
- ▶ Construct a linear model approximating the network around the image  $I$  with weights corresponding to the superpixels.
- ▶ Select the superpixels with weights of large magnitude as the important ones.



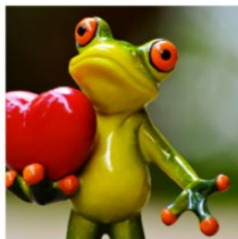
Original Image



Interpretable  
Components



# Superpixels as interpretable components



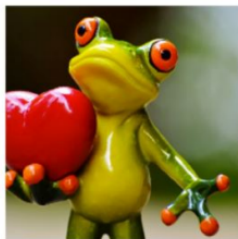
Original Image



Interpretable  
Components

Denote by  $P_1, \dots, P_\ell$  all superpixels of  $I$ .

# Superpixels as interpretable components



Original Image

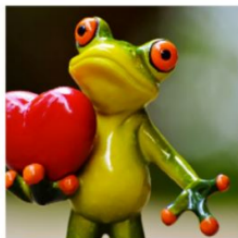


Interpretable  
Components

Denote by  $P_1, \dots, P_\ell$  all superpixels of  $I$ .

Consider binary vectors  $\vec{x} = (x_1, \dots, x_\ell) \in \{0, 1\}^\ell$ .

# Superpixels as interpretable components



Original Image



Interpretable  
Components

Denote by  $P_1, \dots, P_\ell$  all superpixels of  $I$ .

Consider binary vectors  $\vec{x} = (x_1, \dots, x_\ell) \in \{0, 1\}^\ell$ .

Each such vector  $\vec{x}$  determines a "subimage"  $I[\vec{x}]$  of  $I$  obtained by removing all  $P_k$  with  $x_k = 0$ .



- ▶ Let us fix an output neuron  $i$ , we denote by  $\xi_i(J)$  the inner potential of the output neuron  $i$  for the input image  $J$ .

- ▶ Let us fix an output neuron  $i$ , we denote by  $\xi_i(J)$  the inner potential of the output neuron  $i$  for the input image  $J$ .
- ▶ Given the image  $I$  to be interpreted, consider the following training set:

$$\mathcal{T} = \{(\vec{x}_1, \xi_i(I[\vec{x}_1])), \dots, (\vec{x}_p, \xi_i(I[\vec{x}_p]))\}$$

Here  $\vec{x}_h = (x_{h1}, \dots, x_{h\ell})$  are (some) binary vectors of  $\{0, 1\}$ .  
E.g., randomly selected.

- ▶ Let us fix an output neuron  $i$ , we denote by  $\xi_i(J)$  the inner potential of the output neuron  $i$  for the input image  $J$ .
- ▶ Given the image  $I$  to be interpreted, consider the following training set:

$$\mathcal{T} = \{(\vec{x}_1, \xi_i(I[\vec{x}_1])), \dots, (\vec{x}_p, \xi_i(I[\vec{x}_p]))\}$$

Here  $\vec{x}_h = (x_{h1}, \dots, x_{h\ell})$  are (some) binary vectors of  $\{0, 1\}$ .  
E.g., randomly selected.

- ▶ Train a linear model (ADALINE) with weights  $w_0, w_1, \dots, w_\ell$  on  $\mathcal{T}$ .  
Intuitively, the linear model approximates the network on "subimages" of  $I$  obtained by removing some superpixels.
- ▶ Inspect the weights (magnitude and sign).

More precisely, we train a linear model (ADALINE)  $F$  with weights  $\vec{w} = w_0, w_1, \dots, w_\ell$  on  $\mathcal{T}$  minimizing the weighted mean-squared error

$$\mathcal{E}(\vec{w}) = \frac{1}{p} \sum_{k=1}^p \pi_k \cdot (F(\vec{x}_k) - \xi_i(I[\vec{x}_k]))^2 + \Omega(\vec{w})$$

where

- ▶ the weights are defined by

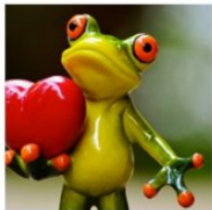
$$\pi_k = \exp\left(\frac{-(1 - \sqrt{1 - (s_k/\ell)})^2}{2\nu^2}\right)$$

Here  $s_k$  is the number of elements in  $\vec{x}_k$  equal to zero,  $\ell$  is the number of superpixels,  $\nu$  determines how much perturbed images are taken into account in the error.

Small  $\nu$  means that  $\pi_k$  is close to zero for  $\vec{x}_k$  with many zeros.







- ▶  $\Omega(\vec{w})$  is a regularization term making the number of non-zero weights as small as possible.

# LIME - example



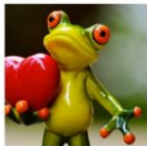
Original Image  
 $P(\text{tree frog}) = 0.54$



Perturbed Instances	$P(\text{tree frog})$
	 0.85
	 0.00001
	 0.52

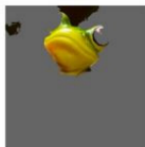


# LIME - example



Original Image  
 $P(\text{tree frog}) = 0.54$

Perturbed Instances	$P(\text{tree frog})$
A perturbed version of the frog image where several pixels are missing, leaving gaps.	A horizontal bar chart showing a value of 0.85. 0.85
A perturbed version of the frog image where almost all pixels are missing, leaving only a few scattered fragments.	A horizontal bar chart showing a value of 0.00001. 0.00001
The original frog image.	A horizontal bar chart showing a value of 0.52. 0.52



Explanation

# LIME - example



(a) Original Image



(b) Explaining *Electric guitar*

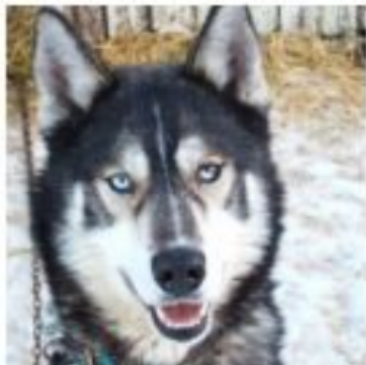


(c) Explaining *Acoustic guitar*



(d) Explaining *Labrador*

## LIME - example



(a) Husky classified as wolf



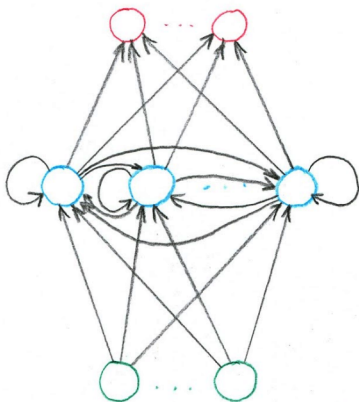
(b) Explanation

# Recurrent Neural Networks

OUTPUT

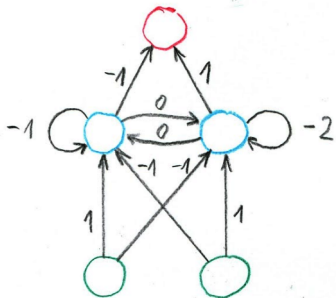
HIDDEN

INPUT



- ▶ Input:  
 $\vec{x} = (x_1, \dots, x_M)$
- ▶ Hidden:  
 $\vec{h} = (h_1, \dots, h_H)$
- ▶ Output:  
 $\vec{y} = (y_1, \dots, y_N)$

# RNN example

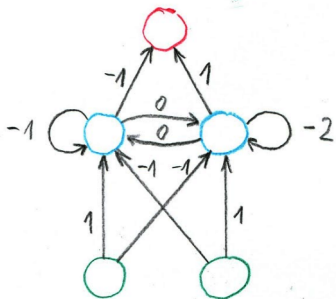


Activation function:

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

<b>y</b>		1	0	1	
<b>h</b>	(0,0)	(1,1)	(1,0)	(0,1)	...
<b>x</b>		(0,0)	(1,0)	(1,1)	

# RNN example

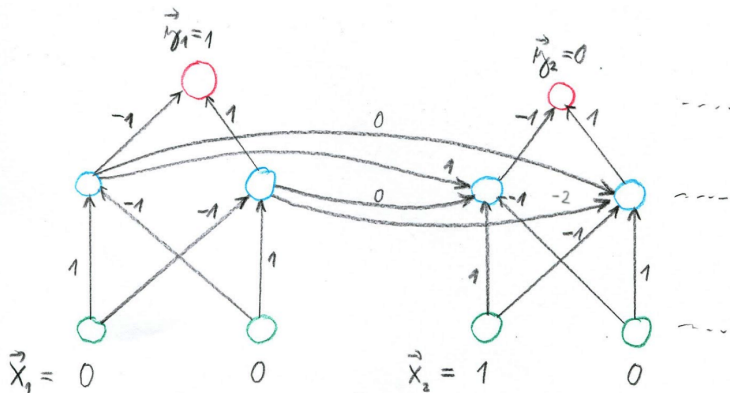


Activation function:

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

<b>y</b>		$\vec{y}_1 = 1$	$\vec{y}_2 = 0$	$\vec{y}_3 = 1$	
<b>h</b>	$\vec{h}_0 = (0, 0)$	$\vec{h}_1 = (1, 1)$	$\vec{h}_2 = (1, 0)$	$\vec{h}_3 = (0, 1)$	...
<b>x</b>		$\vec{x}_1 = (0, 0)$	$\vec{x}_2 = (1, 0)$	$\vec{x}_3 = (1, 1)$	

# RNN example

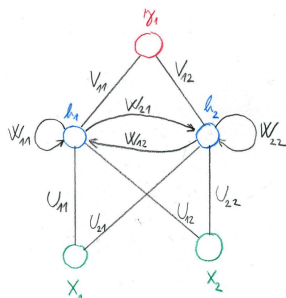


<b>y</b>		$\vec{y}_1 = 1$	$\vec{y}_2 = 0$	$\vec{y}_3 = 1$	
<b>h</b>	$\vec{h}_0 = (0, 0)$	$\vec{h}_1 = (1, 1)$	$\vec{h}_2 = (1, 0)$	$\vec{h}_3 = (0, 1)$	...
<b>x</b>		$\vec{x}_1 = (0, 0)$	$\vec{x}_2 = (1, 0)$	$\vec{x}_3 = (1, 1)$	



# RNN – formally

- ▶  $M$  inputs:  $\vec{x} = (x_1, \dots, x_M)$
- ▶  $H$  hidden neurons:  $\vec{h} = (h_1, \dots, h_H)$
- ▶  $N$  output neurons:  $\vec{y} = (y_1, \dots, y_N)$
- ▶ Weights:
  - ▶  $U_{kk'}$  from input  $x_{k'}$  to hidden  $h_k$
  - ▶  $W_{kk'}$  from hidden  $h_{k'}$  to hidden  $h_k$
  - ▶  $V_{kk'}$  from hidden  $h_{k'}$  to output  $y_k$



## RNN – formally

- ▶ Input sequence:  $\mathbf{x} = \vec{x}_1, \dots, \vec{x}_T$

$$\vec{x}_t = (x_{t1}, \dots, x_{tM})$$

# RNN – formally

- ▶ Input sequence:  $\mathbf{x} = \vec{x}_1, \dots, \vec{x}_T$

$$\vec{x}_t = (x_{t1}, \dots, x_{tM})$$

- ▶ Hidden sequence:  $\mathbf{h} = \vec{h}_0, \vec{h}_1, \dots, \vec{h}_T$

$$\vec{h}_t = (h_{t1}, \dots, h_{tH})$$

We have  $\vec{h}_0 = (0, \dots, 0)$  and

$$\vec{h}_{tk} = \sigma \left( \sum_{k'=1}^M U_{kk'} x_{tk'} + \sum_{k'=1}^H W_{kk'} h_{(t-1)k'} \right)$$

# RNN – formally

- ▶ Input sequence:  $\mathbf{x} = \vec{x}_1, \dots, \vec{x}_T$

$$\vec{x}_t = (x_{t1}, \dots, x_{tM})$$

- ▶ Hidden sequence:  $\mathbf{h} = \vec{h}_0, \vec{h}_1, \dots, \vec{h}_T$

$$\vec{h}_t = (h_{t1}, \dots, h_{tH})$$

We have  $\vec{h}_0 = (0, \dots, 0)$  and

$$\vec{h}_{tk} = \sigma \left( \sum_{k'=1}^M U_{kk'} x_{tk'} + \sum_{k'=1}^H W_{kk'} h_{(t-1)k'} \right)$$

- ▶ Output sequence:  $\mathbf{y} = \vec{y}_1, \dots, \vec{y}_T$

$$\vec{y}_t = (y_{t1}, \dots, y_{tN})$$

where  $y_{tk} = \sigma \left( \sum_{k'=1}^H V_{kk'} h_{tk'} \right)$ .

## RNN – in matrix form

- ▶ Input sequence:  $\mathbf{x} = \vec{x}_1, \dots, \vec{x}_T$

## RNN – in matrix form

- ▶ Input sequence:  $\mathbf{x} = \vec{x}_1, \dots, \vec{x}_T$
- ▶ Hidden sequence:  $\mathbf{h} = \vec{h}_0, \vec{h}_1, \dots, \vec{h}_T$  where

$$\vec{h}_0 = (0, \dots, 0)$$

and

$$\vec{h}_t = \sigma(U\vec{x}_t + W\vec{h}_{t-1})$$

## RNN – in matrix form

- ▶ Input sequence:  $\mathbf{x} = \vec{x}_1, \dots, \vec{x}_T$
- ▶ Hidden sequence:  $\mathbf{h} = \vec{h}_0, \vec{h}_1, \dots, \vec{h}_T$  where

$$\vec{h}_0 = (0, \dots, 0)$$

and

$$\vec{h}_t = \sigma(U\vec{x}_t + W\vec{h}_{t-1})$$

- ▶ Output sequence:  $\mathbf{y} = \vec{y}_1, \dots, \vec{y}_T$  where

$$y_t = \sigma(Vh_t)$$

- ▶  $\vec{h}_t$  is the memory of the network, captures what happened in all previous steps (with decaying quality).
- ▶ RNN **shares weights**  $U, V, W$  along the sequence.  
Note the similarity to convolutional networks where the weights were shared spatially over images, here they are shared temporally over sequences.
- ▶ RNN can deal with **sequences of variable length**.  
Compare with MLP which accepts only fixed-dimension vectors on input.



**The Task:** Design a recurrent network with a single hidden layer which works as a binary adder.

Example of behavior: Input two binary numbers, e.g., 111 and 101 (we assume that the least significant bit is on the *left*).

The input of the network will be:  $(1, 1), (1, 0), (1, 1)$

The output is supposed to be: 0, 0, 1 (we ignore the carry at the end).

## Training set

$$\mathcal{T} = \{(\mathbf{x}_1, \mathbf{d}_1), \dots, (\mathbf{x}_p, \mathbf{d}_p)\}$$

here

- ▶ each  $\mathbf{x}_\ell = \vec{x}_{\ell 1}, \dots, \vec{x}_{\ell T_\ell}$  is an input sequence,
- ▶ each  $\mathbf{d}_\ell = \vec{d}_{\ell 1}, \dots, \vec{d}_{\ell T_\ell}$  is an expected output sequence.

Here each  $\vec{x}_{\ell t} = (x_{\ell t 1}, \dots, x_{\ell t M})$  is an input vector and each  $\vec{d}_{\ell t} = (d_{\ell t 1}, \dots, d_{\ell t N})$  is an expected output vector.

# Error function

In what follows I will consider a training set with a **single element**  $(\mathbf{x}, \mathbf{d})$ . I.e. drop the index  $\ell$  and have

- ▶  $\mathbf{x} = \vec{x}_1, \dots, \vec{x}_T$  where  $\vec{x}_t = (x_{t1}, \dots, x_{tM})$
- ▶  $\mathbf{d} = \vec{d}_1, \dots, \vec{d}_T$  where  $\vec{d}_t = (d_{t1}, \dots, d_{tN})$

The squared error of  $(\mathbf{x}, \mathbf{d})$  is defined by

$$E_{(\mathbf{x}, \mathbf{d})} = \sum_{t=1}^T \sum_{k=1}^N \frac{1}{2} (y_{tk} - d_{tk})^2$$

Recall that we have a sequence of network outputs  $\mathbf{y} = \vec{y}_1, \dots, \vec{y}_T$  and thus  $y_{tk}$  is the  $k$ -th component of  $\vec{y}_t$

## Gradient descent (single training example)

Consider a single training example  $(\mathbf{x}, \mathbf{d})$ .

The algorithm computes a sequence of weight matrices as follows:

## Gradient descent (single training example)

Consider a single training example  $(\mathbf{x}, \mathbf{d})$ .

The algorithm computes a sequence of weight matrices as follows:

- ▶ Initialize all weights randomly close to 0.

## Gradient descent (single training example)

Consider a single training example  $(\mathbf{x}, \mathbf{d})$ .

The algorithm computes a sequence of weight matrices as follows:

- ▶ Initialize all weights randomly close to 0.
- ▶ In the step  $\ell + 1$  (here  $\ell = 0, 1, 2, \dots$ ) compute "new" weights  $U^{(\ell+1)}, V^{(\ell+1)}, W^{(\ell+1)}$  from the "old" weights  $U^{(\ell)}, V^{(\ell)}, W^{(\ell)}$  as follows:

$$U_{kk'}^{(\ell+1)} = U_{kk'}^{(\ell)} - \varepsilon(\ell) \cdot \frac{\delta E_{(\mathbf{x}, \mathbf{d})}}{\delta U_{kk'}}$$

$$V_{kk'}^{(\ell+1)} = V_{kk'}^{(\ell)} - \varepsilon(\ell) \cdot \frac{\delta E_{(\mathbf{x}, \mathbf{d})}}{\delta V_{kk'}}$$

$$W_{kk'}^{(\ell+1)} = W_{kk'}^{(\ell)} - \varepsilon(\ell) \cdot \frac{\delta E_{(\mathbf{x}, \mathbf{d})}}{\delta W_{kk'}}$$

# Gradient descent (single training example)

Consider a single training example  $(\mathbf{x}, \mathbf{d})$ .

The algorithm computes a sequence of weight matrices as follows:

- ▶ Initialize all weights randomly close to 0.
- ▶ In the step  $\ell + 1$  (here  $\ell = 0, 1, 2, \dots$ ) compute "new" weights  $U^{(\ell+1)}, V^{(\ell+1)}, W^{(\ell+1)}$  from the "old" weights  $U^{(\ell)}, V^{(\ell)}, W^{(\ell)}$  as follows:

$$U_{kk'}^{(\ell+1)} = U_{kk'}^{(\ell)} - \varepsilon(\ell) \cdot \frac{\delta E_{(\mathbf{x}, \mathbf{d})}}{\delta U_{kk'}}$$

$$V_{kk'}^{(\ell+1)} = V_{kk'}^{(\ell)} - \varepsilon(\ell) \cdot \frac{\delta E_{(\mathbf{x}, \mathbf{d})}}{\delta V_{kk'}}$$

$$W_{kk'}^{(\ell+1)} = W_{kk'}^{(\ell)} - \varepsilon(\ell) \cdot \frac{\delta E_{(\mathbf{x}, \mathbf{d})}}{\delta W_{kk'}}$$

**The above is THE learning algorithm that modifies weights!**

# Backpropagation

**Computes the derivatives of  $E$ , no weights are modified!**



# Backpropagation

Computes the derivatives of  $E$ , no weights are modified!

$$\frac{\delta E_{(\mathbf{x}, \mathbf{d})}}{\delta U_{kk'}} = \sum_{t=1}^T \frac{\delta E_{(\mathbf{x}, \mathbf{d})}}{\delta h_{tk}} \cdot \sigma' \cdot x_{tk'} \quad k' = 1, \dots, M$$

$$\frac{\delta E_{(\mathbf{x}, \mathbf{d})}}{\delta V_{kk'}} = \sum_{t=1}^T \frac{\delta E_{(\mathbf{x}, \mathbf{d})}}{\delta y_{tk}} \cdot \sigma' \cdot h_{tk'} \quad k' = 1, \dots, H$$

$$\frac{\delta E_{(\mathbf{x}, \mathbf{d})}}{\delta W_{kk'}} = \sum_{t=1}^T \frac{\delta E_{(\mathbf{x}, \mathbf{d})}}{\delta h_{tk}} \cdot \sigma' \cdot h_{(t-1)k'} \quad k' = 1, \dots, H$$

# Backpropagation

Computes the derivatives of  $E$ , no weights are modified!

$$\frac{\delta E(\mathbf{x}, \mathbf{d})}{\delta U_{kk'}} = \sum_{t=1}^T \frac{\delta E(\mathbf{x}, \mathbf{d})}{\delta h_{tk}} \cdot \sigma' \cdot x_{tk'} \quad k' = 1, \dots, M$$

$$\frac{\delta E(\mathbf{x}, \mathbf{d})}{\delta V_{kk'}} = \sum_{t=1}^T \frac{\delta E(\mathbf{x}, \mathbf{d})}{\delta y_{tk}} \cdot \sigma' \cdot h_{tk'} \quad k' = 1, \dots, H$$

$$\frac{\delta E(\mathbf{x}, \mathbf{d})}{\delta W_{kk'}} = \sum_{t=1}^T \frac{\delta E(\mathbf{x}, \mathbf{d})}{\delta h_{tk}} \cdot \sigma' \cdot h_{(t-1)k'} \quad k' = 1, \dots, H$$

Backpropagation:

$$\frac{\delta E(\mathbf{x}, \mathbf{d})}{\delta y_{tk}} = y_{tk} - d_{tk} \quad (\text{assuming squared error})$$

$$\frac{\delta E(\mathbf{x}, \mathbf{d})}{\delta h_{tk}} = \sum_{k'=1}^N \frac{\delta E(\mathbf{x}, \mathbf{d})}{\delta y_{tk'}} \cdot \sigma' \cdot V_{k'k} + \sum_{k'=1}^H \frac{\delta E(\mathbf{x}, \mathbf{d})}{\delta h_{(t+1)k'}} \cdot \sigma' \cdot W_{k'k}$$

# Long-term dependencies

$$\frac{\delta E_{(x,d)}}{\delta h_{tk}} = \sum_{k'=1}^N \frac{\delta E_{(x,d)}}{\delta y_{tk'}} \cdot \sigma' \cdot V_{k'k} + \sum_{k'=1}^H \frac{\delta E_{(x,d)}}{\delta h_{(t+1)k'}} \cdot \sigma' \cdot W_{k'k}$$

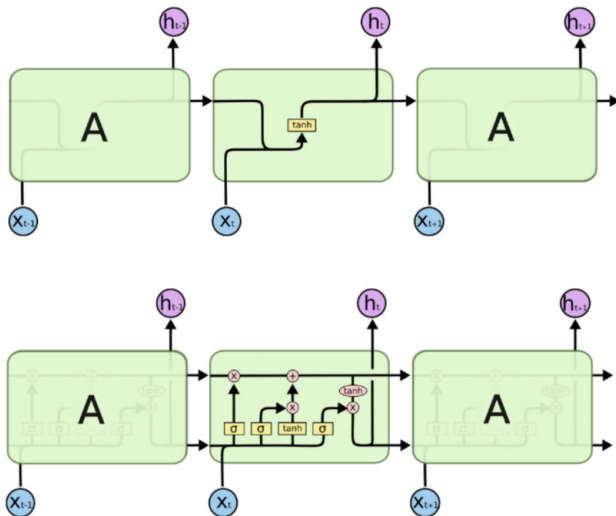
- ▶ Unless  $\sum_{k'=1}^H \sigma' \cdot W_{k'k} \approx 1$ , the gradient either vanishes, or explodes.
- ▶ For a large  $T$  (long-term dependency), the gradient "deeper" in the past tends to be too small (large).
- ▶ A solution: LSTM etc.

LSTM is currently a bit obsolete. The main idea is to decompose  $W$  into several matrices, each responsible for a different task. One is concerned about memory, one is concerned about the output at each step, etc.

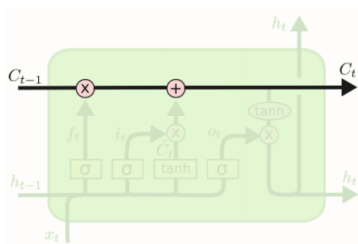
$\vec{h}_t = \vec{o}_t \circ \sigma_h(\vec{C}_t)$	output
$\vec{C}_t = \vec{f}_t \circ \vec{C}_{t-1} + \vec{i}_t \circ \tilde{C}_t$	memory
$\tilde{C}_t = \sigma_h(W_C \cdot \vec{h}_{t-1} + U_C \cdot \vec{x}_t)$	new memory contents
$\vec{o}_t = \sigma_g(W_o \cdot \vec{h}_{t-1} + U_o \cdot \vec{x}_t)$	output gate
$\vec{f}_t = \sigma_g(W_f \cdot \vec{h}_{t-1} + U_f \cdot \vec{x}_t)$	forget gate
$\vec{i}_t = \sigma_g(W_i \cdot \vec{h}_{t-1} + U_i \cdot \vec{x}_t)$	input gate

- ▶  $\circ$  is the component-wise product of vectors
- ▶  $\cdot$  is the matrix-vector product
- ▶  $\sigma_h$  hyperbolic tangents (applied component-wise)
- ▶  $\sigma_g$  logistic sigmoid (applied component-wise)

# RNN vs LSTM



Source: <https://colah.github.io/posts/2015-08-Understanding-LSTMs/>



$$\vec{h}_t = \vec{o}_t \circ \sigma_h(\vec{C}_t)$$

$$\Rightarrow \vec{C}_t = \vec{f}_t \circ \vec{C}_{t-1} + \vec{i}_t \circ \tilde{C}_t$$

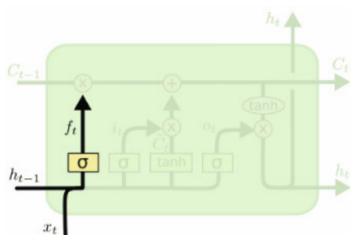
$$\tilde{C}_t = \sigma_h(W_C \cdot \vec{h}_{t-1} + U_C \cdot \vec{x}_t)$$

$$\vec{o}_t = \sigma_g(W_o \cdot \vec{h}_{t-1} + U_o \cdot \vec{x}_t)$$

$$\vec{f}_t = \sigma_g(W_f \cdot \vec{h}_{t-1} + U_f \cdot \vec{x}_t)$$

$$\vec{i}_t = \sigma_g(W_i \cdot \vec{h}_{t-1} + U_i \cdot \vec{x}_t)$$

Source: <https://colah.github.io/posts/2015-08-Understanding-LSTMs/>



$$\vec{h}_t = \vec{o}_t \circ \sigma_h(\vec{C}_t)$$

$$\vec{C}_t = \vec{f}_t \circ \vec{C}_{t-1} + \vec{i}_t \circ \vec{C}_t^{\tilde{}}$$

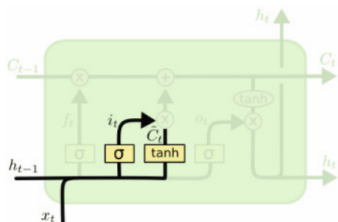
$$\vec{C}_t^{\tilde{}} = \sigma_h(W_C \cdot \vec{h}_{t-1} + U_C \cdot \vec{x}_t)$$

$$\vec{o}_t = \sigma_g(W_o \cdot \vec{h}_{t-1} + U_o \cdot \vec{x}_t)$$

$$\Rightarrow \vec{f}_t = \sigma_g(W_f \cdot \vec{h}_{t-1} + U_f \cdot \vec{x}_t)$$

$$\vec{i}_t = \sigma_g(W_i \cdot \vec{h}_{t-1} + U_i \cdot \vec{x}_t)$$

Source: <https://colah.github.io/posts/2015-08-Understanding-LSTMs/>



$$\vec{h}_t = \vec{o}_t \circ \sigma_h(\vec{C}_t)$$

$$\vec{C}_t = \vec{f}_t \circ \vec{C}_{t-1} + \vec{i}_t \circ \tilde{C}_t$$

$$\Rightarrow \tilde{C}_t = \sigma_h(W_C \cdot \vec{h}_{t-1} + U_C \cdot \vec{x}_t)$$

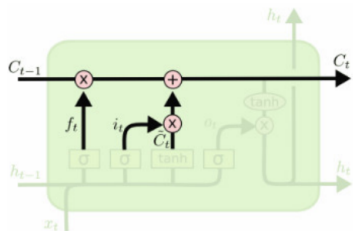
$$\vec{o}_t = \sigma_g(W_o \cdot \vec{h}_{t-1} + U_o \cdot \vec{x}_t)$$

$$\vec{f}_t = \sigma_g(W_f \cdot \vec{h}_{t-1} + U_f \cdot \vec{x}_t)$$

$$\Rightarrow \vec{i}_t = \sigma_g(W_i \cdot \vec{h}_{t-1} + U_i \cdot \vec{x}_t)$$

Source: <https://colah.github.io/posts/2015-08-Understanding-LSTMs/>





$$\vec{h}_t = \vec{o}_t \circ \sigma_h(\vec{C}_t)$$

$$\Rightarrow \vec{C}_t = \vec{f}_t \circ \vec{C}_{t-1} + \vec{i}_t \circ \vec{C}_t$$

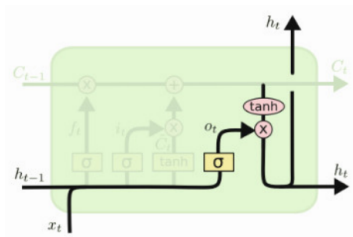
$$\Rightarrow \vec{C}_t = \sigma_h(W_C \cdot \vec{h}_{t-1} + U_C \cdot \vec{x}_t)$$

$$\vec{o}_t = \sigma_g(W_o \cdot \vec{h}_{t-1} + U_o \cdot \vec{x}_t)$$

$$\Rightarrow \vec{f}_t = \sigma_g(W_f \cdot \vec{h}_{t-1} + U_f \cdot \vec{x}_t)$$

$$\Rightarrow \vec{i}_t = \sigma_g(W_i \cdot \vec{h}_{t-1} + U_i \cdot \vec{x}_t)$$

Source: <https://colah.github.io/posts/2015-08-Understanding-LSTMs/>



$$\Rightarrow \vec{h}_t = \vec{o}_t \circ \sigma_h(\vec{C}_t)$$

$$\vec{C}_t = \vec{f}_t \circ \vec{C}_{t-1} + \vec{i}_t \circ \tilde{C}_t$$

$$\tilde{C}_t = \sigma_h(W_C \cdot \vec{h}_{t-1} + U_C \cdot \vec{x}_t)$$

$$\Rightarrow \vec{o}_t = \sigma_g(W_o \cdot \vec{h}_{t-1} + U_o \cdot \vec{x}_t)$$

$$\vec{f}_t = \sigma_g(W_f \cdot \vec{h}_{t-1} + U_f \cdot \vec{x}_t)$$

$$\vec{i}_t = \sigma_g(W_i \cdot \vec{h}_{t-1} + U_i \cdot \vec{x}_t)$$

Source: <https://colah.github.io/posts/2015-08-Understanding-LSTMs/>

- ▶ LSTM (almost) solves the vanishing gradient problem w.r.t. the "internal" state of the network.
- ▶ Learns to control its own memory (via forget gate).
- ▶ Revolution in machine translation and text processing.

... but the development goes on ...

# Time-series Forecasting with LSTM

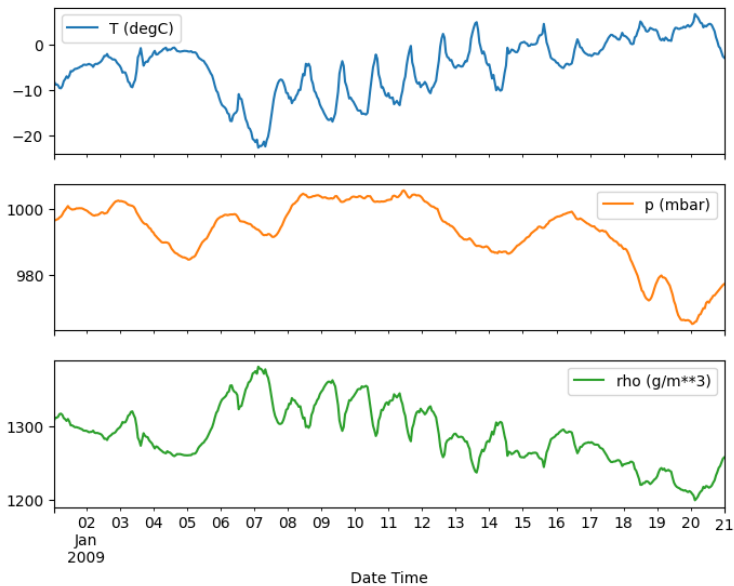
(see: [https://www.tensorflow.org/tutorials/structured\\_data/time\\_series](https://www.tensorflow.org/tutorials/structured_data/time_series))

- ▶ Weather time series dataset
- ▶ 14 different features such as air temperature, atmospheric pressure, and humidity
- ▶ collected every 10 minutes, beginning in 2003 (only 2009 - 2016 considered in the example)

	p (mbar)	T (degC)	Tpot (K)	Tdew (degC)	rh (%)	VPmax (mbar)	VPact (mbar)	VPdef (mbar)	sh (g/kg)	H2OC (mmol/mol)	rho (g/m**3)	wv (m/s)	max. wv (m/s)	wd (deg)
5	996.50	-8.05	265.38	-8.78	94.4	3.33	3.14	0.19	1.96	3.15	1307.86	0.21	0.63	192.7
11	996.62	-8.88	264.54	-9.77	93.2	3.12	2.90	0.21	1.81	2.91	1312.25	0.25	0.63	190.3
17	996.84	-8.81	264.59	-9.66	93.5	3.13	2.93	0.20	1.83	2.94	1312.18	0.18	0.63	167.2
23	996.99	-9.05	264.34	-10.02	92.6	3.07	2.85	0.23	1.78	2.85	1313.61	0.10	0.38	240.0
29	997.46	-9.63	263.72	-10.65	92.2	2.94	2.71	0.23	1.69	2.71	1317.19	0.40	0.88	157.0

**The Task:** Predict the temperature for the next hour.

# Weather Data



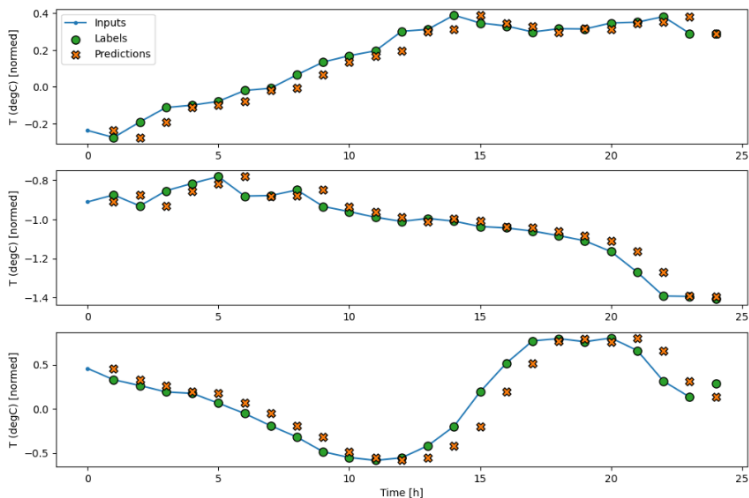
# Preprocessing (omitted)

Before applying any prediction model, proper preprocessing is essential for time series data.

- ▶ Train-validation-test Split
- ▶ Data Normalization
  
- ▶ Detrending
- ▶ Seasonal Adjustment
- ▶ Smoothing

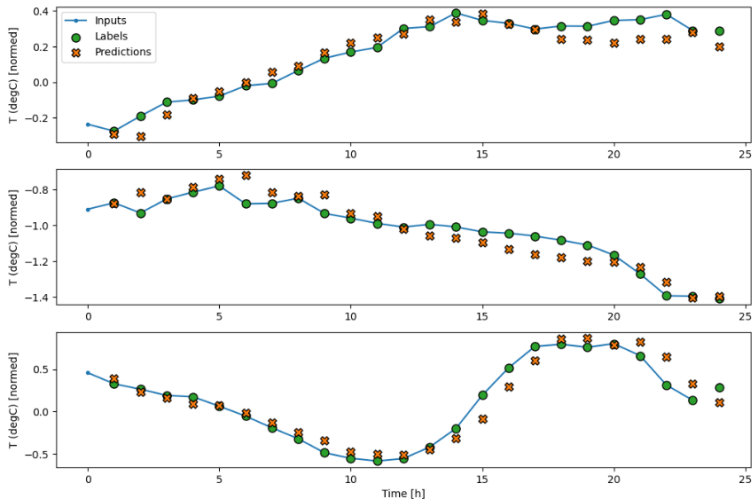
# Baseline model

The **baseline**: Predict that the temperature stays constant.



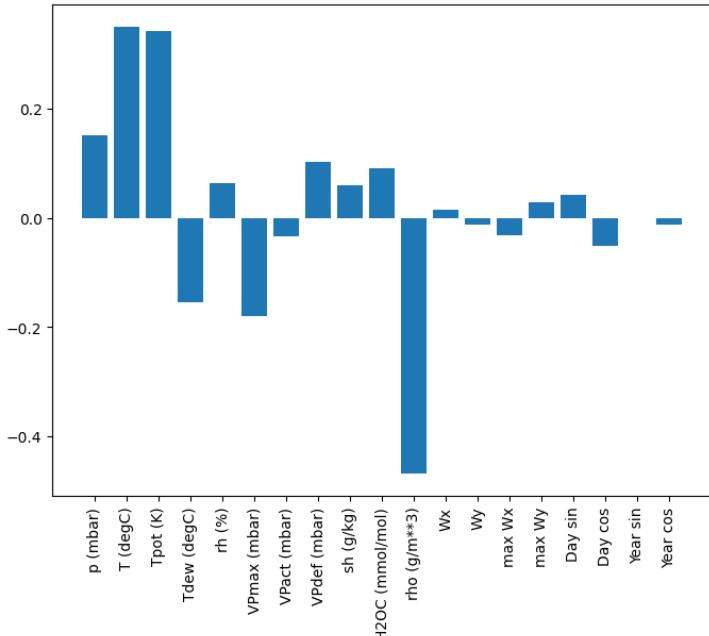
# Simple linear model

**The linear model:** Consider the current values of all variables and predict the temperature using linear regression.





# Linear explained - weights

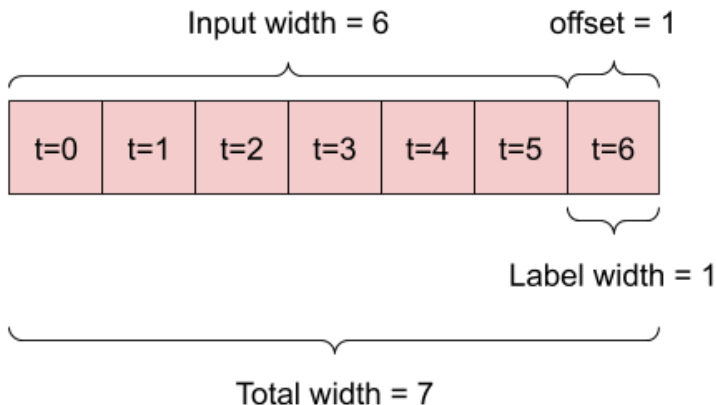


## Data Windowing

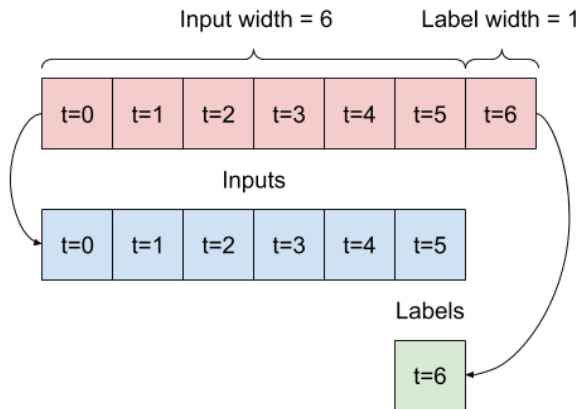
Assume that the samples are taken hourly (subsample the 10-minute samples).

Consider windowed inputs to the model.

E.g., predict one hour given 6 hours from the past:



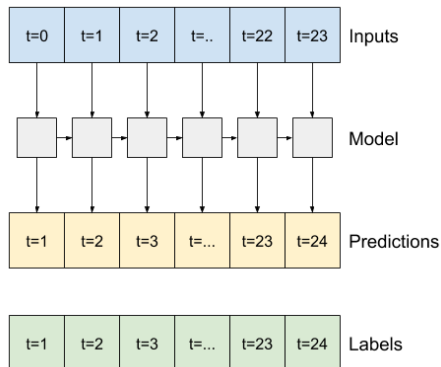
# Data Windowing



# LSTM

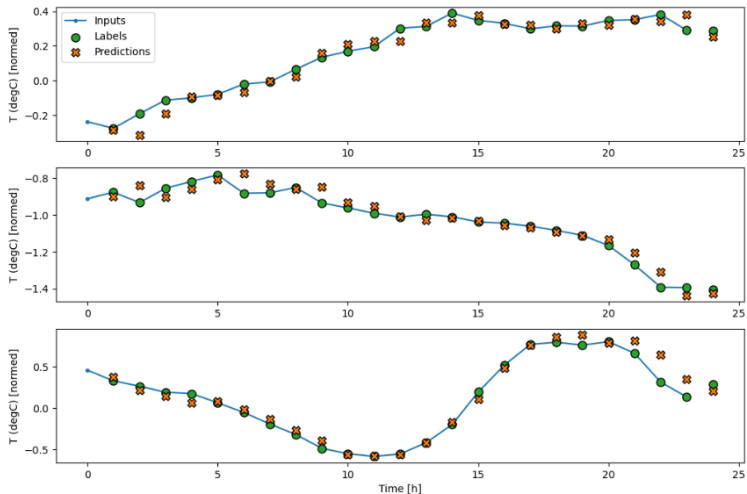
Given 24 hours in the past, predict the next hour with LSTM.

A possible LSTM architecture:



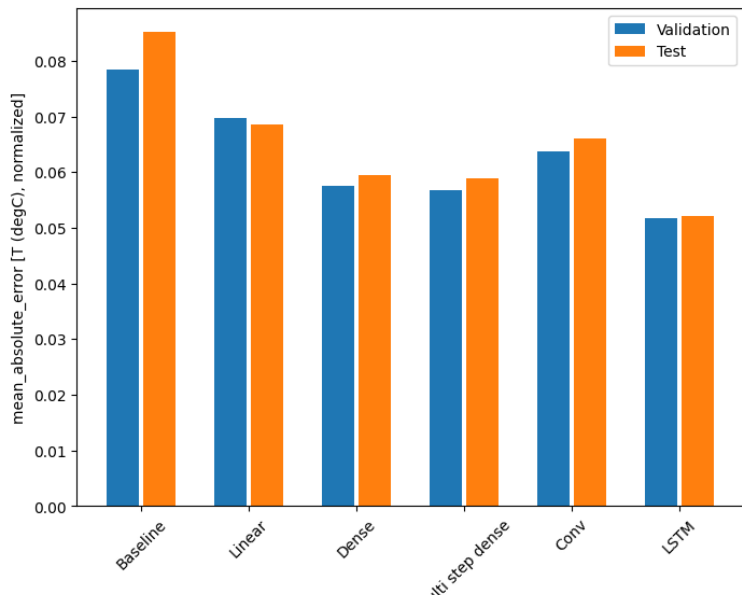
The used LSTM had the memory dimension equal to 32.

# LSTM forecasting



# Model comparison

MAE = mean absolute error

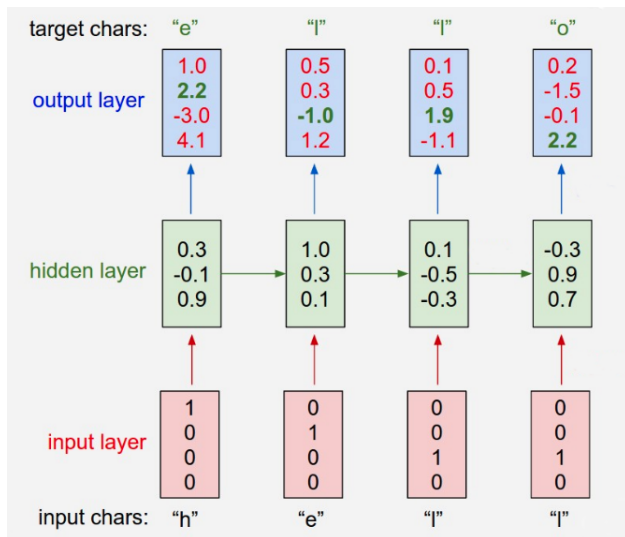


# Time-series prediction summary

- ▶ The presented approach is very basic!
- ▶ Omitted lots of important ideas:
  - ▶ Preprocessing - extremely important!
  - ▶ Cross-validation - tricky!
  - ▶ Classical models ARIMA etc. - very deep and advanced area!
  - ▶ Proper evaluation, explainability, ... (a whole new course possible!)
- ▶ Read books, e.g.
  - ▶ Hyndman and Athanasopoulos. *Forecasting: Principles and Practice*. Online: <https://otexts.com/fpp2/>
  - ▶ Manu Joseph. *Modern Time Series Forecasting with Python*. Packt Publishing. 2022

# RNN text generator (a little obsolete example)

Generating texts letter by letter.





- ▶ Generating Shakespeare letter by letter.
- ▶ Trained on Shakespeare's plays (4.4MB).

VIOLA: Why, Salisbury must find his flesh and thought That which I am not apt, not a man and in fire, To show the reining of the raven and the wars To grace my hand reproach within, and not a fair are hand, That Caesar and my goodly father's world; When I was heaven of presence and our fleets, We spare with hours, but cut thy council I am great, Murdered and by thy master's ready there My power to give thee but so much as hell: Some service in the noble bondman here, Would show him to her wine.

KING LEAR: O, if you were a feeble sight, the courtesy of your law, Your sight and several breath, will wear the gods With his heads, and my hands are wonder'd at the deeds, So drop upon your lordship's head, and your opinion Shall be against your honour.

## ▶ Hutter Prize 100MB dataset from Wikipedia (96MB)

Naturalism and decision for the majority of Arab countries' capitalide was grounded by the Irish language by [[John Clair]], [[An Imperial Japanese Revolt]], associated with Guangzham's sovereignty. His generals were the powerful ruler of the Portugal in the [[Protestant Immineners]], which could be said to be directly in Cantonese Communication, which followed a ceremony and set inspired prison, training. The emperor travelled back to [[Antioch, Perth, October 25|21]] to note, the Kingdom of Costa Rica, unsuccessful fashioned the [[Thrales]], [[Cynth's Dajoard]], known in western [[Scotland]], near Italy to the conquest of India with the conflict. Copyright was the succession of independence in the slop of Syrian influence that was a famous German movement based on a more popular servicious, non-doctrinal and sexual power post. Many governments recognize the military housing of the [[Civil Liberalization and Infantry Resolution 265 National Party in Hungary]], that is sympathetic to be to the [[Punjab Resolution]] (PJS)[http:

//www.humah.yahoo.com/guardian.cfm/7754800786d17551963s89.htm

Official economics Adjoint for the Nazism, Montgomery was swear to advance to the resources for those Socialism's rule. was starting to signing a

## Xml hallucination:

```
<page>
  <title>Antichrist</title>
  <id>865</id>
  <revision>
    <id>15900676</id>
    <timestamp>2002-08-03T18:14:12Z</timestamp>
    <contributor>
      <username>Paris</username>
      <id>23</id>
    </contributor>
    <minor />
    <comment>Automated conversion</comment>
    <text xml:space="preserve">
      #REDIRECT [[Christianity]]</text>
    </revision>
  </page>
```

- ▶ Algebraic geometry textbook.
- ▶ LaTeX source (16MB).
- ▶ Almost compilable.

*Proof.* Omitted. □

**Lemma 0.1.** *Let  $\mathcal{C}$  be a set of the construction.*

*Let  $\mathcal{C}$  be a gerber covering. Let  $\mathcal{F}$  be a quasi-coherent sheaves of  $\mathcal{O}$ -modules. We have to show that*

$$\mathcal{O}_{\mathcal{O}_X} = \mathcal{O}_X(\mathcal{L})$$

*Proof.* This is an algebraic space with the composition of sheaves  $\mathcal{F}$  on  $X_{\acute{e}tale}$  we have

$$\mathcal{O}_X(\mathcal{F}) = \{morph_1 \times_{\mathcal{O}_X} (\mathcal{G}, \mathcal{F})\}$$

where  $\mathcal{G}$  defines an isomorphism  $\mathcal{F} \rightarrow \mathcal{F}$  of  $\mathcal{O}$ -modules. □

**Lemma 0.2.** *This is an integer  $\mathcal{Z}$  is injective.*

*Proof.* See Spaces, Lemma ?? □

**Lemma 0.3.** *Let  $S$  be a scheme. Let  $X$  be a scheme and  $X$  is an affine open covering. Let  $\mathcal{U} \subset \mathcal{X}$  be a canonical and locally of finite type. Let  $X$  be a scheme. Let  $X$  be a scheme which is equal to the formal complex.*

*The following to the construction of the lemma follows.*

*Let  $X$  be a scheme. Let  $X$  be a scheme covering. Let*

$$b : X \rightarrow Y' \rightarrow Y \rightarrow Y \rightarrow Y' \times_X Y \rightarrow X.$$

*be a morphism of algebraic spaces over  $S$  and  $Y$ .*

*Proof.* Let  $X$  be a nonzero scheme of  $X$ . Let  $X$  be an algebraic space. Let  $\mathcal{F}$  be a quasi-coherent sheaf of  $\mathcal{O}_X$ -modules. The following are equivalent

- (1)  $\mathcal{F}$  is an algebraic space over  $S$ .
- (2) If  $X$  is an affine open covering.

Consider a common structure on  $X$  and  $X$  the functor  $\mathcal{O}_X(U)$  which is locally of

- ▶ Trained on all source files of Linux kernel concatenated into a single file (474MB of C code).

```

/*
 * Increment the size file of the new incorrect UI_FILTER group information
 * of the size generatively.
 */
static int indicate_policy(void)
{
    int error;
    if (fd == MARN_EPT) {
        /*
         * The kernel blank will coeld it to userspace.
         */
        if (ss->segment < mem_total)
            unblock_graph_and_set_blocked();
        else
            ret = 1;
        goto bail;
    }
    segaddr = in_SB(in.addr);
    selector = seg / 16;
    setup_works = true;
    for (i = 0; i < blocks; i++) {
        seq = buf[i++];
        bpf = bd->bd.next + i * search;
        if (fd) {
            current = blocked;
        }
    }
    rw->name = "Getjbbregs";
    bprm_self_clearl(&iv->version);
    regs->new = blocks[(BPF_STATS << info->historidac) | PFMR_CLOBATHINC_SECONDS << 12];
    return segtable;
}

```

```
/*
 * Copyright (c) 2006-2010, Intel Mobile Communications. All rights reserved.
 *
 * This program is free software; you can redistribute it and/or modify it
 * under the terms of the GNU General Public License version 2 as published by
 * the Free Software Foundation.
 *
 * This program is distributed in the hope that it will be useful,
 * but WITHOUT ANY WARRANTY; without even the implied warranty of
 * MERCHANTABILITY or FITNESS FOR A PARTICULAR PURPOSE. See the
 *
 * GNU General Public License for more details.
 *
 * You should have received a copy of the GNU General Public License
 * along with this program; if not, write to the Free Software Foundation,
 * Inc., 675 Mass Ave, Cambridge, MA 02139, USA.
 */

#include <linux/kexec.h>
#include <linux/errno.h>
#include <linux/io.h>
#include <linux/platform_device.h>
#include <linux/multi.h>
#include <linux/ckevent.h>

#include <asm/io.h>
#include <asm/prom.h>
#include <asm/e820.h>
#include <asm/system_info.h>
#include <asm/setew.h>
#include <asm/pgproto.h>
```



# Evolution of Shakespeare

## 100 iter.:

tyntd-iafhatawiaoihrdemot lytdws e ,tfti, astai f ogoh eoase rrranbyne 'nhthnee e  
plia tklrqd t o idoe ns,smtt h ne etie h,hregtrs nigtiike,aoaenns lng

## 300 iter.:

"Tmont thithey" fomesscerliund  
Keushey. Thom here  
sheulke, anmerenith ol sivh I lalterthend Bleipile shuw y fil on aseterlome  
coaniogennc Phe lism thond hon at. MeiDimorotion in ther thize."

## 500 iter.:

we counter. He stutn co des. His stanted out one ofler that concossions and was  
to gearang reay Jotrets and with fre colt ofp paitt thin wall. Which das stimn

## 700 iter.:

Aftair fall unsuch that the hall for Prince Velzonski's that me of  
her hearly, and behs to so arwage fiving were to it beloge, pavu say falling misfort  
how, and Gogition is so overelical and offer.

## 1200 iter.:

"Kite vouch!" he repeated by her  
door. "But I would be done and quarts, feeling, then, son is people..."

## 2000 iter.:

"Why do what that day," replied Natasha, and wishing to himself the fact the  
princess, Princess Mary was easier, fed in had oftened him.  
Pierre aking his soul came to the packs and drove up his father-in-law women.

# Attention

Consider the following task: Given a sequence of vectors

$$\mathbf{x} = \vec{x}_1, \dots, \vec{x}_T$$

generate a new sequence

$$\mathbf{y} = \vec{y}_1, \dots, \vec{y}_{T'}$$

of possibly different lengths (i.e., possibly  $T \neq T'$ ).

E.g., a machine translation task,  $\mathbf{x}$  is an embedding of an English sentence,  $\mathbf{y}$  is a sequence of probability distributions on a German vocabulary.

# Attention

Consider two recurrent networks:

- ▶ Enc the encoder
  - ▶ Hidden state  $\vec{h}_0$  initialized by standard methods for recurrent networks
  - ▶ Reads  $\vec{x}_1, \dots, \vec{x}_T$ , does not output anything but produces a sequence of hidden states  $\vec{h}_1, \dots, \vec{h}_T$

# Attention

Consider two recurrent networks:

- ▶ Enc the encoder
  - ▶ Hidden state  $\vec{h}_0$  initialized by standard methods for recurrent networks
  - ▶ Reads  $\vec{x}_1, \dots, \vec{x}_T$ , does not output anything but produces a sequence of hidden states  $\vec{h}_1, \dots, \vec{h}_T$
- ▶ Dec the decoder
  - ▶ The initial hidden state is  $\vec{h}_T$
  - ▶ Does not read anything but outputs the sequence  $\vec{y}_1, \dots, \vec{y}_T$ .  
This is a simplification. Typically, Dec reads  $\vec{y}_0, \vec{y}_1, \dots, \vec{y}_{T-1}$  where  $\vec{y}_0$  is a special vector embedding a separator.

# Attention

Consider two recurrent networks:

- ▶ Enc the encoder
  - ▶ Hidden state  $\vec{h}_0$  initialized by standard methods for recurrent networks
  - ▶ Reads  $\vec{x}_1, \dots, \vec{x}_T$ , does not output anything but produces a sequence of hidden states  $\vec{h}_1, \dots, \vec{h}_T$
- ▶ Dec the decoder
  - ▶ The initial hidden state is  $\vec{h}_T$
  - ▶ Does not read anything but outputs the sequence  $\vec{y}_1, \dots, \vec{y}_T$   
This is a simplification. Typically, Dec reads  $\vec{y}_0, \vec{y}_1, \dots, \vec{y}_{T-1}$  where  $\vec{y}_0$  is a special vector embedding a separator.

Trained on pairs of sentences, able to learn a fine translation between major languages (if the recurrent networks are LSTM).

Is not perfect because all info about  $\mathbf{x} = \vec{x}_1, \dots, \vec{x}_T$  is squeezed into the single state vector  $\vec{h}_T$ .

In particular, the network tends to forget the context of each word.

## Attention in Recurrent Networks

What if we provide the decoder with an information about the *relevant context* of the generated word?

## Attention in Recurrent Networks

What if we provide the decoder with an information about the *relevant context* of the generated word?

We use the same encoder  $\text{Enc}$  producing the sequence of hidden states:  $\vec{h}_1, \dots, \vec{h}_T$

## Attention in Recurrent Networks

What if we provide the decoder with an information about the *relevant context* of the generated word?

We use the same encoder Enc producing the sequence of hidden states:  $\vec{h}_1, \dots, \vec{h}_T$

The decoder Dec is still a recurrent network but

- ▶ the hidden state  $\vec{h}'_0$  initialized by  $\vec{h}_T$  and a sequence of hidden states  $\vec{h}'_0, \dots, \vec{h}'_T$ , is computed,



# Attention in Recurrent Networks

What if we provide the decoder with an information about the *relevant context* of the generated word?

We use the same encoder Enc producing the sequence of hidden states:  $\vec{h}_1, \dots, \vec{h}_T$

The decoder Dec is still a recurrent network but

- ▶ the hidden state  $\vec{h}'_0$  initialized by  $\vec{h}_T$  and a sequence of hidden states  $\vec{h}'_0, \dots, \vec{h}'_{T'}$ , is computed,
- ▶ reads a sequence of context vectors  $\vec{c}_1, \dots, \vec{c}_{T'}$ , where

$$\vec{c}_i = \sum_{j=1}^T \alpha_{ij} \vec{h}_j \quad \text{where} \quad \alpha_{ij} = \frac{\exp(\mathbf{e}_{ij})}{\sum_{k=1}^T \exp(\mathbf{e}_{ik})}$$

where  $\mathbf{e}_{ij} = \text{MLP}(\vec{h}'_{i-1}, \vec{h}_j)$

- ▶ outputs the sequence  $\vec{y}_1, \dots, \vec{y}_{T'}$

# Do We Still Need the Recurrence?

- ▶ The attention mechanism extracts the information from the sequence quite well.

## Do We Still Need the Recurrence?

- ▶ The attention mechanism extracts the information from the sequence quite well.
- ▶ Is there a reason for reading the input sequence sequentially?

## Do We Still Need the Recurrence?

- ▶ The attention mechanism extracts the information from the sequence quite well.
- ▶ Is there a reason for reading the input sequence sequentially?
- ▶ Could we remove the recurrent network itself and preserve only the attention?

---

# Attention Is All You Need

---

**Ashish Vaswani\***  
Google Brain  
avaswani@google.com

**Noam Shazeer\***  
Google Brain  
noam@google.com

**Niki Parmar\***  
Google Research  
nikip@google.com

**Jakob Uszkoreit\***  
Google Research  
usz@google.com

**Llion Jones\***  
Google Research  
llion@google.com

**Aidan N. Gomez\* †**  
University of Toronto  
aidan@cs.toronto.edu

**Lukasz Kaiser\***  
Google Brain  
lukaszkaizer@google.com

**Illia Polosukhin\* ‡**  
illia.polosukhin@gmail.com

## Abstract

The dominant sequence transduction models are based on complex recurrent or convolutional neural networks that include an encoder and a decoder. The best performing models also connect the encoder and decoder through an attention mechanism. We propose a new simple network architecture, the Transformer, based solely on attention mechanisms, dispensing with recurrence and convolutions entirely. Experiments on two machine translation tasks show these models to be superior in quality while being more parallelizable and requiring significantly less time to train. Our model achieves 28.4 BLEU on the WMT 2014 English-to-German translation task, improving over the existing best results, including ensembles, by over 2 BLEU. On the WMT 2014 English-to-French translation task, our model establishes a new single-model state-of-the-art BLEU score of 41.0 after training for 3.5 days on eight GPUs, a small fraction of the training costs of the best models from the literature.

# Self-Attention Layer (is all you need)

Fix an input sequence:  $\vec{x}_1, \dots, \vec{x}_T$

Consider three learnable matrices:  $W_q, W_k, W_v$

Generate sequences of *queries*, *keys*, and *values*:

- ▶  $\vec{q}_1, \dots, \vec{q}_T$  where  $\vec{q}_i = W_q \vec{x}_i$  for all  $i = 1, \dots, T$
- ▶  $\vec{k}_1, \dots, \vec{k}_T$  where  $\vec{k}_i = W_k \vec{x}_i$  for all  $i = 1, \dots, T$
- ▶  $\vec{v}_1, \dots, \vec{v}_T$  where  $\vec{v}_i = W_v \vec{x}_i$  for all  $i = 1, \dots, T$

# Self-Attention Layer (is all you need)

Fix an input sequence:  $\vec{x}_1, \dots, \vec{x}_T$

Consider three learnable matrices:  $W_q, W_k, W_v$

Generate sequences of *queries*, *keys*, and *values*:

- ▶  $\vec{q}_1, \dots, \vec{q}_T$  where  $\vec{q}_i = W_q \vec{x}_i$  for all  $i = 1, \dots, T$
- ▶  $\vec{k}_1, \dots, \vec{k}_T$  where  $\vec{k}_i = W_k \vec{x}_i$  for all  $i = 1, \dots, T$
- ▶  $\vec{v}_1, \dots, \vec{v}_T$  where  $\vec{v}_i = W_v \vec{x}_i$  for all  $i = 1, \dots, T$

Define a vector score for all  $i, j \in \{1, \dots, T\}$  by

$$e_{ij} = \vec{q}_i \cdot \vec{k}_j$$

Intuitively,  $e_{ij}$  measures how much the input at the position  $i$  is related to the input at the position  $j$ , in other words, how much the query fits the key.

Define

$$\alpha_{ij} = \frac{\exp(e_{ij} / \sqrt{d_{\text{attn}}})}{\sum_{k=1}^T \exp(e_{ik} / \sqrt{d_{\text{attn}}})} \quad d_{\text{attn}} \text{ is the dimension of } \vec{v}_i$$

I.e., we apply the good old softmax to  $(e_{i1}, \dots, e_{iT}) / \sqrt{d_{\text{attn}}}$

# Self-Attention Layer (is all you need)

Define a vector score for all  $i, j \in \{1, \dots, T\}$  by

$$e_{ij} = \vec{q}_i \cdot \vec{k}_j$$

Intuitively,  $e_{ij}$  measures how much the input at the position  $i$  is related to the input at the position  $j$ , in other words, how much the query fits the key.

Define

$$\alpha_{ij} = \frac{\exp(e_{ij} / \sqrt{d_{\text{attn}}})}{\sum_{k=1}^T \exp(e_{ik} / \sqrt{d_{\text{attn}}})} \quad d_{\text{attn}} \text{ is the dimension of } \vec{v}_i$$

I.e., we apply the good old softmax to  $(e_{i1}, \dots, e_{iT}) / \sqrt{d_{\text{attn}}}$

Define a sequence of outputs  $\vec{y}_1, \dots, \vec{y}_T$  by

$$\vec{y}_i = \sum_{j=1}^T \alpha_{ij} \cdot \vec{v}_j$$



# Language Model

A sequence of *tokens*  $a_1, \dots, a_T \in \Sigma^*$

E.g. words from a vocabulary  $\Sigma$ .

**The goal:** Maximize

$$\prod_{k=1}^T P(a_k \mid a_1, \dots, a_{k-1}; W) \quad (= P(a_1, \dots, a_T; W))$$

where

- ▶  $P$  is the conditional probability measure over  $\Sigma$  modeled using a neural network with weights  $W$ .

# Language Model

A sequence of *tokens*  $a_1, \dots, a_T \in \Sigma^*$

E.g. words from a vocabulary  $\Sigma$ .

**The goal:** Maximize

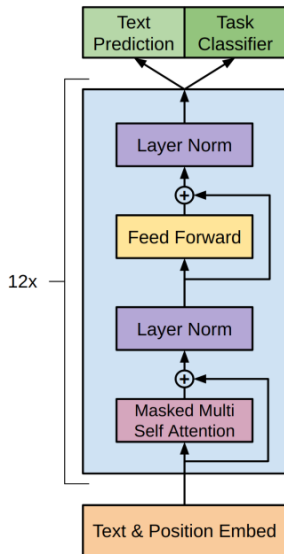
$$\prod_{k=1}^T P(a_k \mid a_1, \dots, a_{k-1}; W) \quad (= P(a_1, \dots, a_T; W))$$

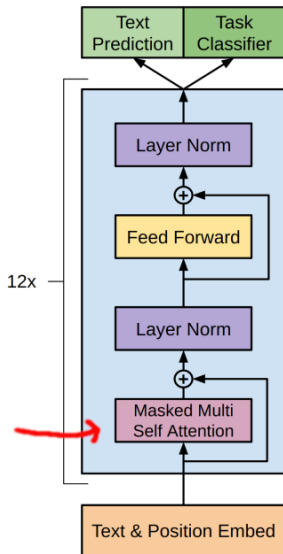
where

- ▶  $P$  is the conditional probability measure over  $\Sigma$  modeled using a neural network with weights  $W$ .

Can be used to generate text:

Given  $a_1, \dots, a_k$ , sample  $a_{k+1}$  from  $P(a_{k+1} \mid a_1, \dots, a_k; W)$





## Masked Self-Attention Layer (is all you need)

Assume an attention mechanism which given an input sequence  $\vec{x}_1, \dots, \vec{x}_T$  generates  $\vec{y}_1, \dots, \vec{y}_T$ .

**The Problem:** How to generate  $\vec{y}_k$  only based on  $\vec{x}_1, \dots, \vec{x}_{k-1}$  ?

# Masked Self-Attention Layer (is all you need)

Assume an attention mechanism which given an input sequence  $\vec{x}_1, \dots, \vec{x}_T$  generates  $\vec{y}_1, \dots, \vec{y}_T$ .

**The Problem:** How to generate  $\vec{y}_k$  only based on  $\vec{x}_1, \dots, \vec{x}_{k-1}$  ?

Define a vector score for all  $i, j \in \{1, \dots, T\}$  by

$$e_{ij} = \begin{cases} \vec{q}_i \cdot \vec{k}_j & \text{if } j < i \\ -\infty & \text{otherwise.} \end{cases}$$

This means that

$$\alpha_{ij} = \begin{cases} \frac{\exp(e_{ij} / \sqrt{d_{\text{attn}}})}{\sum_{k=1}^T \exp(e_{ik} / \sqrt{d_{\text{attn}}})} & \text{if } j < i \\ 0 & \text{otherwise.} \end{cases}$$

# Masked Self-Attention Layer (is all you need)

Assume an attention mechanism which given an input sequence  $\vec{x}_1, \dots, \vec{x}_T$  generates  $\vec{y}_1, \dots, \vec{y}_T$ .

**The Problem:** How to generate  $\vec{y}_k$  only based on  $\vec{x}_1, \dots, \vec{x}_{k-1}$  ?

Define a vector score for all  $i, j \in \{1, \dots, T\}$  by

$$e_{ij} = \begin{cases} \vec{q}_i \cdot \vec{k}_j & \text{if } j < i \\ -\infty & \text{otherwise.} \end{cases}$$

This means that

$$\alpha_{ij} = \begin{cases} \frac{\exp(e_{ij} / \sqrt{d_{\text{attn}}})}{\sum_{k=1}^T \exp(e_{ik} / \sqrt{d_{\text{attn}}})} & \text{if } j < i \\ 0 & \text{otherwise.} \end{cases}$$

Define a sequence of outputs  $\vec{y}_1, \dots, \vec{y}_T$  by

$$\vec{y}_i = \sum_{j=1}^T \alpha_{ij} \cdot \vec{v}_j$$

# Multi-head Self-Attention Layer (is all you need)

Assume the number of *heads* is  $H$ .

For  $h = 1, \dots, H$  the  $h$ -th head is an attention mechanism which given the input  $\vec{x}_1, \dots, \vec{x}_T$  produces

$$\vec{y}_1^h, \dots, \vec{y}_T^h$$

Note that the output may be different, which means that, in particular, the matrices  $W_q, W_k, W_v$  may be different for each head.

Assume that all vectors  $\vec{y}_k^h$  are of the same dimension  $d_{mid}$  and consider a learnable matrix  $W_{out}$  of dimensions  $d_{out} \times (H \cdot d_{mid})$ .



# Multi-head Self-Attention Layer (is all you need)

Assume the number of *heads* is  $H$ .

For  $h = 1, \dots, H$  the  $h$ -th head is an attention mechanism which given the input  $\vec{x}_1, \dots, \vec{x}_T$  produces

$$\vec{y}_1^h, \dots, \vec{y}_T^h$$

Note that the output may be different, which means that, in particular, the matrices  $W_q, W_k, W_v$  may be different for each head.

Assume that all vectors  $\vec{y}_k^h$  are of the same dimension  $d_{mid}$  and consider a learnable matrix  $W_{out}$  of dimensions  $d_{out} \times (H \cdot d_{mid})$ .

The multi-head attention produces the following output:

$$\vec{y}_1, \dots, \vec{y}_T$$

where

$$\vec{y}_k = W_{out} \cdot (\vec{y}_k^1 \odot \vec{y}_k^2 \odot \dots \odot \vec{y}_k^H)$$

Here,  $\odot$  is a concatenation of vectors.

# Multi-head Self-Attention Summary

Input: A sequence  $\vec{x}_1, \dots, \vec{x}_T$

Output: A sequence  $\vec{y}_1, \dots, \vec{y}_T$

I.e., a sequence of the same length. The dimensions of  $\vec{y}_k$  and  $\vec{x}_k$  do not have to be equal.

# Multi-head Self-Attention Summary

Input: A sequence  $\vec{x}_1, \dots, \vec{x}_T$

Output: A sequence  $\vec{y}_1, \dots, \vec{y}_T$

I.e., a sequence of the same length. The dimensions of  $\vec{y}_k$  and  $\vec{x}_k$  do not have to be equal.

## **Attention:**

Learnable parameters: Matrices  $W_q, W_k, W_v$ .

These matrices are used to compute queries, keys, and values from  $\vec{x}_1, \dots, \vec{x}_T$ . Output  $\vec{y}_1, \dots, \vec{y}_T$  is computed using values "scaled" by the query-key attention.

# Multi-head Self-Attention Summary

Input: A sequence  $\vec{x}_1, \dots, \vec{x}_T$

Output: A sequence  $\vec{y}_1, \dots, \vec{y}_T$

I.e., a sequence of the same length. The dimensions of  $\vec{y}_k$  and  $\vec{x}_k$  do not have to be equal.

## Attention:

Learnable parameters: Matrices  $W_q, W_k, W_v$ .

These matrices are used to compute queries, keys, and values from  $\vec{x}_1, \dots, \vec{x}_T$ . Output  $\vec{y}_1, \dots, \vec{y}_T$  is computed using values "scaled" by the query-key attention.

## Multi-head attention:

Learnable parameters:

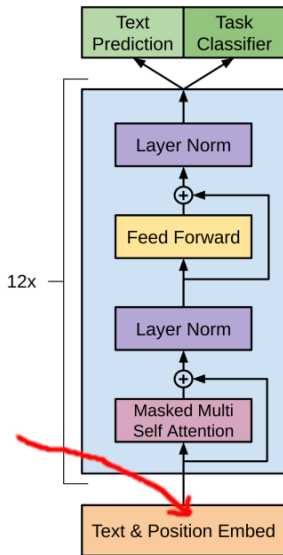
- ▶ Matrices  $W_q^h, W_k^h, W_v^h$  where  $h = 1, \dots, H$  and  $H$  is the number of heads.

Each attention head operates independently on the input  $\vec{x}_1, \dots, \vec{x}_T$ .

- ▶ Matrix  $W_{out}$ .

Linearly transforms the concatenated results of the attention heads.

# GPT - transformer



## Positional encoding

**The Goal:** To encode a position (index)  $k \in \{1, \dots, T\}$  into a vector  $\vec{P}_k$  of real numbers.

## Positional encoding

**The Goal:** To encode a position (index)  $k \in \{1, \dots, T\}$  into a vector  $\vec{P}_k$  of real numbers.

Assume that  $\vec{P}_k$  should have a dimension  $d$ .

Given a position  $k \in \{1, \dots, T\}$  and  $i \in \{0, \dots, d/2\}$  define

$$P_{k,2i} = \sin\left(\frac{k}{n^{2i/d}}\right)$$

$$P_{k,(2i+1)} = \cos\left(\frac{k}{n^{2i/d}}\right)$$

Here  $n = 10000$ .

A user defined constant, the original paper suggests  $n = 10000$ .

## Positional encoding

**The Goal:** To encode a position (index)  $k \in \{1, \dots, T\}$  into a vector  $\vec{P}_k$  of real numbers.

Assume that  $\vec{P}_k$  should have a dimension  $d$ .

Given a position  $k \in \{1, \dots, T\}$  and  $i \in \{0, \dots, d/2\}$  define

$$P_{k,2i} = \sin\left(\frac{k}{n^{2i/d}}\right)$$

$$P_{k,(2i+1)} = \cos\left(\frac{k}{n^{2i/d}}\right)$$

Here  $n = 10000$ .

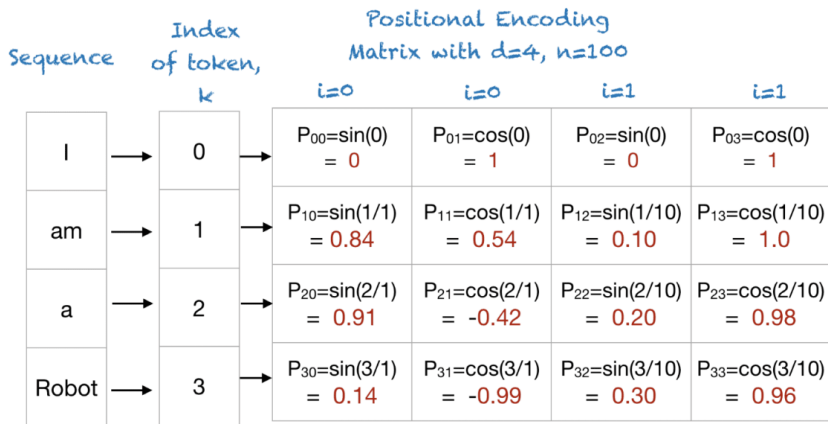
A user defined constant, the original paper suggests  $n = 10000$ .

Given an input sequence  $\vec{x}_1, \dots, \vec{x}_T$  we add the position embedding to each  $\vec{x}_k$  obtaining a new input sequence  $\vec{x}'_1, \dots, \vec{x}'_T$  where

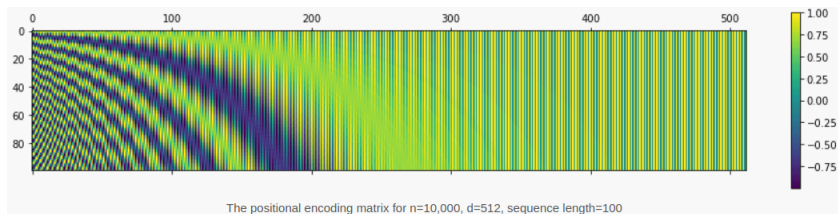
$$\vec{x}'_k = \vec{x}_k + \vec{P}_k$$



# Positional encoding/embedding



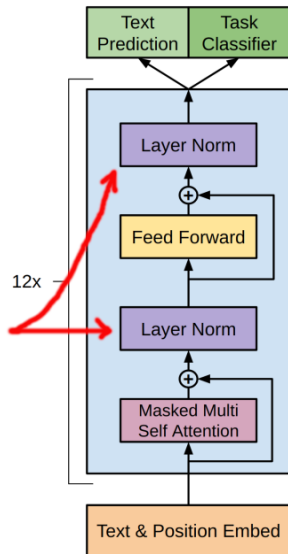
# Positional encoding/embedding



- ▶ Vertically: Sinusoidal functions
- ▶ Horizontally: Decreasing frequency

For any offset  $o \in \{1, \dots, T\}$  there is a linear transformation  $M$  such that for any  $k \in \{1, \dots, T - o\}$  we have  $M\vec{P}_k = \vec{P}_{k+o}$ .  
Intuitively, just rotate each component of the  $\vec{P}_k$  appropriately.

# GPT-2 - transformer



# Layer normalization

Given a vector  $\vec{x} \in \mathbb{R}^d$ , the *layer normalization* computes:

$$\vec{x}' = \gamma \cdot \frac{(\vec{x} - \mu)}{\sigma} + \beta$$

Here

- ▶  $\mu = \frac{1}{d} \sum_{i=1}^d x_i$  and  $\sigma^2 = \frac{1}{d} \sum_{i=1}^d (x_i - \mu)^2$
- ▶  $\gamma, \beta \in \mathbb{R}^d$  are vectors of trainable parameters

# Layer normalization

Given a vector  $\vec{x} \in \mathbb{R}^d$ , the *layer normalization* computes:

$$\vec{x}' = \gamma \cdot \frac{(\vec{x} - \mu)}{\sigma} + \beta$$

Here

- ▶  $\mu = \frac{1}{d} \sum_{i=1}^d x_i$  and  $\sigma^2 = \frac{1}{d} \sum_{i=1}^d (x_i - \mu)^2$
- ▶  $\gamma, \beta \in \mathbb{R}^d$  are vectors of trainable parameters

## In Transformer:

The input to the layer normalization is a sequence of vectors:  $\vec{x}_1, \dots, \vec{x}_T$ . The layer normalization is applied to each  $\vec{x}_k$ , producing a sequence of "normalized" vectors.

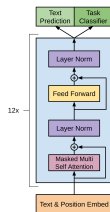
# GPT - learning

A sequence of tokens  $a_1, \dots, a_T \in \Sigma$  and their one-hot encodings  $\vec{u}_1, \dots, \vec{u}_T \in \{0, 1\}^{|\Sigma|}$

We assume that  $a_1$  is a special token marking the start of the sequence.

Embed to vectors and add the position encoding ( $W_e$  is an embedding matrix):

$$\vec{x}_k = W_e \cdot \vec{u}_k + P_k \in \mathbb{R}^{set^d}$$



# GPT - learning

A sequence of tokens  $a_1, \dots, a_T \in \Sigma$  and their one-hot encodings  $\vec{u}_1, \dots, \vec{u}_T \in \{0, 1\}^{|\Sigma|}$

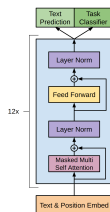
We assume that  $a_1$  is a special token marking the start of the sequence.

Embed to vectors and add the position encoding ( $W_e$  is an embedding matrix):

$$\vec{x}_k = W_e \cdot \vec{u}_k + P_k \in \mathbb{R}^{set^d}$$

Apply the network (with the transformer block repeated 12x) to  $\vec{x}_1, \dots, \vec{x}_T$  and obtain  $\vec{y}_1, \dots, \vec{y}_T$

(Here assume that each  $\vec{y}_k \in [0, 1]^{|\Sigma|}$  is a probability distribution on  $\Sigma$ )



# GPT - learning

A sequence of tokens  $a_1, \dots, a_T \in \Sigma$  and their one-hot encodings  $\vec{u}_1, \dots, \vec{u}_T \in \{0, 1\}^{|\Sigma|}$

We assume that  $a_1$  is a special token marking the start of the sequence.

Embed to vectors and add the position encoding ( $W_e$  is an embedding matrix):

$$\vec{x}_k = W_e \cdot \vec{u}_k + P_k \in \mathbb{R}^{set^d}$$

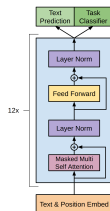
Apply the network (with the transformer block repeated 12x) to  $\vec{x}_1, \dots, \vec{x}_T$  and obtain  $\vec{y}_1, \dots, \vec{y}_T$

(Here assume that each  $\vec{y}_k \in [0, 1]^{|\Sigma|}$  is a probability distribution on  $\Sigma$ )

Compute the error:

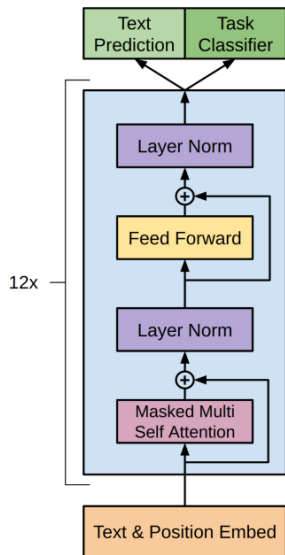
$$-\sum_{\ell=1}^{T-1} \log(\vec{y}_\ell[a_{\ell+1}])$$

Here  $\vec{y}_\ell[a_{k+1}]$  is the probability of  $a_{k+1}$  in the distribution  $\vec{y}_k$ .





# GPT - inference



<https://transformer.huggingface.co/doc/distil-gpt2>

A sequence of tokens  $a_1, \dots, a_\ell \in \Sigma$  and their one-hot encodings  $\vec{u}_1, \dots, \vec{u}_\ell \in \{0, 1\}^{|\Sigma|}$

Embed to vectors and add the position encoding:

$$\vec{x}_k = W_e \cdot \vec{u}_k + P_k \in Rset^d$$

Apply the network to  $\vec{x}_1, \dots, \vec{x}_\ell$  and obtain  $\vec{y}_1, \dots, \vec{y}_\ell$

(Assume that each  $\vec{y}_k \in [0, 1]^\Sigma$  is a probability distribution on  $\Sigma$ )

Sample the next token from

$$a_{\ell+1} \sim \vec{y}_\ell$$

# Neural networks summary

## Architectures:

- ▶ Multi-layer perceptron (MLP):
  - ▶ dense connections between layers
- ▶ Convolutional networks (CNN):
  - ▶ local receptors, feature maps
  - ▶ pooling
- ▶ Recurrent networks (RNN):
  - ▶ self-loops but still feed-forward through time
- ▶ Transformer
  - ▶ Attention, query-key-value

## Training:

- ▶ gradient descent algorithm + heuristics