

The background features a dark gray grid pattern. In the top right and bottom left corners, there are decorative wavy lines in a bright purple color, creating a modern, tech-oriented aesthetic.

iTMO

**Neural Networks
and Deep Learning**

Computer Vision

Neural Networks and Deep Learning

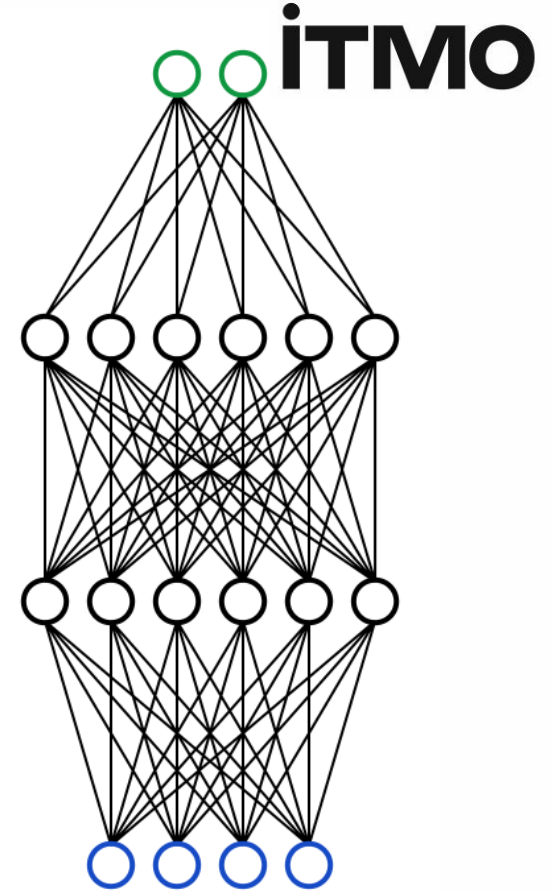
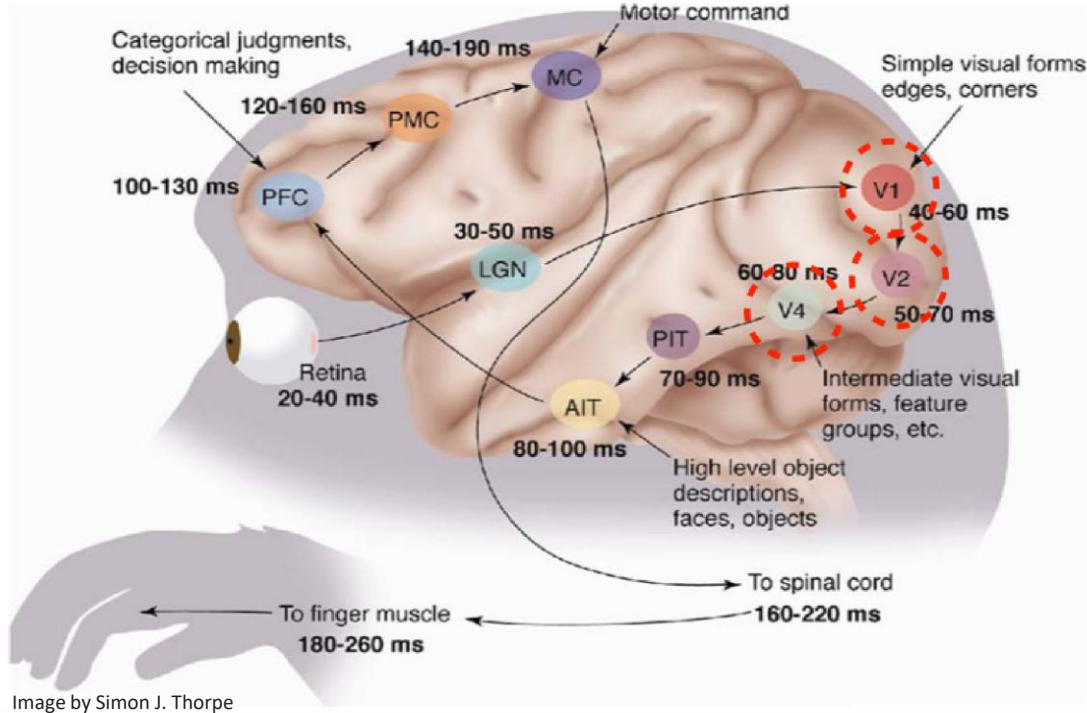


Outline

- Biological neural network
- Artificial neuron
- Artificial neural network
- Neural network training

Neural Networks and Deep Learning

Biological neural network



The number of neurons in the human brain is estimated to be around $10^{10} - 10^{11}$

Biological neural network

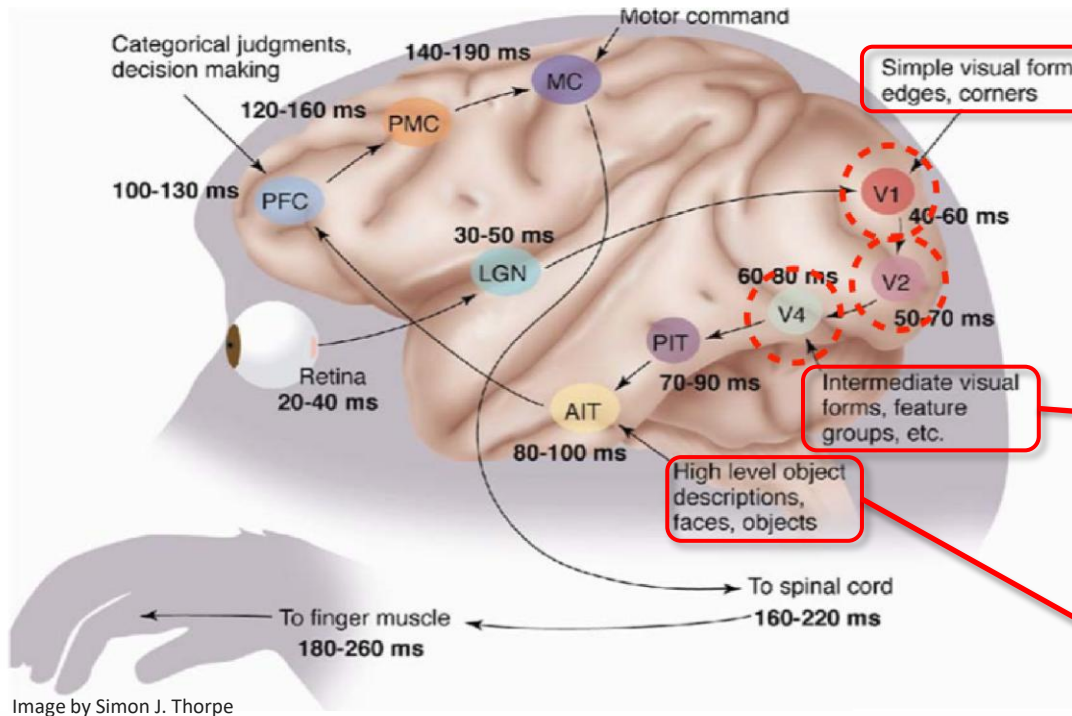
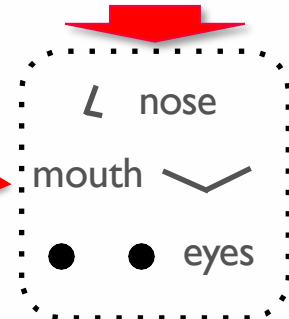
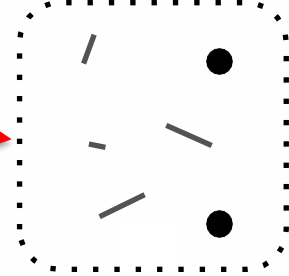


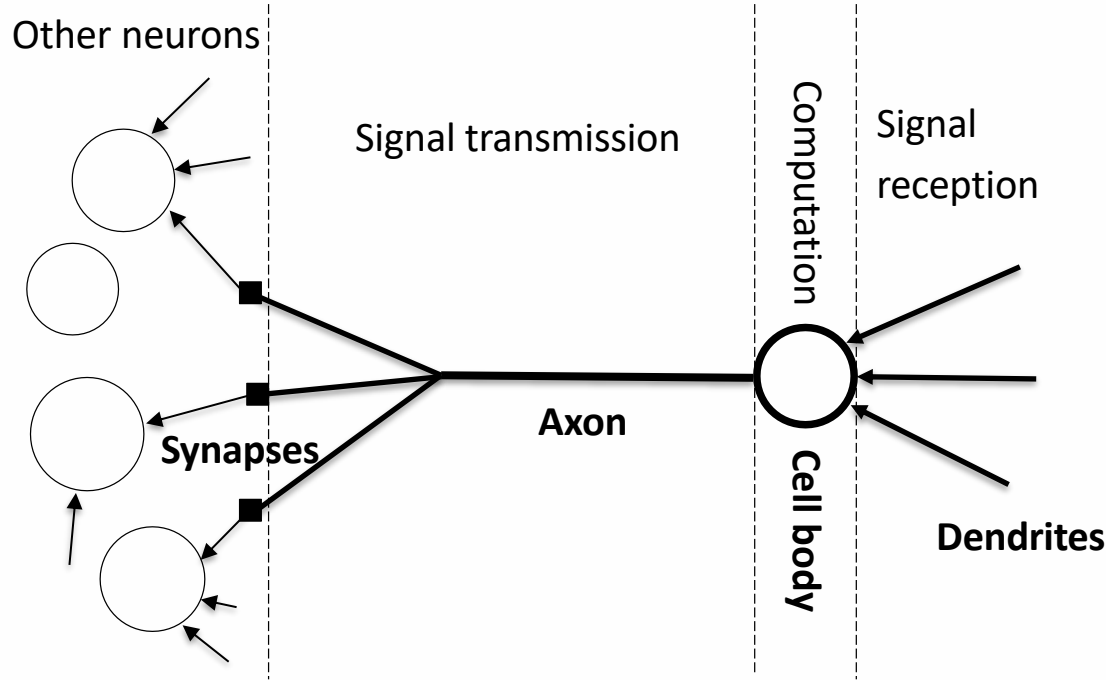
Image by Simon J. Thorpe

ITMO



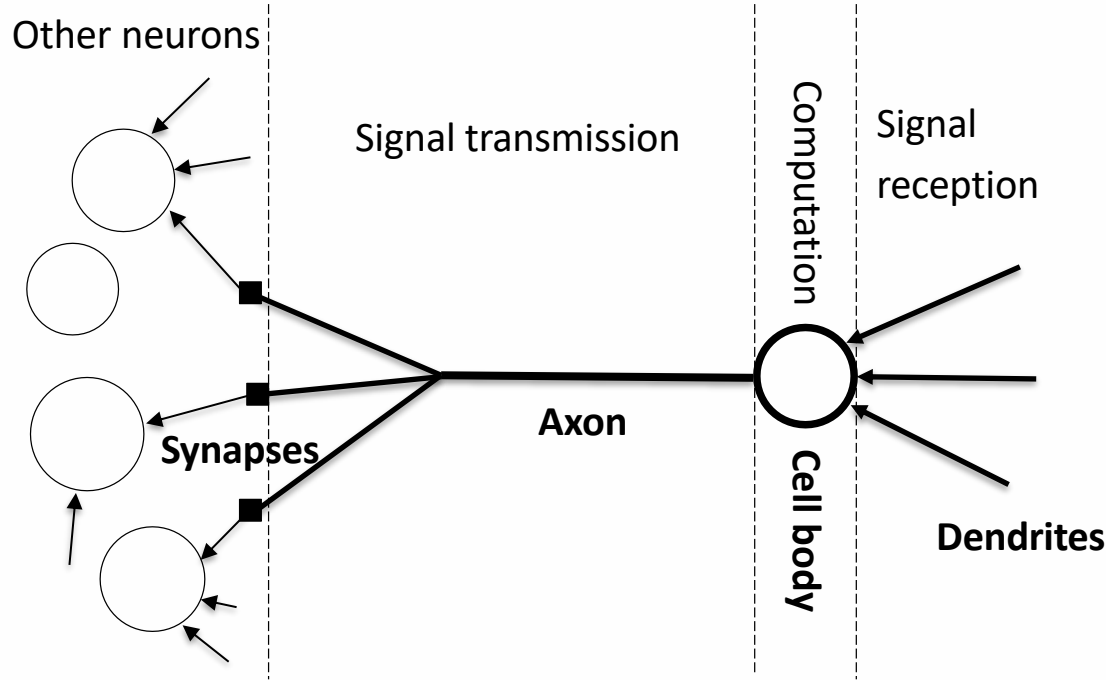
The number of neurons in the human brain is estimated to be around $10^{10} - 10^{11}$

Biological neural network



- Neurons receive information from other neurons through their dendrites
- Neurons process the information in their cell body (soma)
- Neurons send information through an axon
- The connection between the axon branches and other neurons' dendrites are called synapses

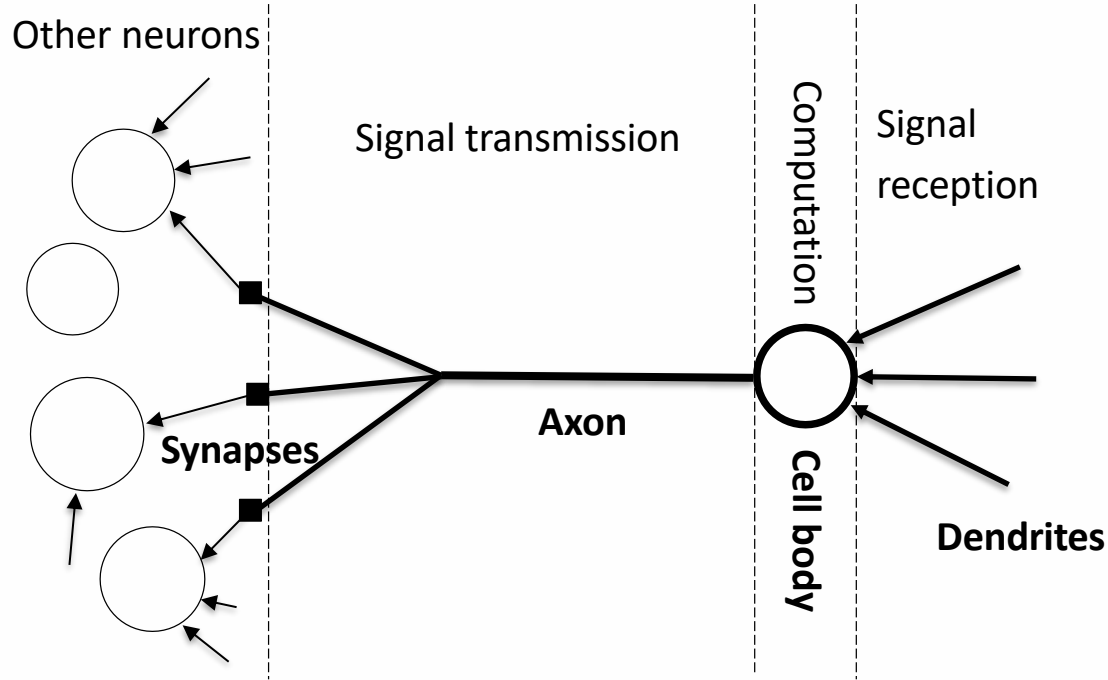
Biological neuron



An action potential is an electrical impulse that travels through the axon

- this is how neurons communicate
- neuron generates a spike in the electric potential of the axon
- an action potential is generated at neuron only if the pattern of spikes it receives from other neurons is above some threshold

Biological neuron



Neurons generate several spikes every seconds

- the neuron activity is characterized by its firing rate which is the frequency of the spikes
- neurons have a spontaneous firing rate, so they always fire a little bit, but they will increase the rate if receive the right stimulus

Firing rates of different input neurons are combined and influence the firing rate of other neurons

- depending on the dendrite and axon, a neuron can either work to increase (excite) or decrease (inhibit) the firing rate of another neuron

Artificial neuron

Neuron pre-activation (or input activation):

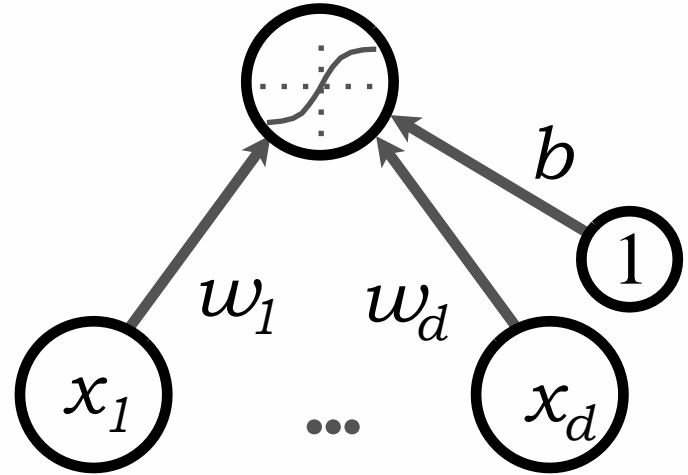
- $a(\mathbf{x}) = b + \sum_i w_i x_i = b + \mathbf{w}^T \mathbf{x}$

Neuron (output) activation:

- $h(\mathbf{x}) = g(a(\mathbf{x})) = g(b + \sum_i w_i x_i)$

where

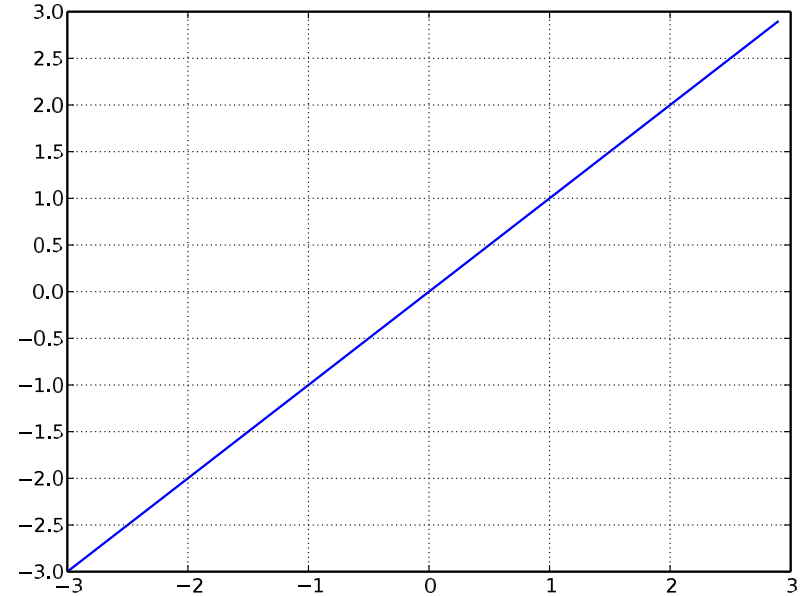
- \mathbf{w} are connection weights
- b is neuron bias
- $g(a(\mathbf{x}))$ is the activation function



Activation function

Linear activation function

- Performs no input squashing
- Has infinite values

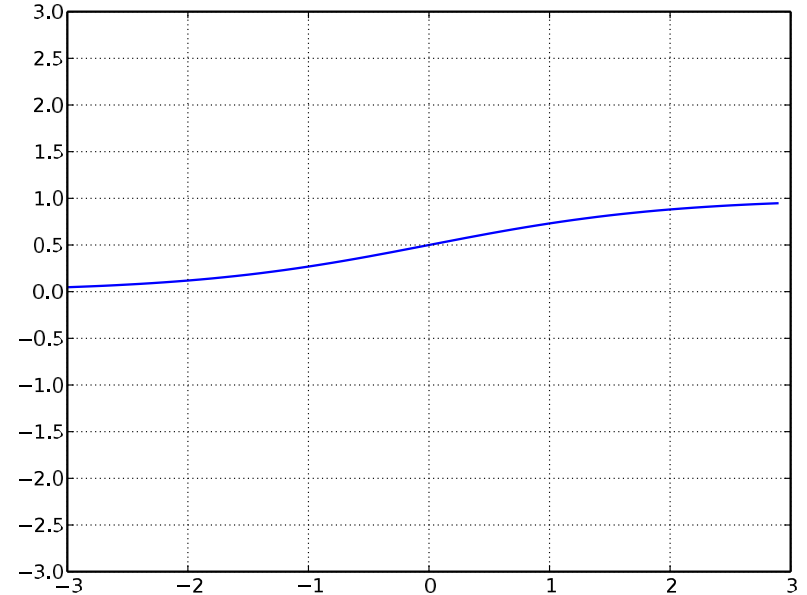


$$g(a) = a$$

Activation function

Sigmoid activation function

- Squashes the neuron pre-activation value to $[0, 1]$ range
- Has finite values
- Always positive
- Bounded
- Strictly increasing

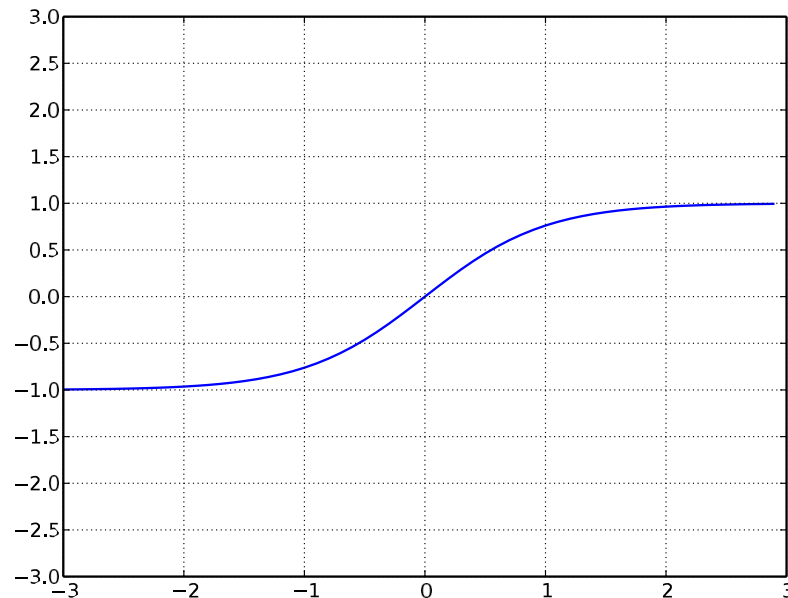


$$g(a) = \text{sigm}(a) = \frac{1}{1 + e^{-a}}$$

Activation function

Hyperbolic tangent activation function

- Squashes the neuron pre-activation value to $[-1, 1]$ range
- Has finite values
- Can be positive or negative
- Bounded
- Strictly increasing

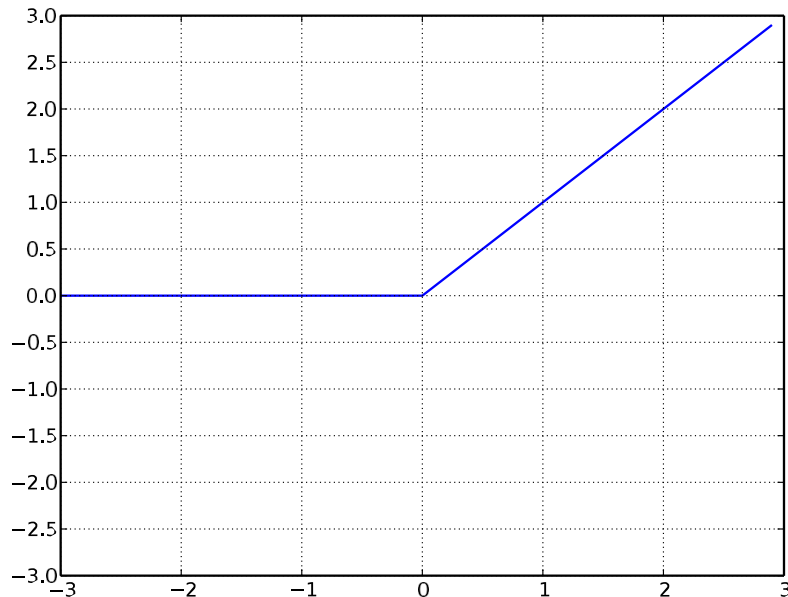


$$g(a) = \tanh(a) = \frac{e^a - e^{-a}}{e^a + e^{-a}} = \frac{e^{2a} - 1}{e^{2a} + 1}$$

Activation function

Rectified linear activation function

- Bounded below by 0
- Has infinite values (not upper bounded)
- Always non-negative
- Tends to give neurons with sparse activities



$$g(a) = \text{reclin}(a) = \max(0, a)$$

Capacity of a single neuron

Range is determined by the activation function $g(a)$

- $h(x) = g(a(x)) = g(b + \sum_i w_i x_i)$

Bias b only changes the position of the riff

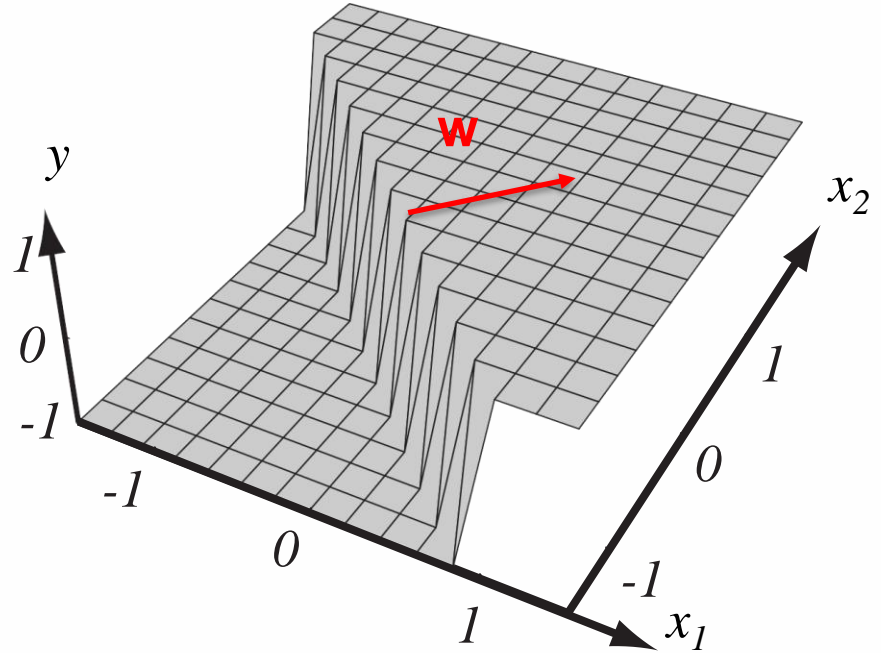
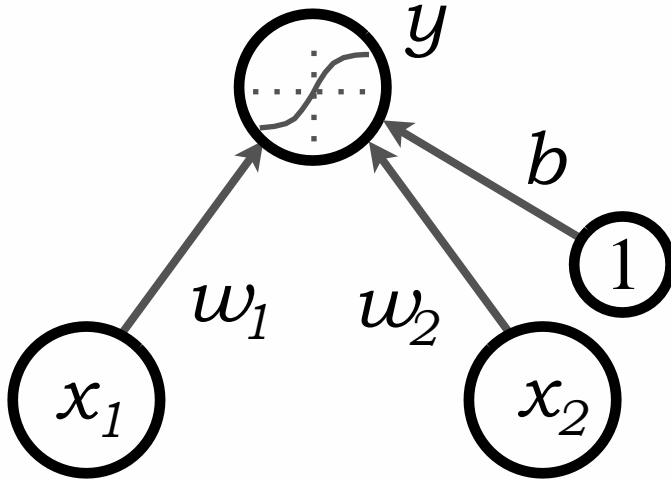


Image by Pascal Vincent

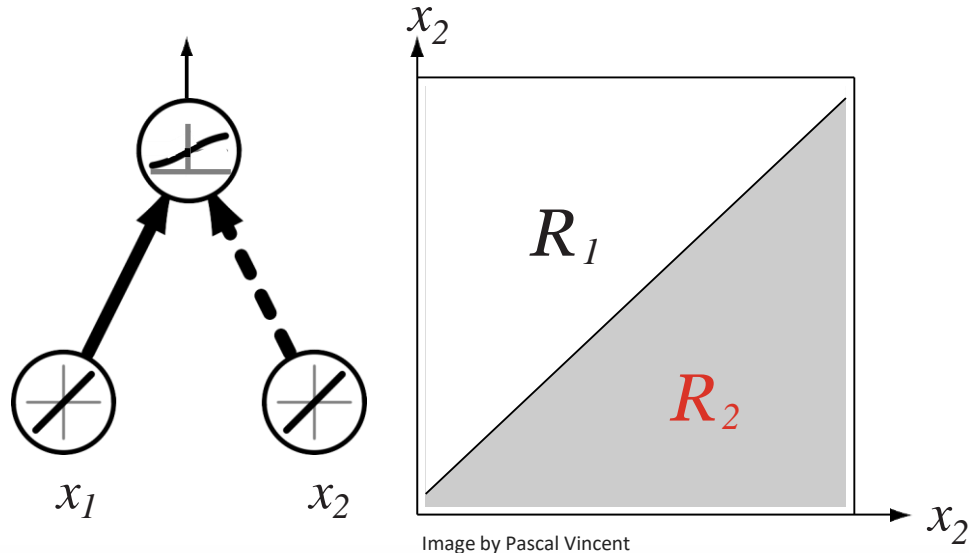
Capacity of a single neuron

Single neuron can do a binary classification

With sigmoid activation function it can be an estimation of probability

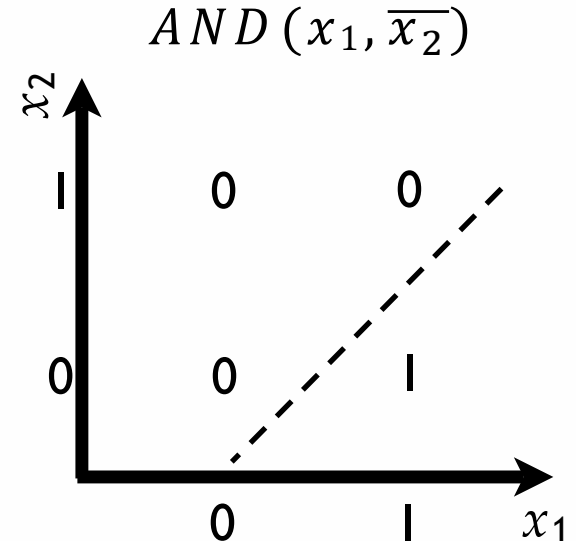
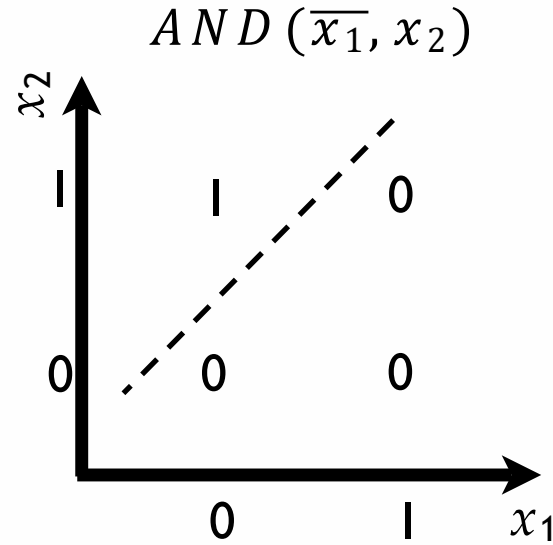
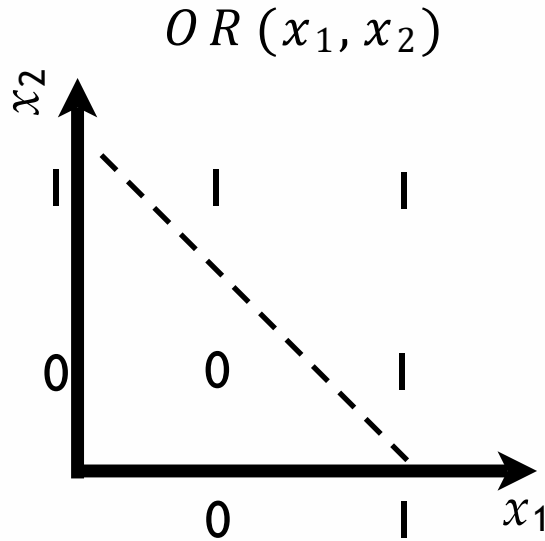
- Logical regression classifier
- $f(\mathbf{x}) = p(y = 1|\mathbf{x})$
- If greater than 0.5
predicts class 1,
- Otherwise,
predicts class 2

Decision boundary is linear



Capacity of a single neuron

Single neuron can solve linearly separable problems

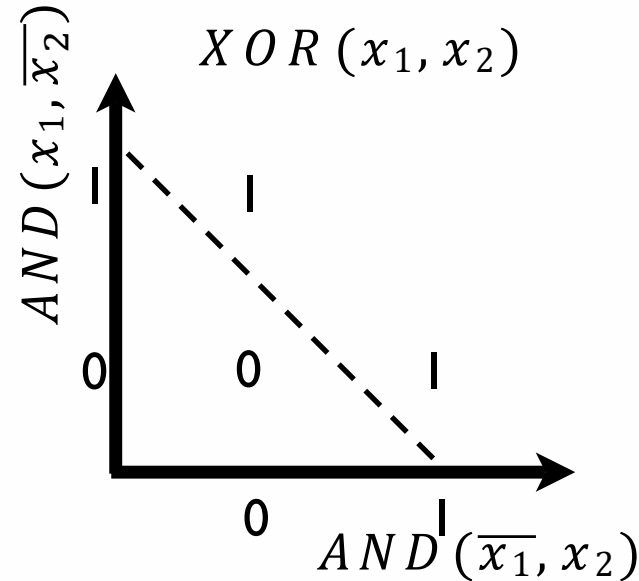
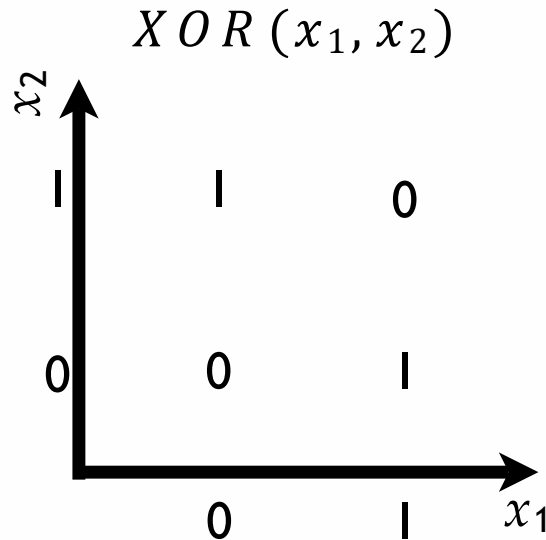


Capacity of a single neuron

Single neuron can't solve not linearly separable problems

Unless we transform an input into a better representation

- $XOR(x_1, x_2) = OR(AND(x_1, \overline{x_2}), AND(\overline{x_1}, x_2))$



Single hidden layer

Hidden layer pre-activation:

- $a(\mathbf{x})_i = b_i^{(1)} + \sum_j w_{i,j}^{(1)} x_j = b_i^{(1)} + \mathbf{w}_i^{(1)T} \mathbf{x}$
- $\mathbf{a}(\mathbf{x}) = \mathbf{b}^{(1)} + \mathbf{W}^{(1)} \mathbf{x}$

Hidden layer activation:

- $\mathbf{h}(\mathbf{x}) = \mathbf{g}(\mathbf{a}(\mathbf{x}))$

Output layer activation:

- $f(\mathbf{x}) = o \left(b^{(2)} + \mathbf{w}^{(1)T} \mathbf{h}(\mathbf{x}) \right)$

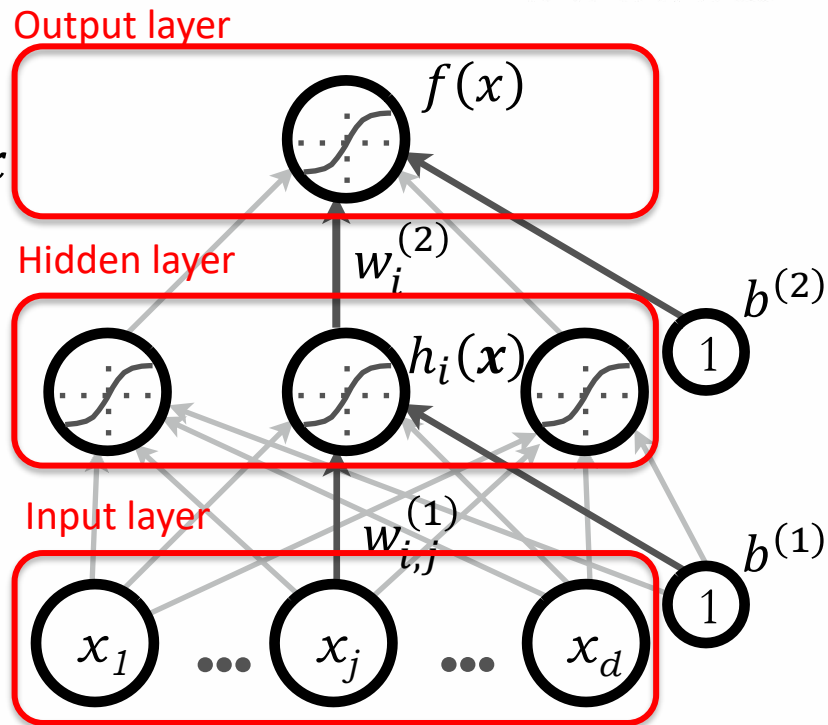


Image by Hugo Larochelle

Multilayer neural network

For a network with L hidden layers:

k -th hidden layer pre-activation:

- $\mathbf{a}^{(k)}(\mathbf{x}) = \mathbf{b}^{(k)} + \mathbf{W}^{(k)} \mathbf{h}^{(k-1)}(\mathbf{x})$

k -th hidden layer activation:

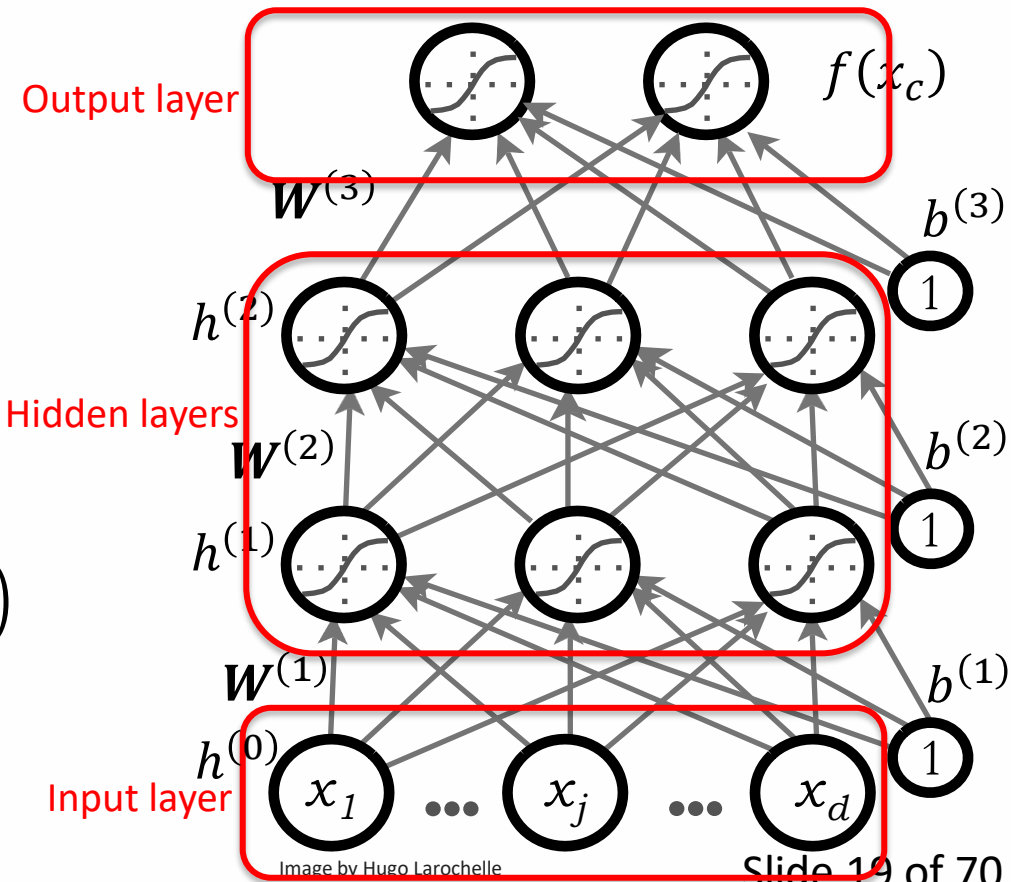
- $\mathbf{h}^{(k)}(\mathbf{x}) = \mathbf{g}(\mathbf{a}^{(k)}(\mathbf{x}))$

Output layer activation:

- $f(\mathbf{x}) = \mathbf{h}^{(L+1)}(\mathbf{x}) = \mathbf{o}(\mathbf{a}^{(L+1)}(\mathbf{x}))$

Input

- $\mathbf{h}^{(0)}(\mathbf{x}) = \mathbf{x}$



Multiple outputs

For multiple class classification we need multiple outputs:

- 1 output per class
- Each output estimates the probability of belonging to a class:
 - $f(\mathbf{x})_c = p(y = c|\mathbf{x})$
- Predicted class is the one with the highest probability

Softmax activation function:

- $\mathbf{o}(\mathbf{a}) = \text{softmax}(\mathbf{a}) = \left[\frac{e^{a_1}}{\sum_c e^{a_c}} \quad \cdots \quad \frac{e^{a_N}}{\sum_c e^{a_c}} \right]^T$
 - where N is the number of classes
- Strictly positive
- Sums to 1

Exponent, $\exp(a) = e^a$

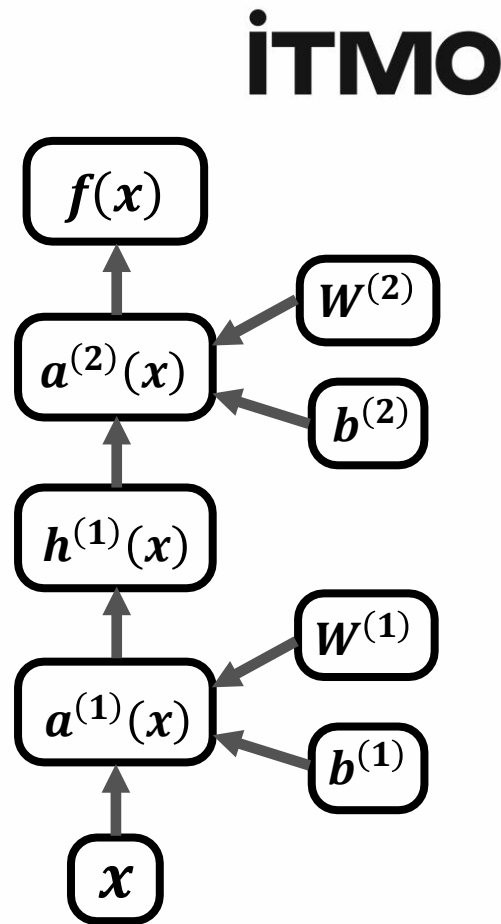
Forward propagation

Calculation of the neural network is called the **forward propagation**

It can be represented as an acyclic flow graph

Can be easily implemented in modular way:

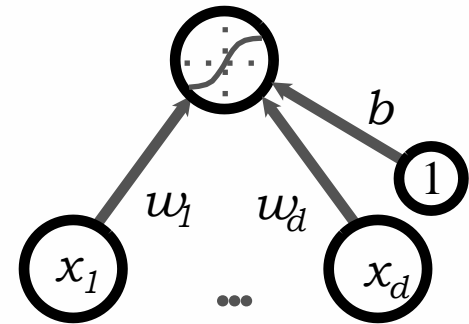
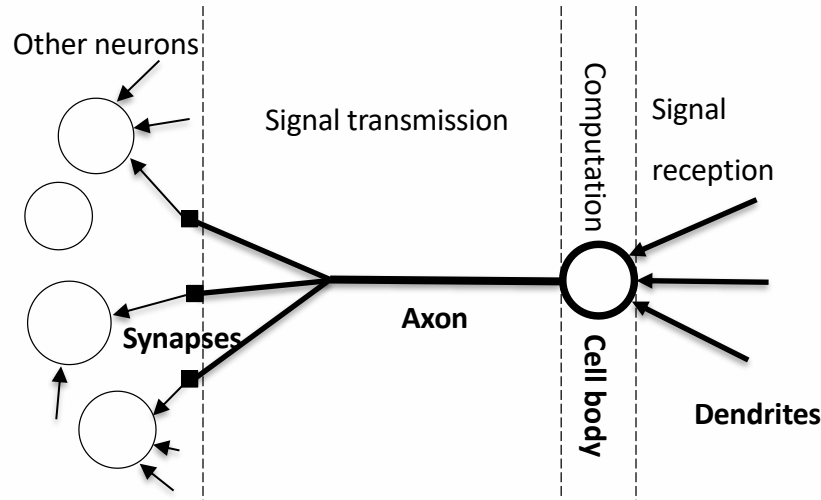
- Each block is a function with arguments of its children
- Calling functions in a right order implements a forward propagation



Biological vs artificial neuron

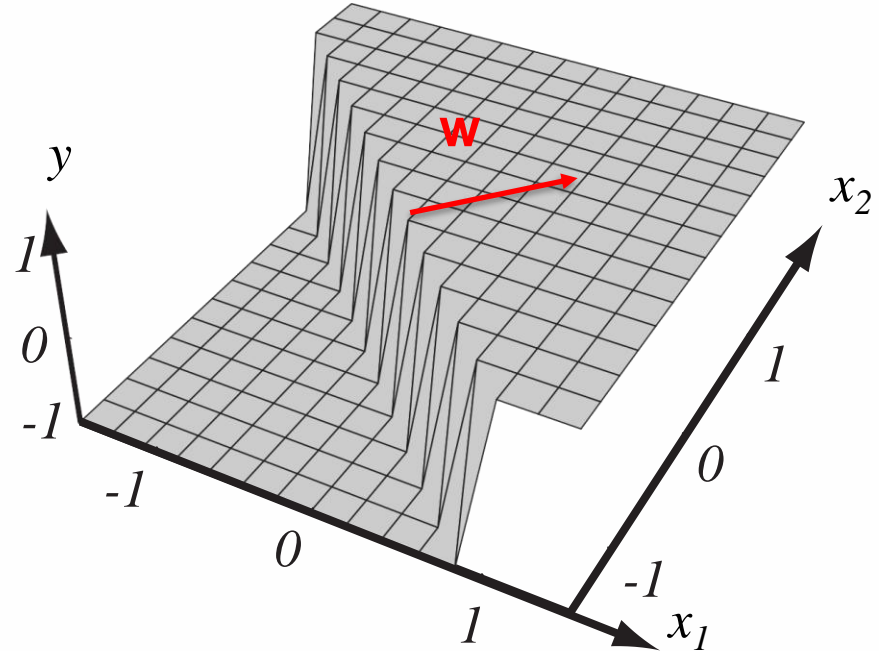
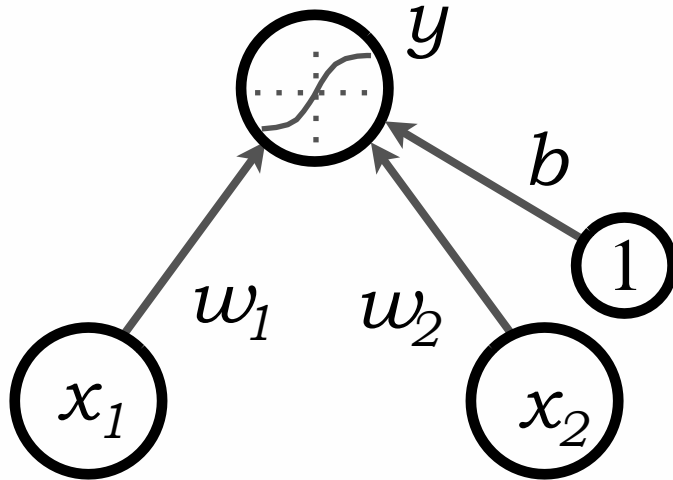
The artificial neuron approximates the biological one:

- the activation corresponds to a “sort of ” firing rate
- the weights between neurons model whether neurons excite or inhibit each other
- the activation function and bias model the thresholded behavior of action potentials



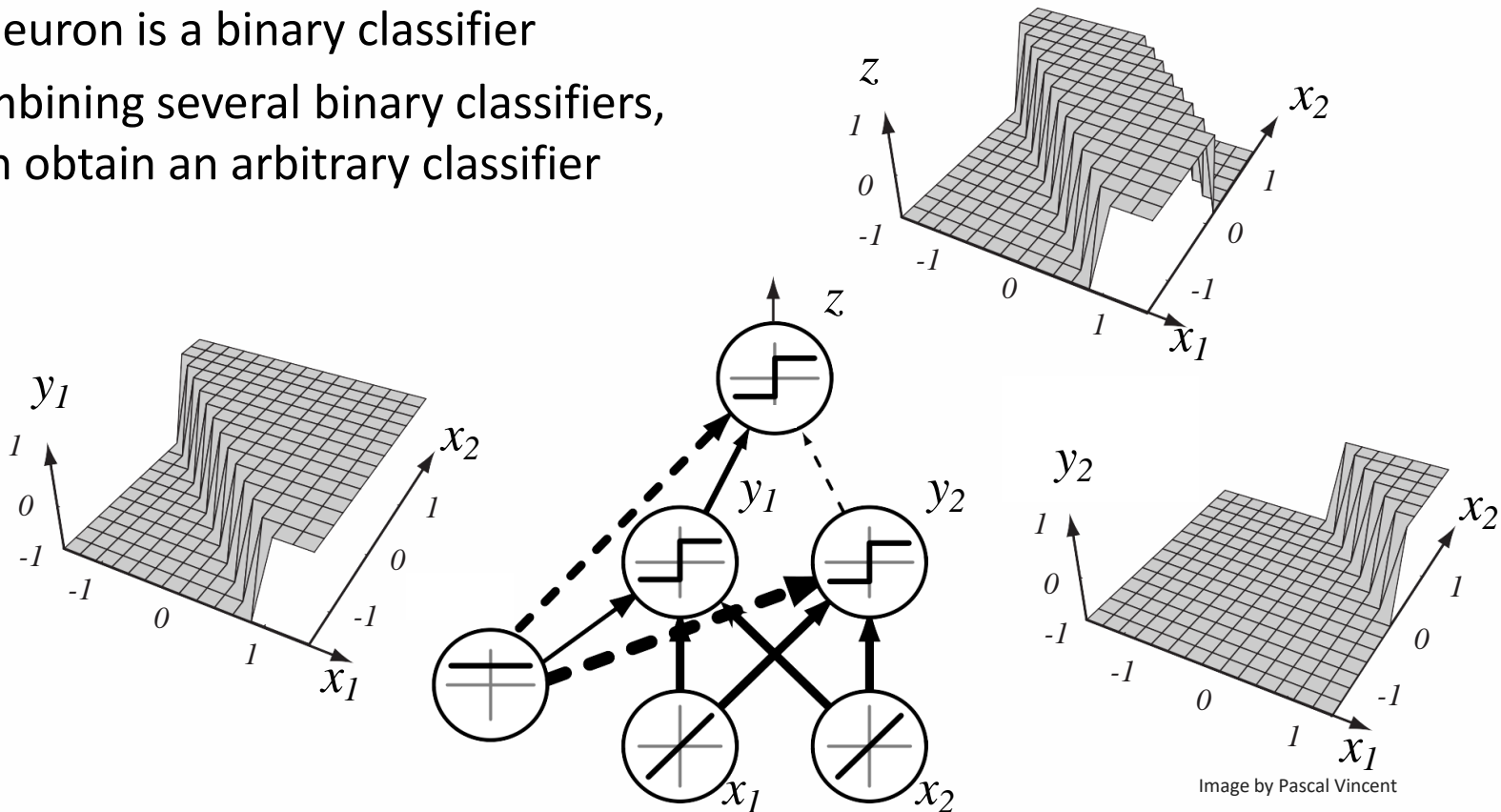
Capacity of a single neuron

- Single neuron can do a binary classification
- With sigmoid activation function it can be an estimation of probability



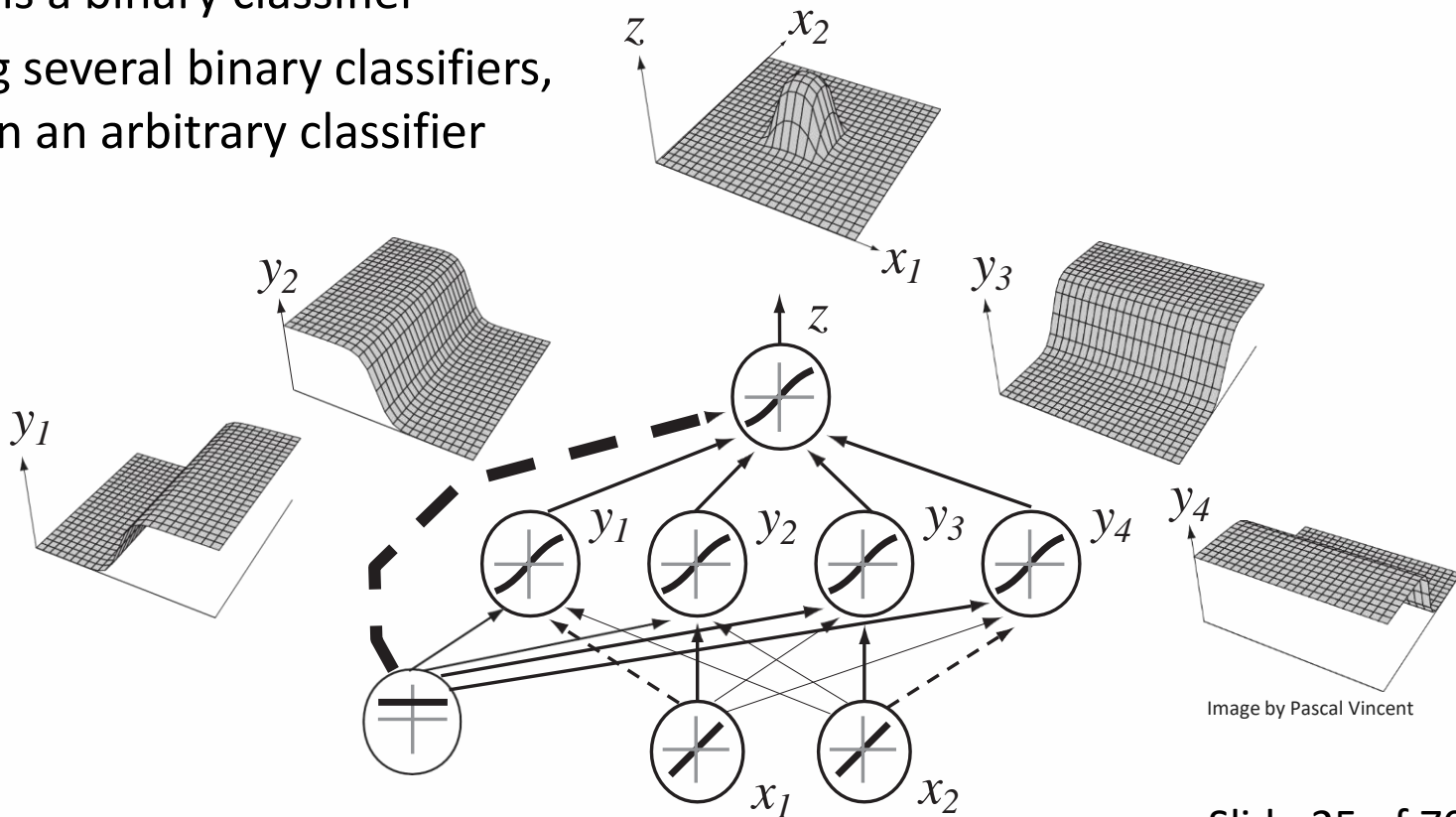
Capacity of a single hidden layer neural network

- Each neuron is a binary classifier
- By combining several binary classifiers, we can obtain an arbitrary classifier



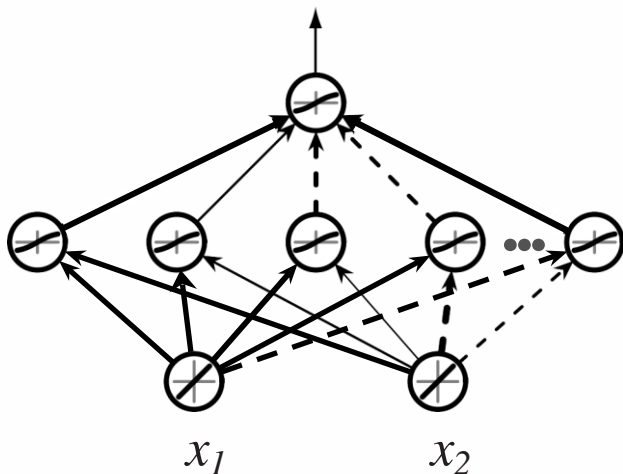
Capacity of a single hidden layer neural network

- Each neuron is a binary classifier
- By combining several binary classifiers, we can obtain an arbitrary classifier



Capacity of a single hidden layer neural network

- Each neuron is a binary classifier
- By combining several binary classifiers, we can obtain an arbitrary classifier



iTMO

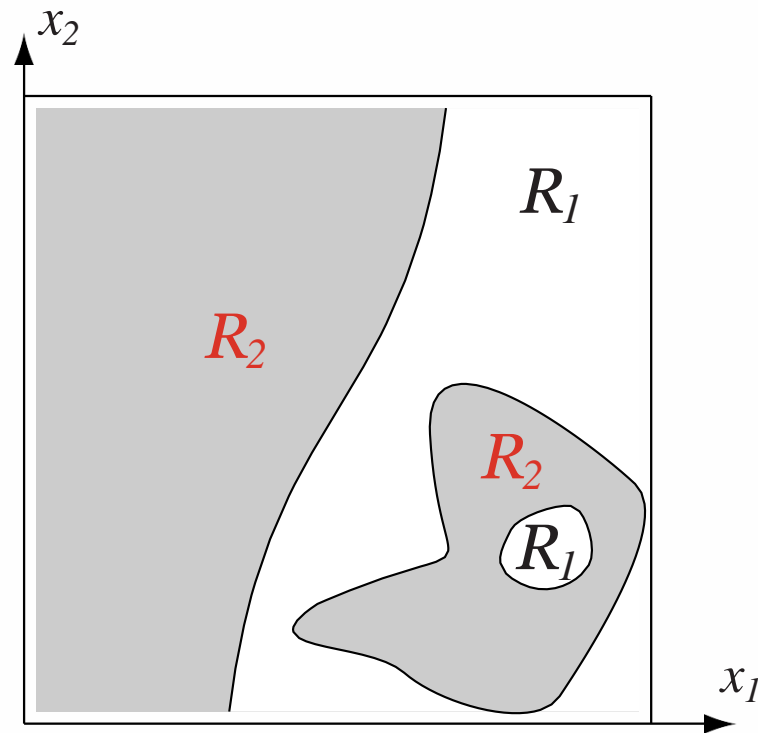


Image by Pascal Vincent

Capacity of a single hidden layer neural network

- Each neuron is a binary classifier
- By combining several binary classifiers, we can obtain an arbitrary classifier
- Universal approximation theorem (Hornik, 1991):
 - “A single hidden layer neural network with a linear output unit can approximate any continuous function arbitrarily well, given enough hidden units”
 - This can be applied for sigmoid, hyperbolic tangent and many other hidden layer activation functions
- It doesn't mean there is a learning algorithm that can find the necessary neural network parameter values

Deep learning

Use the supervised machine learning for empirical risk minimization

- $\arg \min_{\theta} \frac{1}{T} \sum_t l(f(\mathbf{x}^{(t)}, \theta), y^{(t)}) + \lambda \Omega(\theta)$
 - $l(f(\mathbf{x}^{(t)}, \theta), y^{(t)})$ is a loss function
 - $\Omega(\theta)$ is a regularizer, which penalizes some certain values of θ ; $\theta = \{W_{i,j}^k, b^k\}$
 - $\{(\mathbf{x}^{(t)}, y^{(t)})\}$ is the dataset
- Learning is the task of optimization of the empirical risk
 - We would like to optimize a classification error, but it's not smooth
 - Loss function is an upper bound of what we should optimize

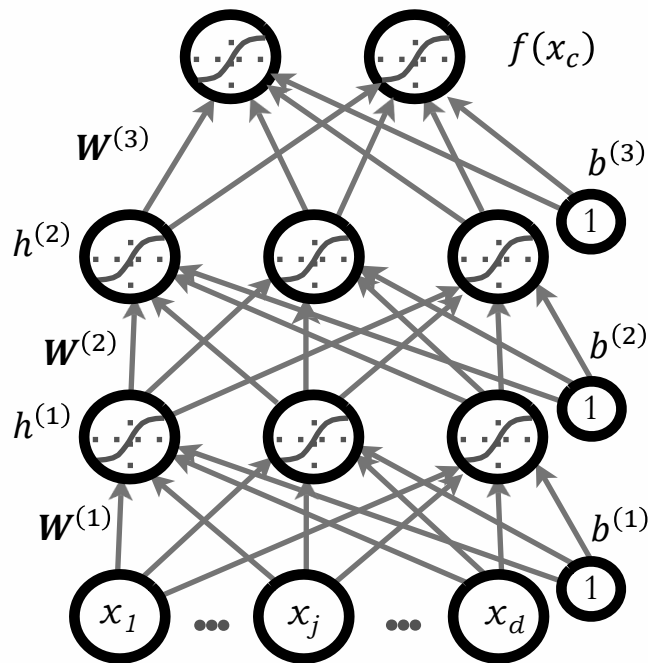


Image by Hugo Larochelle

Stochastic gradient descent (SGD) method

- Initialize θ_0 values: $\theta_0 \equiv \{W^{(1)}, b^{(1)}, \dots, W^{(L+1)}, b^{(L+1)}\}$
- Do N iterations:
 - for each training example $(x^{(t)}, y^{(t)})$
 - $\Delta_e = -\nabla_{\theta} l(f(x^{(t)}, \theta_e), y^{(t)}) - \lambda \nabla_{\theta} \Omega(\theta_e)$
 - $\theta_{e+1} \leftarrow \theta_e + \alpha_e \Delta_e$

} Training epoch

To implement the SGD method, we need:

- a loss function $l(f(x^{(t)}, \theta), y^{(t)})$
- a method to compute the parameters gradient $\nabla_{\theta} l(f(x^{(t)}, \theta), y^{(t)})$
- a regularizer function $\Omega(\theta)$ and a method to compute its gradient $\nabla_{\theta} \Omega(\theta)$
- a method to initialize parameter values θ

Loss function

To implement the SGD method, we need:

- a loss function $l(f(\mathbf{x}^{(t)}, \boldsymbol{\theta}), y^{(t)})$
- As we already learned:
 - Neural network output estimates the probability $f_c(\mathbf{x}) = p(y = c|\mathbf{x})$
 - We want to maximize the probabilities $y^{(t)}$ for a training set items $x^{(t)}$
- We will minimize the log-likelihood
 - $l(f(\mathbf{x}), y) = -\sum_c 1_{(y=c)} \ln(f(\mathbf{x})_c) = -\ln(f(\mathbf{x})_y)$
 - Natural logarithm is used for numerical stability and math simplicity
 - Sometimes is referred as a cross-entropy

Loss function derivative

To implement the SGD method, we need:

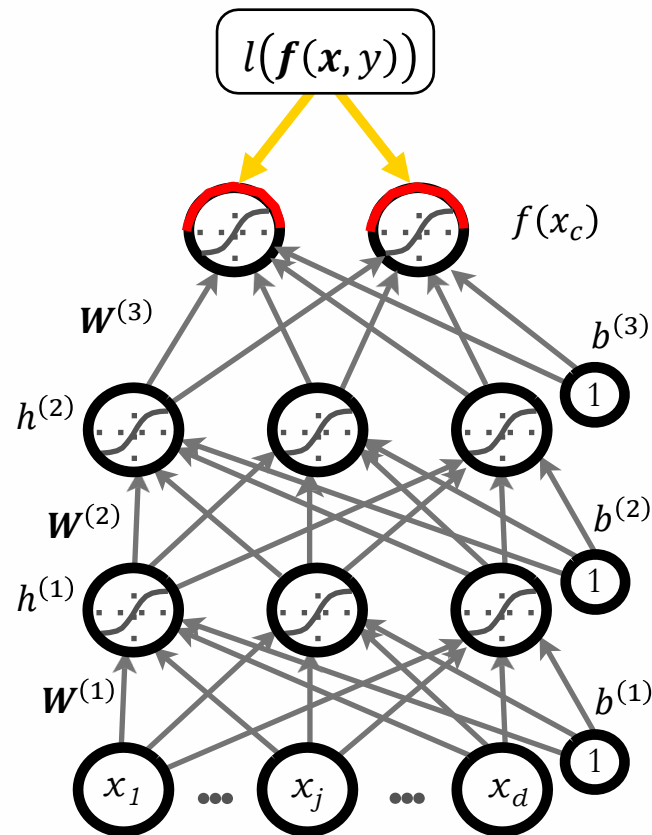
- a method to compute the parameters gradient

Partial derivative of output activation:

$$\bullet \frac{\partial}{\partial f(x)_c} (-\ln(f(x)_y)) = \frac{-1_{(y=c)}}{f(x)_y}$$

Gradient of output activation:

$$\begin{aligned} \bullet \nabla_{f(x)} (-\ln(f(x)_y)) &= \\ &= \frac{-1}{f(x)_y} \begin{bmatrix} 1_{(y=1)} \\ \dots \\ 1_{(y=N)} \end{bmatrix} = \frac{-\mathbf{e}(y)}{f(x)_y} \end{aligned}$$



Loss function derivative

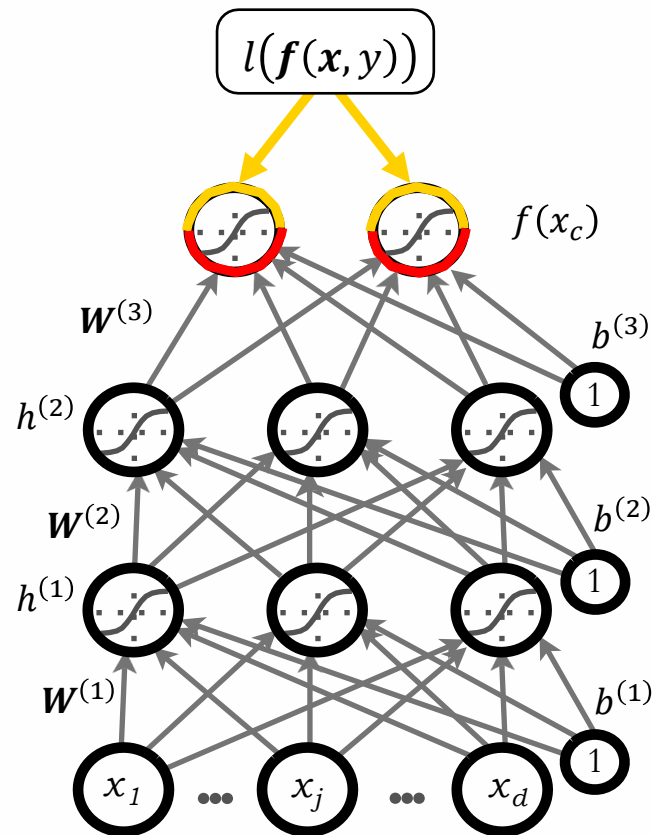
To implement the SGD method, we need:

- a method to compute the parameters gradient

Partial derivative of output pre-activation:

$$\begin{aligned} & \bullet \frac{\partial}{\partial a^{(L+1)}(x)_c} (-\ln(f(x)_y)) = \\ &= \frac{\partial(-\ln(f(x)_y))}{\partial f(x)_y} \frac{\partial f(x)_y}{\partial a^{(L+1)}(x)_c} = \\ &= \frac{-1}{f(x)_y} \frac{\partial f(x)_y}{\partial a^{(L+1)}(x)_c} = \\ &= \frac{-1}{f(x)_y} \frac{\partial}{\partial a^{(L+1)}(x)_c} \text{softmax} \left(a^{(L+1)}(x) \right)_y = \end{aligned}$$

Chain rule: $\frac{\partial y}{\partial x} = \frac{\partial y}{\partial t} \frac{\partial t}{\partial x}$



Loss function derivative

$$= \frac{-1}{f(\mathbf{x})_y} \frac{\partial}{\partial a^{(L+1)}(\mathbf{x})_c} \frac{e^{(a^{(L+1)}(\mathbf{x})_y)}}{\sum_{c'} e^{(a^{(L+1)}(\mathbf{x})_{c'})}} =$$

$$\text{Derivative of ratio: } \frac{\partial}{\partial x} \left(\frac{g(x)}{h(x)} \right) = \frac{\partial g(x)}{\partial x} \frac{1}{h(x)} - \frac{g(x)}{h(x)^2} \frac{\partial h(x)}{\partial x}$$

Loss function derivative

$$\begin{aligned} &= \frac{-1}{f(\mathbf{x})_y} \frac{\partial}{\partial a^{(L+1)}(\mathbf{x})_c} \frac{e^{(a^{(L+1)}(\mathbf{x})_y)}}{\sum_{c'} e^{(a^{(L+1)}(\mathbf{x})_{c'})}} = \\ &= \frac{-1}{f(\mathbf{x})_y} \left(\frac{\frac{\partial}{\partial a^{(L+1)}(\mathbf{x})_c} e^{(a^{(L+1)}(\mathbf{x})_y)}}{\sum_{c'} e^{(a^{(L+1)}(\mathbf{x})_{c'})}} - \left(\frac{e^{(a^{(L+1)}(\mathbf{x})_y)} \frac{\partial}{\partial a^{(L+1)}(\mathbf{x})_c} \sum_{c'} e^{(a^{(L+1)}(\mathbf{x})_{c'})}}{\left(\sum_{c'} e^{(a^{(L+1)}(\mathbf{x})_{c'})} \right)^2} \right) \right) = \end{aligned}$$

Derivative of ratio: $\frac{\partial}{\partial x} \left(\frac{g(x)}{h(x)} \right) = \frac{\partial g(x)}{\partial x} \frac{1}{h(x)} - \frac{g(x)}{h(x)^2} \frac{\partial h(x)}{\partial x}$
--

Loss function derivative

$$\begin{aligned} &= \frac{-1}{f(\mathbf{x})_y} \frac{\partial}{\partial a^{(L+1)}(\mathbf{x})_c} \frac{e^{(a^{(L+1)}(\mathbf{x})_y)}}{\sum_{c'} e^{(a^{(L+1)}(\mathbf{x})_{c'})}} = \\ &= \frac{-1}{f(\mathbf{x})_y} \left(\frac{\frac{\partial}{\partial a^{(L+1)}(\mathbf{x})_c} e^{(a^{(L+1)}(\mathbf{x})_y)}}{\sum_{c'} e^{(a^{(L+1)}(\mathbf{x})_{c'})}} - \left(\frac{e^{(a^{(L+1)}(\mathbf{x})_y)} \frac{\partial}{\partial a^{(L+1)}(\mathbf{x})_c} \sum_{c'} e^{(a^{(L+1)}(\mathbf{x})_{c'})}}{\left(\sum_{c'} e^{(a^{(L+1)}(\mathbf{x})_{c'})} \right)^2} \right) \right) = \\ &= \frac{-1}{f(\mathbf{x})_y} \left(\frac{1_{(y=c)} e^{(a^{(L+1)}(\mathbf{x})_y)}}{\sum_{c'} e^{(a^{(L+1)}(\mathbf{x})_{c'})}} - \frac{e^{(a^{(L+1)}(\mathbf{x})_y)}}{\sum_{c'} e^{(a^{(L+1)}(\mathbf{x})_{c'})}} \frac{e^{(a^{(L+1)}(\mathbf{x})_c)}}{\sum_{c'} e^{(a^{(L+1)}(\mathbf{x})_{c'})}} \right) = \end{aligned}$$

Derivative of ratio: $\frac{\partial}{\partial x} \left(\frac{g(x)}{h(x)} \right) = \frac{\partial g(x)}{\partial x} \frac{1}{h(x)} - \frac{g(x)}{h(x)^2} \frac{\partial h(x)}{\partial x}$

Loss function derivative

$$\begin{aligned}
 &= \frac{-1}{f(\mathbf{x})_y} \frac{\partial}{\partial a^{(L+1)}(\mathbf{x})_c} \frac{e^{(a^{(L+1)}(\mathbf{x})_y)}}{\sum_{c'} e^{(a^{(L+1)}(\mathbf{x})_{c'})}} = \\
 &= \frac{-1}{f(\mathbf{x})_y} \left(\frac{\frac{\partial}{\partial a^{(L+1)}(\mathbf{x})_c} e^{(a^{(L+1)}(\mathbf{x})_y)}}{\sum_{c'} e^{(a^{(L+1)}(\mathbf{x})_{c'})}} - \left(\frac{e^{(a^{(L+1)}(\mathbf{x})_y)} \frac{\partial}{\partial a^{(L+1)}(\mathbf{x})_c} \sum_{c'} e^{(a^{(L+1)}(\mathbf{x})_{c'})}}{\left(\sum_{c'} e^{(a^{(L+1)}(\mathbf{x})_{c'})} \right)^2} \right) \right) = \\
 &= \frac{-1}{f(\mathbf{x})_y} \left(\frac{1_{(y=c)} e^{(a^{(L+1)}(\mathbf{x})_y)}}{\sum_{c'} e^{(a^{(L+1)}(\mathbf{x})_{c'})}} - \frac{e^{(a^{(L+1)}(\mathbf{x})_y)}}{\sum_{c'} e^{(a^{(L+1)}(\mathbf{x})_{c'})}} \frac{e^{(a^{(L+1)}(\mathbf{x})_c)}}{\sum_{c'} e^{(a^{(L+1)}(\mathbf{x})_{c'})}} \right) = \\
 &= \frac{-1}{f(\mathbf{x})_y} \left(1_{(y=c)} \text{softmax} \left(a^{(L+1)}(\mathbf{x}) \right)_y - \text{softmax} \left(a^{(L+1)}(\mathbf{x}) \right)_y \text{softmax} \left(a^{(L+1)}(\mathbf{x}) \right)_c \right) = \\
 &= \frac{-1}{f(\mathbf{x})_y} (1_{(y=c)} f(\mathbf{x})_y - f(\mathbf{x})_y f(\mathbf{x})_c) =
 \end{aligned}$$

Derivative of ratio: $\frac{\partial}{\partial x} \left(\frac{g(x)}{h(x)} \right) = \frac{\partial g(x)}{\partial x} \frac{1}{h(x)} - \frac{g(x)}{h(x)^2} \frac{\partial h(x)}{\partial x}$

Loss function derivative

$$\begin{aligned} &= \frac{-1}{f(\mathbf{x})_y} \frac{\partial}{\partial a^{(L+1)}(\mathbf{x})_c} \frac{e^{(a^{(L+1)}(\mathbf{x})_y)}}{\sum_{c'} e^{(a^{(L+1)}(\mathbf{x})_{c'})}} = \\ &= \frac{-1}{f(\mathbf{x})_y} \left(\frac{\frac{\partial}{\partial a^{(L+1)}(\mathbf{x})_c} e^{(a^{(L+1)}(\mathbf{x})_y)}}{\sum_{c'} e^{(a^{(L+1)}(\mathbf{x})_{c'})}} - \left(\frac{e^{(a^{(L+1)}(\mathbf{x})_y)} \frac{\partial}{\partial a^{(L+1)}(\mathbf{x})_c} \sum_{c'} e^{(a^{(L+1)}(\mathbf{x})_{c'})}}{\left(\sum_{c'} e^{(a^{(L+1)}(\mathbf{x})_{c'})} \right)^2} \right) \right) = \\ &= \frac{-1}{f(\mathbf{x})_y} \left(\frac{1_{(y=c)} e^{(a^{(L+1)}(\mathbf{x})_y)}}{\sum_{c'} e^{(a^{(L+1)}(\mathbf{x})_{c'})}} - \frac{e^{(a^{(L+1)}(\mathbf{x})_y)}}{\sum_{c'} e^{(a^{(L+1)}(\mathbf{x})_{c'})}} \frac{e^{(a^{(L+1)}(\mathbf{x})_c)}}{\sum_{c'} e^{(a^{(L+1)}(\mathbf{x})_{c'})}} \right) = \\ &= \frac{-1}{f(\mathbf{x})_y} \left(1_{(y=c)} \text{softmax} \left(a^{(L+1)}(\mathbf{x}) \right)_y - \text{softmax} \left(a^{(L+1)}(\mathbf{x}) \right)_y \text{softmax} \left(a^{(L+1)}(\mathbf{x}) \right)_c \right) = \\ &= \frac{-1}{f(\mathbf{x})_y} (1_{(y=c)} f(\mathbf{x})_y - f(\mathbf{x})_y f(\mathbf{x})_c) = -(1_{(y=c)} - f(\mathbf{x})_c) \end{aligned}$$

Derivative of ratio: $\frac{\partial}{\partial x} \left(\frac{g(x)}{h(x)} \right) = \frac{\partial g(x)}{\partial x} \frac{1}{h(x)} - \frac{g(x)}{h(x)^2} \frac{\partial h(x)}{\partial x}$
--

Loss function derivative

To implement the SGD method, we need:

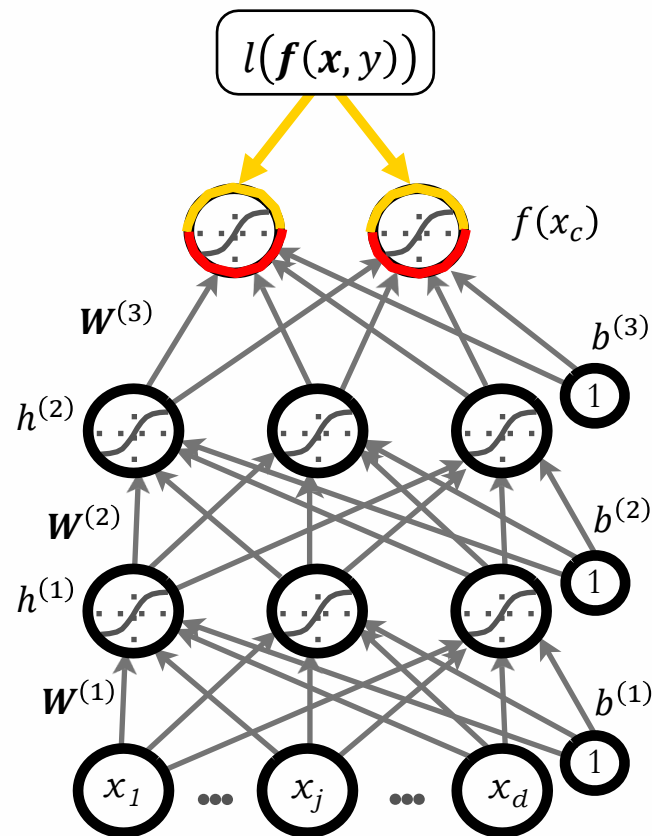
- a method to compute the parameters gradient

Partial derivative of output pre-activation:

$$\begin{aligned} \bullet \quad \frac{\partial}{\partial a^{(L+1)}(\mathbf{x})_c} (-\ln(f(\mathbf{x})_y)) &= \\ &= -(1_{(y=c)} - f(\mathbf{x})_c) \end{aligned}$$

Gradient of output pre-activation:

$$\begin{aligned} \bullet \quad \nabla_{a^{(L+1)}(\mathbf{x})} (-\ln(f(\mathbf{x})_y)) &= \\ &= -(\mathbf{e}(y) - \mathbf{f}(\mathbf{x})) \end{aligned}$$



Loss function derivative

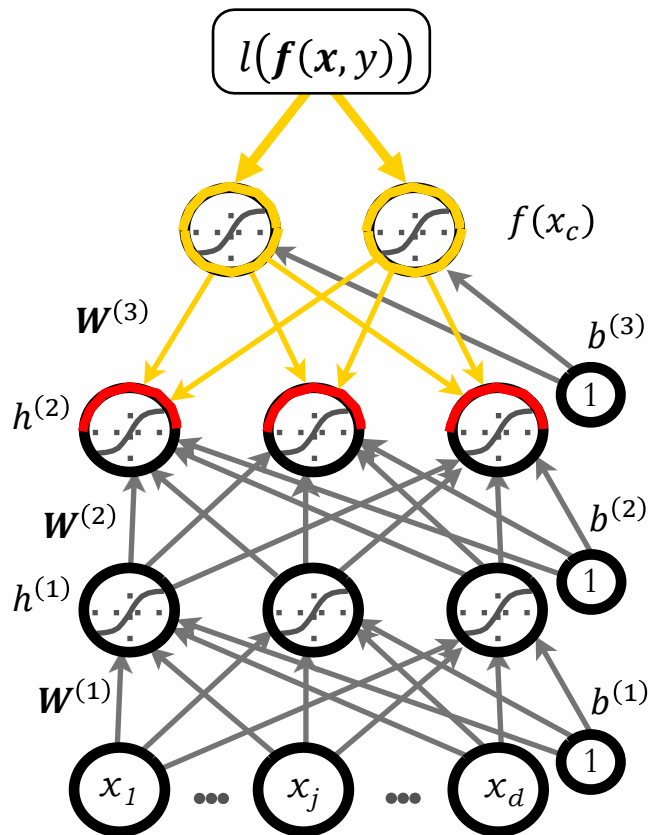
To implement the SGD method, we need:

- a method to compute the parameters gradient

Partial derivative of hidden layer activation:

- $\frac{\partial}{\partial h^{(k)}(\mathbf{x})_j} (-\ln(f(\mathbf{x})_y))$

$$a^{(k+1)}(\mathbf{x})_i = b_i^{(k+1)} + \sum_j w_{i,j}^{(k+1)} h^{(k)}(\mathbf{x})_j$$



Loss function derivative

To implement the SGD method, we need:

- a method to compute the parameters gradient

Partial derivative of hidden layer activation:

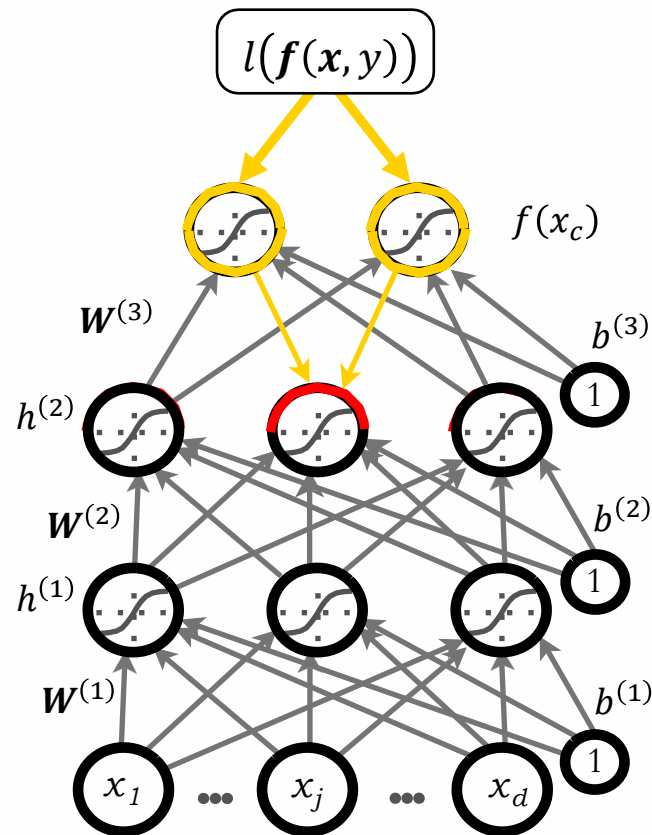
$$\bullet \frac{\partial}{\partial h^{(k)}(x)_j} (-\ln(f(x)_y))$$

Multivariable chain rule:

$$\begin{aligned} \bullet \frac{\partial p(a)}{\partial a} &= \frac{\partial p(q_1(a), \dots, q_n(a))}{\partial a} = \\ &= \sum_i \frac{\partial p(a)}{\partial q_i(a)} \frac{\partial q_i(a)}{\partial a} \end{aligned}$$

where a is a unit in layer

- $q_i(a)$ is a pre-activation in a layer above
- $p(a)$ is a loss function



Loss function derivative

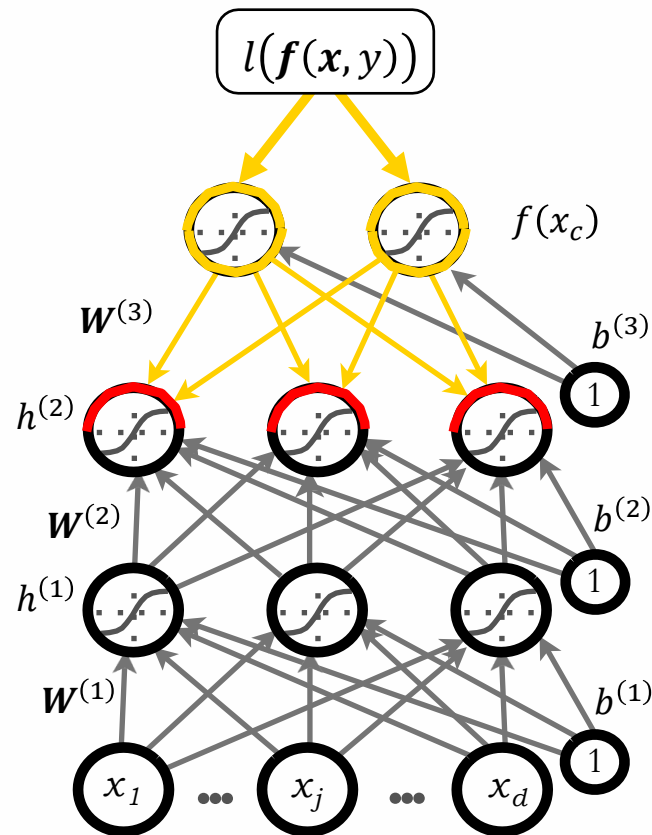
To implement the SGD method, we need:

- a method to compute the parameters gradient

Partial derivative of hidden layer activation:

$$\begin{aligned} & \bullet \frac{\partial}{\partial h^{(k)}(\mathbf{x})_j} (-\ln(f(\mathbf{x})_y)) = \\ &= \sum_i \frac{\partial(-\ln(f(\mathbf{x})_y))}{\partial a^{(k+1)}(\mathbf{x})_i} \frac{\partial a^{(k+1)}(\mathbf{x})_i}{\partial h^{(k)}(\mathbf{x})_j} = \\ &= \sum_i \frac{\partial(-\ln(f(\mathbf{x})_y))}{\partial a^{(k+1)}(\mathbf{x})_i} \mathbf{W}_{i,j}^{(k+1)} = \\ &= \mathbf{W}_{:,j}^{(k+1)T} \left(\nabla_{a^{(k+1)}(\mathbf{x})} (-\ln(f(\mathbf{x})_y)) \right) \end{aligned}$$

$$a^{(k+1)}(\mathbf{x})_i = b_i^{(k+1)} + \sum_j \mathbf{W}_{i,j}^{(k+1)} h^{(k)}(\mathbf{x})_j$$



Loss function derivative

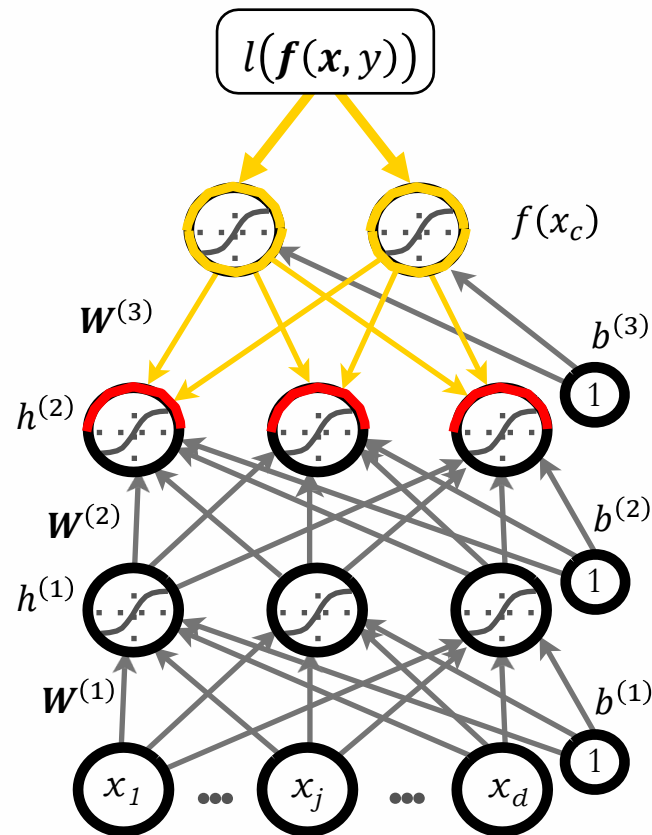
To implement the SGD method, we need:

- a method to compute parameters gradient

Gradient of hidden layer activation:

$$\begin{aligned} \bullet \quad \nabla_{\mathbf{h}^{(k)}(\mathbf{x})}(-\ln(f(\mathbf{x})_y)) &= \\ &= \mathbf{W}^{(k+1)T} \left(\nabla_{\mathbf{a}^{(k+1)}(\mathbf{x})}(-\ln(f(\mathbf{x})_y)) \right) \end{aligned}$$

$$a^{(k+1)}(\mathbf{x})_i = b_i^{(k+1)} + \sum_j W_{i,j}^{(k+1)} h^{(k)}(\mathbf{x})_j$$



Loss function derivative

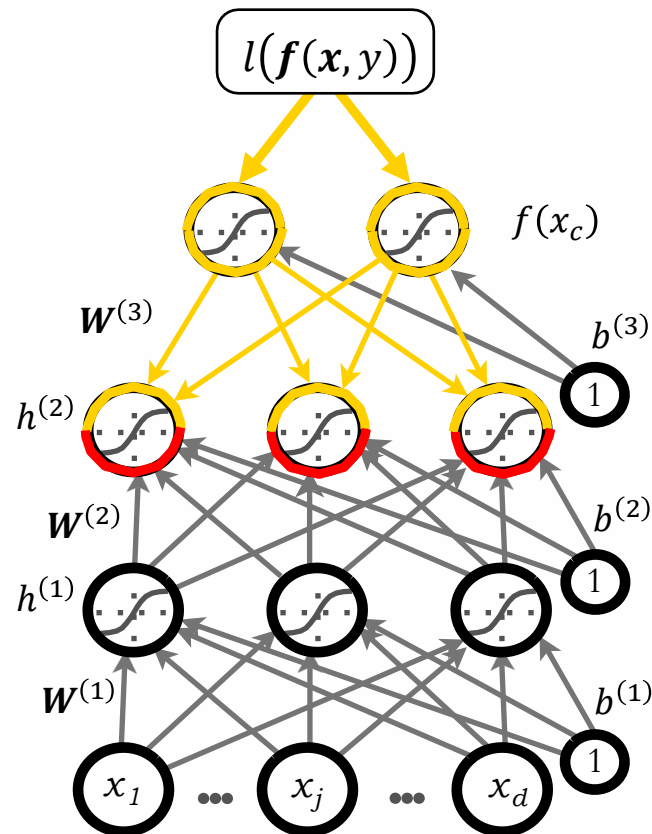
To implement the SGD method, we need:

- a method to compute the parameters gradient

Partial derivative of hidden layer pre-activation:

$$\begin{aligned} \bullet \quad & \frac{\partial}{\partial a^{(k)}(\mathbf{x})_j} (-\ln(f(\mathbf{x})_y)) = \\ &= \frac{\partial(-\ln(f(\mathbf{x})_y))}{\partial h^{(k)}(\mathbf{x})_j} \frac{\partial h^{(k)}(\mathbf{x})_j}{\partial a^{(k)}(\mathbf{x})_j} = \\ &= \frac{\partial(-\ln(f(\mathbf{x})_y))}{\partial h^{(k)}(\mathbf{x})_j} g'(a^{(k)}(\mathbf{x})_j) \end{aligned}$$

$$h^{(k)}(\mathbf{x})_j = g(a^{(k)}(\mathbf{x})_j)$$



Loss function derivative

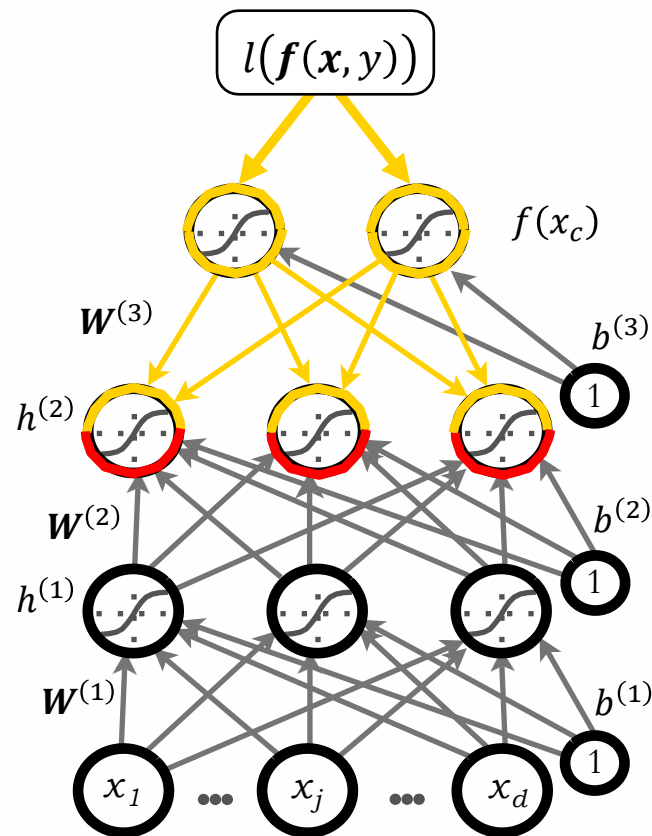
To implement the SGD method, we need:

- a method to compute the parameters gradient

Partial gradient of hidden layer pre-activation:

$$\begin{aligned} & \bullet \nabla_{a^{(k)}(x)}(-\ln(f(x)_y)) = \\ & = \left(\nabla_{h^{(k)}(x)}(-\ln(f(x)_y)) \right)^T \nabla_{a^{(k)}(x)} h^{(k)}(x) = \\ & = \left(\nabla_{h^{(k)}(x)}(-\ln(f(x)_y)) \right)^T \cdot \\ & \bullet \left[\dots \quad g'(a^{(k)}(x)_j) \quad \dots \right] \end{aligned}$$

$$h^{(k)}(x)_j = g(a^{(k)}(x)_j)$$



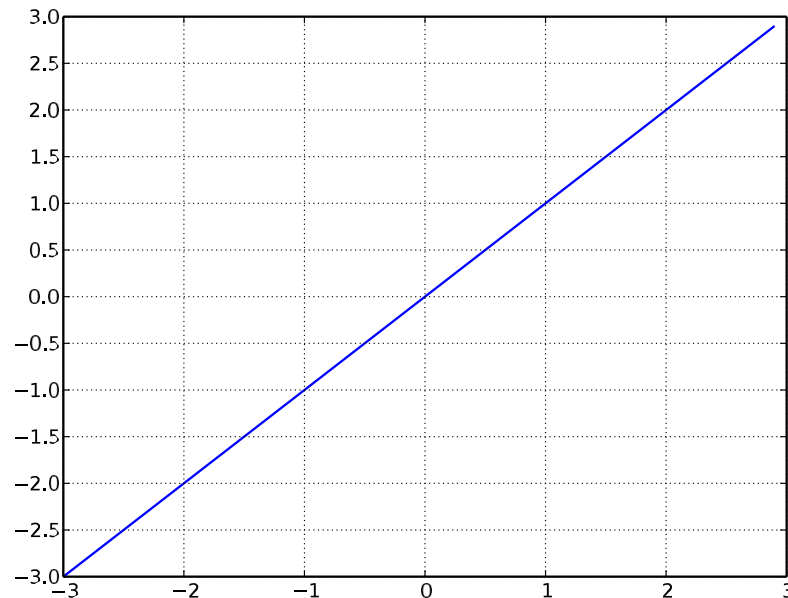
Activation function

To implement the SGD method, we need:

- a method to compute the parameters gradient

Linear activation function

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



$$g(a) = a$$

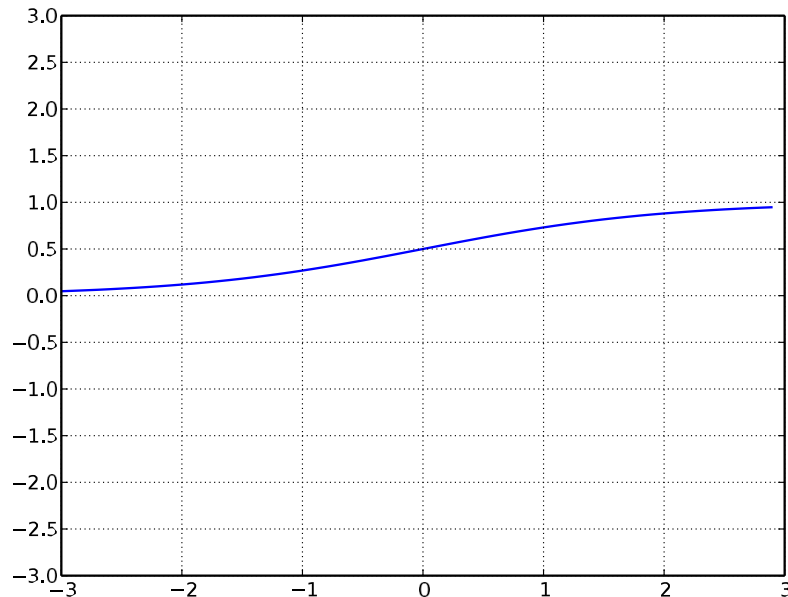
Activation function

To implement the SGD method, we need:

- a method to compute the parameters gradient

Sigmoid activation function

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



$$g(a) = \text{sigm}(a) = \frac{1}{1 + e^{-a}}$$

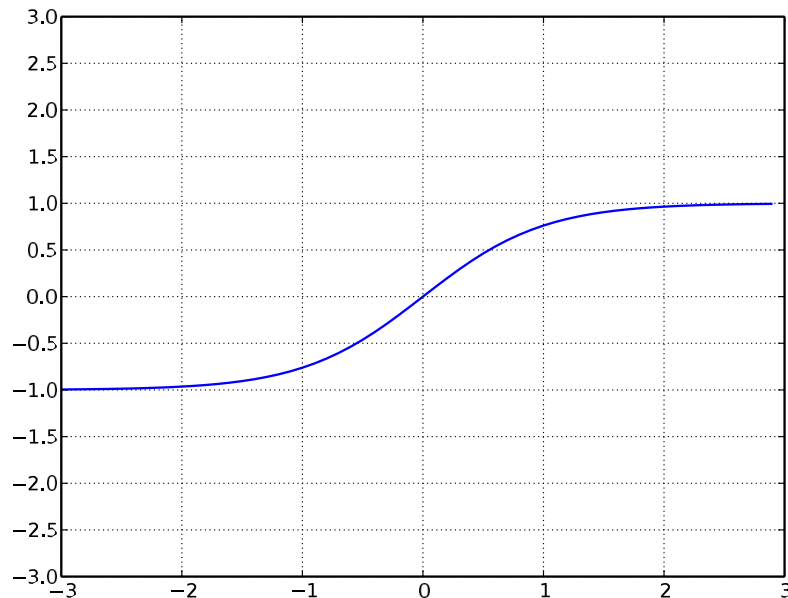
Activation function

To implement the SGD method, we need:

- a method to compute the parameters gradient

Hyperbolic tangent activation function

- Partial derivative
- $g'(a) = 1 - g(a)^2$



$$g(a) = \tanh(a) = \frac{e^a - e^{-a}}{e^a + e^{-a}} = \frac{e^{2a} - 1}{e^{2a} + 1}$$

Loss function derivative

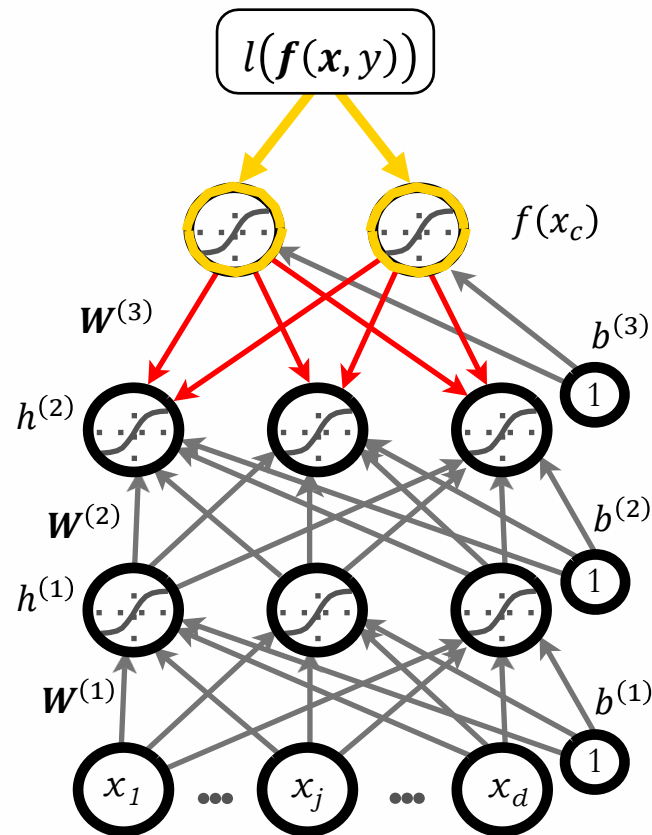
To implement the SGD method, we need:

- a method to compute the parameters gradient

Partial derivative for weights:

$$\begin{aligned} \bullet \quad \frac{\partial}{\partial W_{i,j}^{(k)}} (-\ln(f(\mathbf{x})_y)) &= \\ &= \frac{\partial(-\ln(f(\mathbf{x})_y))}{\partial a^{(k)}(\mathbf{x})_i} \frac{\partial a^{(k)}(\mathbf{x})_i}{\partial W_{i,j}^{(k)}} = \\ &= \frac{\partial(-\ln(f(\mathbf{x})_y))}{\partial a^{(k)}(\mathbf{x})_i} h^{(k-1)}(\mathbf{x})_j \end{aligned}$$

$$a^{(k)}(\mathbf{x})_i = b_i^{(k)} + \sum_j W_{i,j}^{(k)} h^{(k-1)}(\mathbf{x})_j$$



Loss function derivative

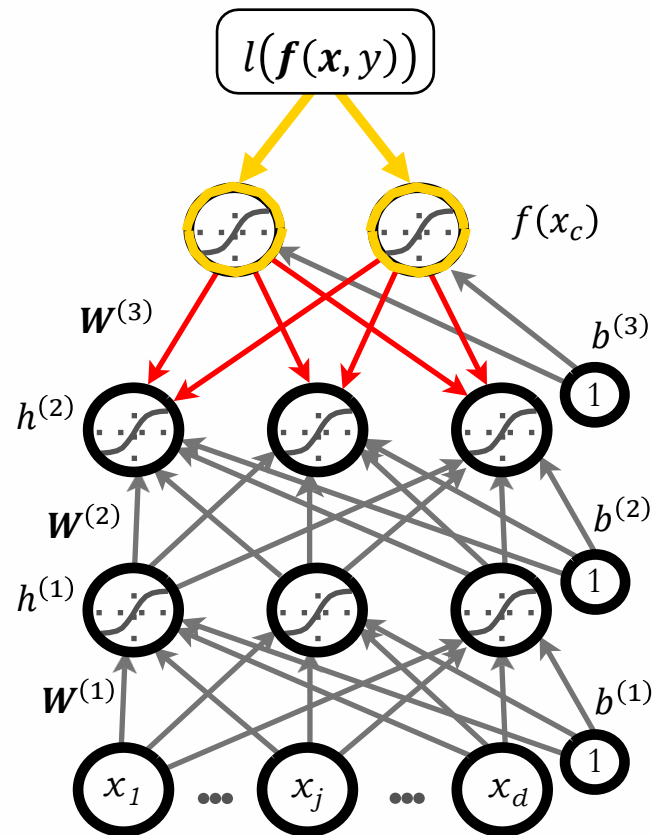
To implement the SGD method, we need:

- a method to compute the parameters gradient

Gradient for weights:

$$\begin{aligned} \nabla_{\mathbf{W}^{(k)}} \left(-\ln(f(\mathbf{x})_y) \right) &= \\ &= \nabla_{\mathbf{a}^{(k)}(\mathbf{x})} \left(-\ln(f(\mathbf{x})_y) \right) \mathbf{h}^{(k-1)}(\mathbf{x})^T \end{aligned}$$

$$\mathbf{a}^{(k)}(\mathbf{x})_i = b_i^{(k)} + \sum_j \mathbf{W}_{i,j}^{(k)} \mathbf{h}^{(k-1)}(\mathbf{x})_j$$



Loss function derivative

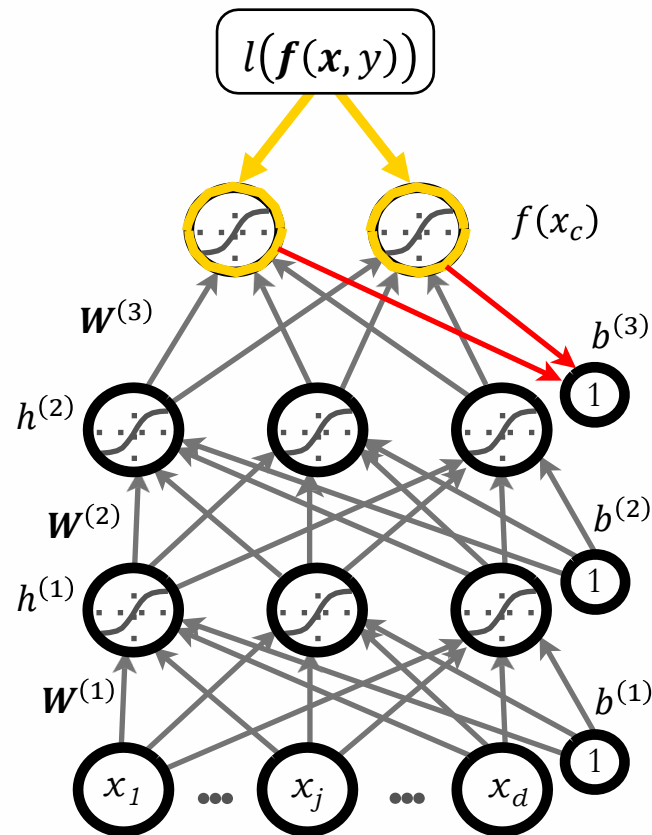
To implement the SGD method, we need:

- a method to compute the parameters gradient

Partial derivative for bias:

$$\begin{aligned} \bullet \quad \frac{\partial}{\partial b_i^{(k)}} (-\ln(f(\mathbf{x})_y)) &= \\ &= \frac{\partial(-\ln(f(\mathbf{x})_y))}{\partial a^{(k)}(\mathbf{x})_i} \frac{\partial a^{(k)}(\mathbf{x})_i}{\partial b_i^{(k)}} = \\ &= \frac{\partial(-\ln(f(\mathbf{x})_y))}{\partial a^{(k)}(\mathbf{x})_i} \end{aligned}$$

$$a^{(k)}(\mathbf{x})_i = b_i^{(k)} + \sum_j w_{i,j}^{(1)} h^{(k-1)}(\mathbf{x})_j$$



Loss function derivative

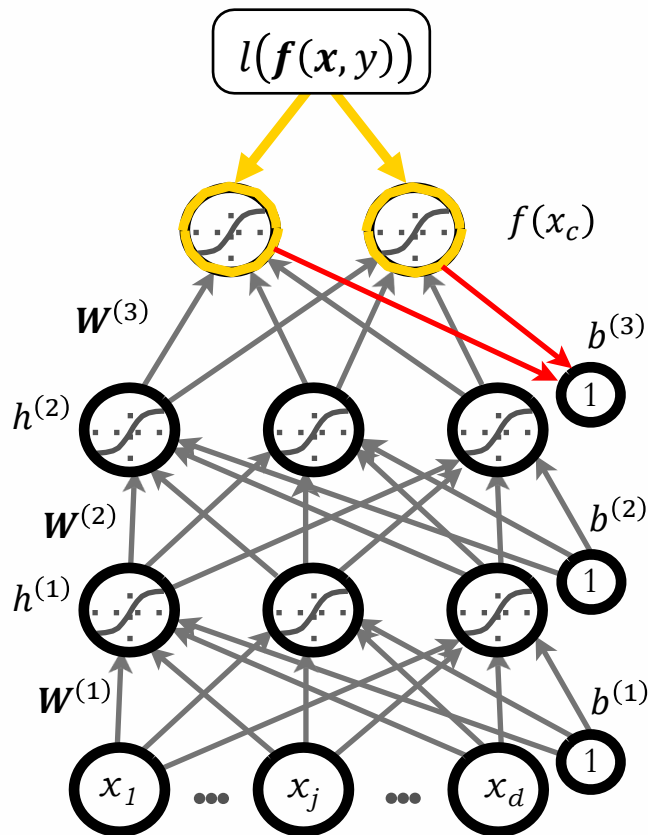
To implement the SGD method, we need:

- a method to compute the parameters gradient

Gradient for bias:

$$\begin{aligned} \nabla_{b^{(k)}}(-\ln(f(\mathbf{x})_y)) &= \\ &= \nabla_{a^{(k)}(\mathbf{x})}(-\ln(f(\mathbf{x})_y)) \end{aligned}$$

$$a^{(k)}(\mathbf{x})_i = b_i^{(k)} + \sum_j W_{i,j}^{(1)} h^{(k-1)}(\mathbf{x})_j$$



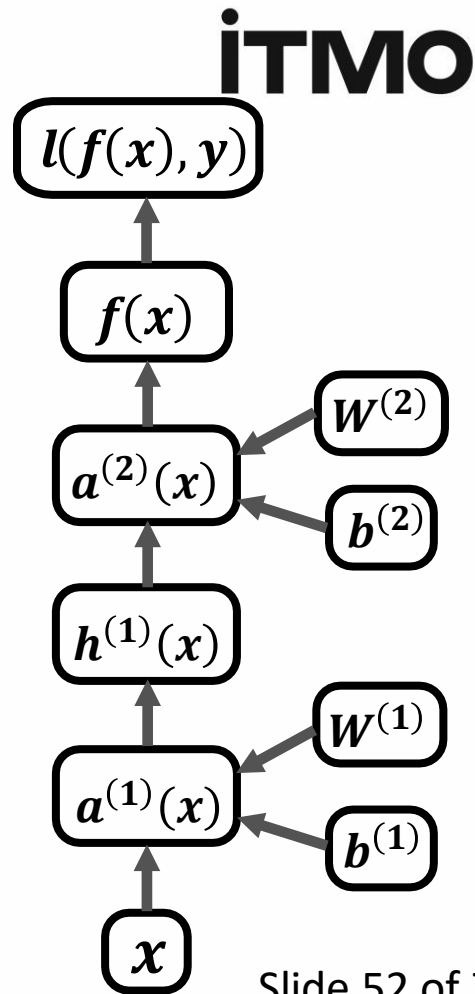
Forward propagation

Calculation of the neural network is a forward propagation

It can be represented as an acyclic flow graph

Can be easily implemented in modular way:

- Each block is a function with arguments of its children
- Calling functions in a right order implements a forward propagation

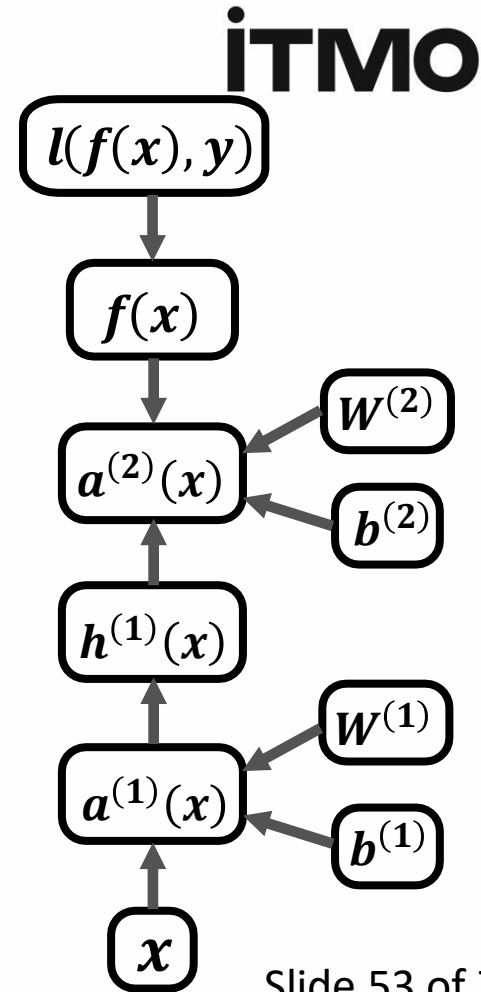


Backward propagation

Calculation of the gradient is a backward propagation

It calculates the gradient of the loss with respect to each children

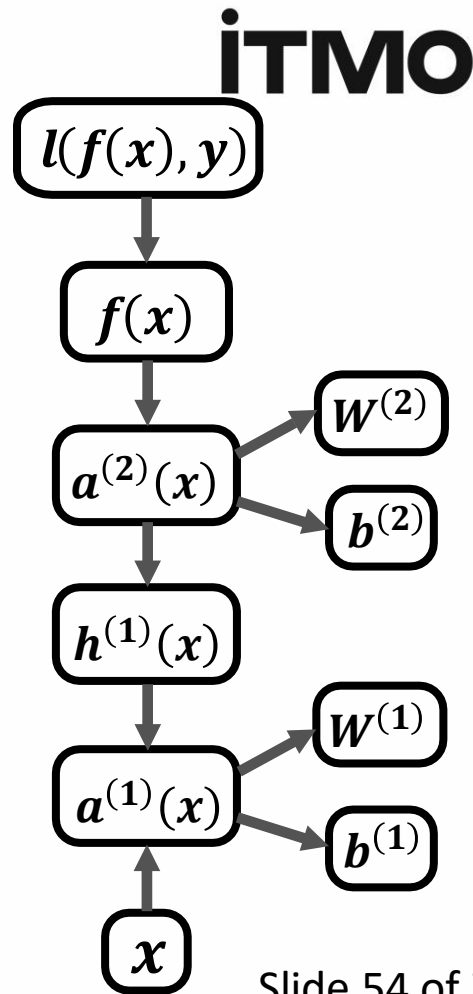
1. Compute the output layer gradient
 - $\nabla_{a^{(L+1)}(x)}(-\ln(f(x)_y)) = -(e(y) - f(x))$
2. For each k from $L + 1$ to 1
 - a. Compute the gradients of hidden layer parameters
 - b. Compute the gradient of the hidden layer below
- Backward propagation is executed until parameters θ are reached



Backward propagation

For each k from $L + 1$ to 1

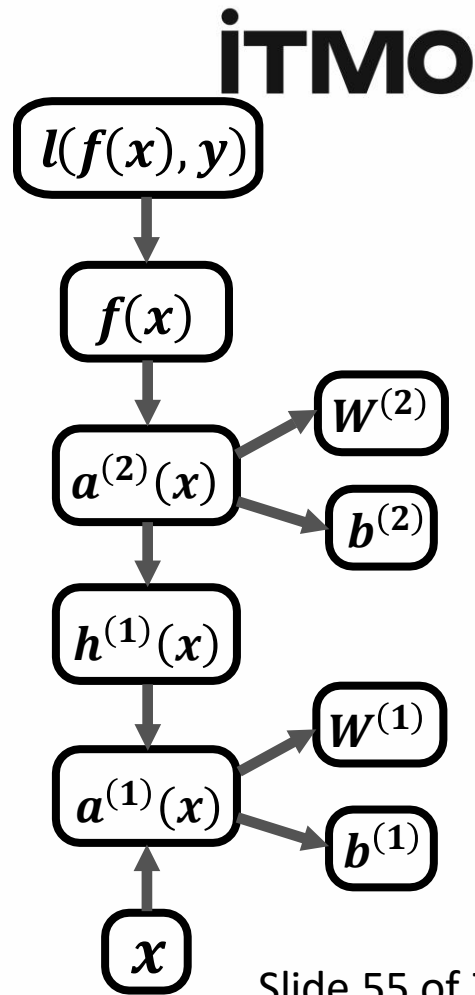
- Compute the gradients of hidden layer parameters
 - $\nabla_{W^{(k)}} (-\ln(f(x)_y)) = \nabla_{a^{(k)}(x)} (-\ln(f(x)_y)) h^{(k-1)}(x)^T$
 - $\nabla_{b^{(k)}} (-\ln(f(x)_y)) = \nabla_{a^{(k)}(x)} (-\ln(f(x)_y))$
- Compute the gradient of the hidden layer below
 - activation: $\nabla_{h^{(k-1)}(x)} (-\ln(f(x)_y)) =$
 $= W^{(k)T} \left(\nabla_{a^{(k)}(x)} (-\ln(f(x)_y)) \right)$
 - pre-activation: $\nabla_{a^{(k)}(x)} (-\ln(f(x)_y)) =$
 $= \left(\nabla_{h^{(k+1)}(x)} (-\ln(f(x)_y)) \right)^T \cdot$
 $\cdot \begin{bmatrix} \dots & g'(a^{(k)}(x)_j) & \dots \end{bmatrix}$



Backward propagation

The consistency of the forward and backward propagation methods

- The gradient can be estimated as:
 - $\frac{\partial f(x)}{\partial x} \approx \frac{f(x+\epsilon) - f(x-\epsilon)}{2\epsilon}$
 - where $f(x)$ is a loss function
 - x is a parameter
 - $f(x + \epsilon)$ is a loss value is the parameter x is increased by ϵ
 - $f(x - \epsilon)$ is a loss value is the parameter x is decreased by ϵ



Regularization

To implement the SGD method, we need:

- a regularizer function $\Omega(\boldsymbol{\theta})$ and a method to compute its gradient

L2 regularization

- $$\Omega(\boldsymbol{\theta}) = \sum_k \sum_i \sum_j \left(W_{i,j}^{(k)} \right)^2 = \sum_k \|W^{(k)}\|_F^2$$

Gradient

- $$\nabla_{W^{(k)}} \Omega(\boldsymbol{\theta}) = 2W^{(k)}$$

Can be applied on weights, but not biases

Regularization

To implement the SGD method, we need:

- a regularizer function $\Omega(\boldsymbol{\theta})$ and a method to compute its gradient

L1 regularization

- $\Omega(\boldsymbol{\theta}) = \sum_k \sum_i \sum_j |W_{i,j}^{(k)}|$

Gradient

- $\nabla_{\mathbf{W}^{(k)}} \Omega(\boldsymbol{\theta}) = \text{sign}(\mathbf{W}^{(k)})$
- where $\text{sign}(\mathbf{W}^{(k)}) = 1_{W_{i,j}^{(k)} > 0} - 1_{W_{i,j}^{(k)} < 0}$

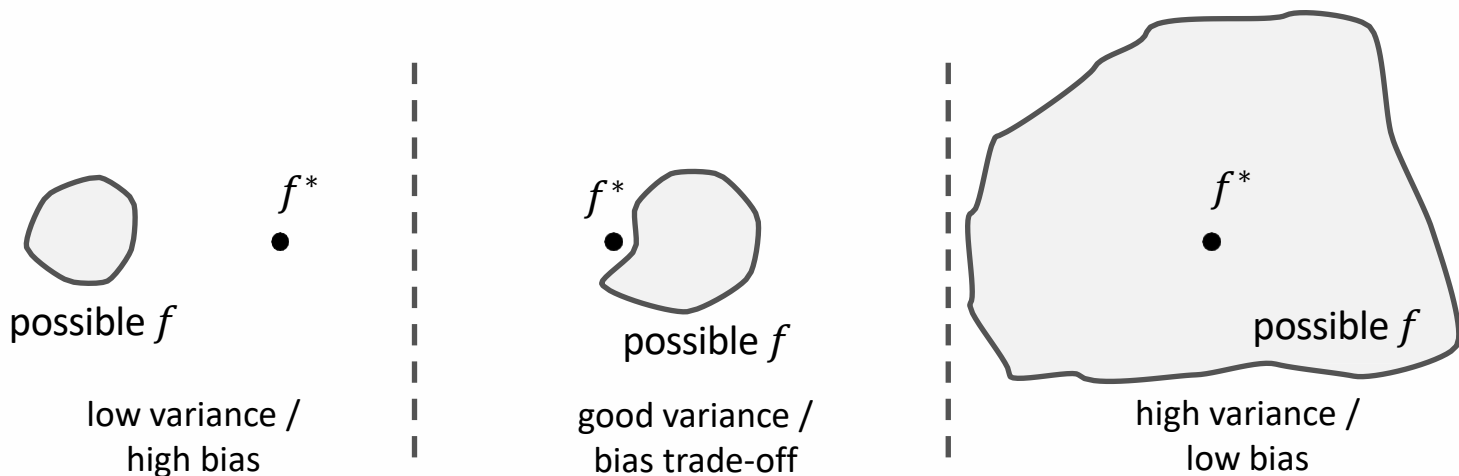
Can be applied on weights, but not biases

Unlike L2 it will push certain weights to be exactly 0

Regularization

Regularization parameter

- Variance of trained model – does it vary a lot if the training set changes
- Bias of trained model – is the average model close to the true solution f^*
- Variance / bias ratio is controlled by the λ multiplier



Parameters initialization

To implement the SGD method, we need:

- a method to initialize parameter values

Biases

- $b_i^{(k)} = 0$

Weights

- Can't initialize weights to 0 with *tanh* activation
 - as all gradients would then be equal to 0
- Can't initialize all weights to the same value
 - all hidden units in a layer will always behave the same
- Need to break symmetry

Parameters initialization

To implement the SGD method, we need:

- a method to initialize parameter values

Biases

- $b_i^{(k)} = 0$

Weights

- Idea: sample around 0 but break symmetry
 - sample each $W_{i,j}^{(k)}$ uniformly from a range $[-b, b]$
 - where $b = \frac{\sqrt{6}}{\sqrt{H_k + H_{k+1}}}$, and H_k is the number of units in $h^{(k)}(x)$
 - other values of b could also work well

Model parameters and hyper-parameters

Stochastic gradient descent method allows training a neural network

- Choose model parameters $(W_{i,j}^{(k)}, b^{(k)})$

Model hyper-parameters are parameters that define a neural network

- The neural network shape:
 - the number of hidden layers
 - the number of units in each hidden layer
- Regularization parameters
- Number of epochs
- Learning coefficients

How to select hyper-parameters?

Model dataset

Dataset

- Training dataset
 - used to train a model
- Validation dataset
 - used to select a model hyper-parameters
- Testing dataset
 - used to estimate the generalization of a trained neural network

Generalization is a behavior of the neural network on unseen examples

- the goal of the machine learning

Choosing hyper-parameters

Grid search

- Specify a list of possible values for each hyper-parameter
- Check all possible combinations of selected hyper-parameters
- Choose the best combination

Random search

- Specify a range for each hyper-parameter
- Sample a combination of hyper-parameters taking their ranges into an account
- Choose the best sample

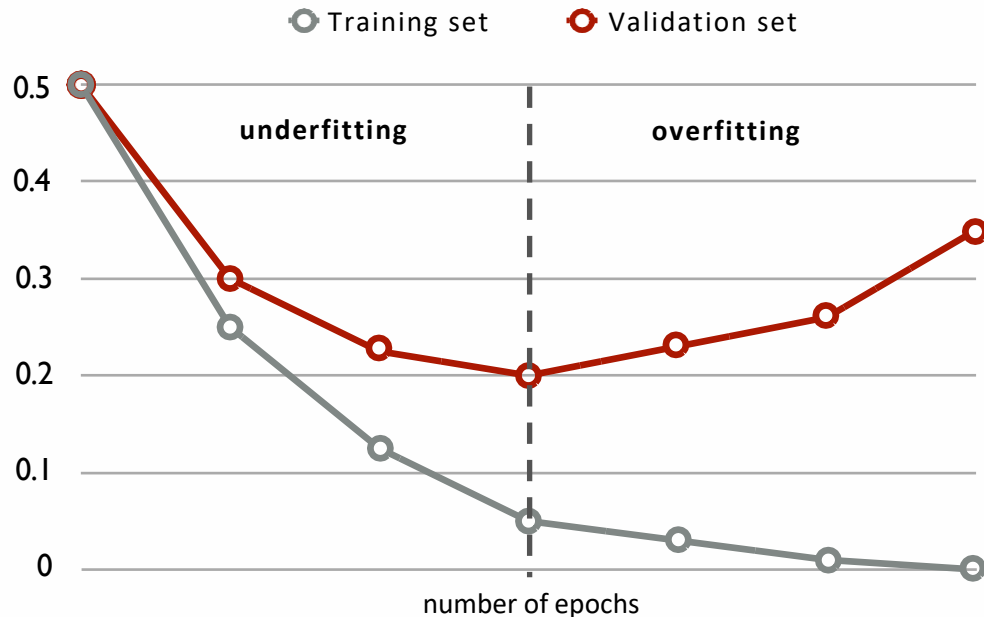
Validation set is used to estimate the neural network performance

Choosing hyper-parameters

Stop increasing a hyper-parameter when validation set error increases

Overfitting – model captured too much specific parameters of the training set and loses generalization

Underfitting – model has not captured enough training set parameters



Validation set is used to estimate the neural network performance

Model convergence

SGD convergence conditions

The learning coefficients α_t should satisfy the condition

- $\sum_{t=1}^{\infty} \alpha_t = \infty$
- $\sum_{t=1}^{\infty} \alpha_t^2 < \infty$

Decreasing strategy

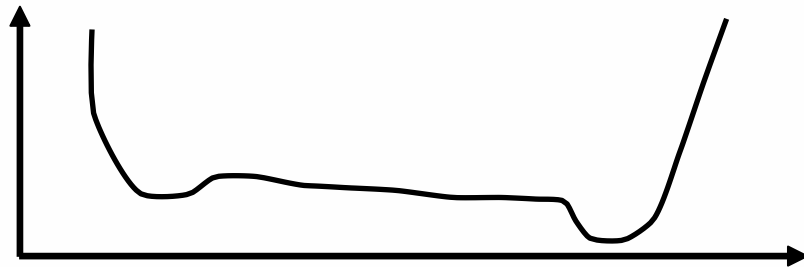
- $\alpha_t = \frac{\alpha}{1+\delta t}$
- $\alpha_t = \frac{\alpha}{t^\delta}, \delta \in (0.5; 1]$
- where α_t is a learning coefficient of the i^{th} epoch

It's advised to use the constant learning coefficient for first few epoches

Model convergence

SGD optimization

- Optimization can get stuck in local minima or plateau
 - there is no single global minimum
- The gradient is getting vanished with increasing the number of layers



Possible optimization strategies

- Use batch of examples instead of one when calculating a gradient
 - gradient is an average of batch gradients
- Use an exponential average of previous gradients
 - $\overline{\nabla}_{\theta}^{(t)} = \nabla_{\theta} l(f(x^{(t)}), y^{(t)}) + \beta \overline{\nabla}_{\theta}^{(t-1)}$

Neural Networks and Deep Learning



Neural networks are mathematical model of biological neural networks

They can be used for wide range of the categorization and decision-making tasks

There are other learning strategies as well

- Unsupervised learning
- Reinforcement learning
- etc.

There are various neural network models

- Convolutional neural networks
- Recurrent neural networks
- Long short-term memory networks
- etc.

Test

Lecture 13-14 Test

iTMO



Please scan the code to start the test

**THANK YOU
FOR YOUR TIME!**

it^{'s}**MO** *re than a*
UNIVERSITY

Andrei Zhdanov
adzhdanov@itmo.ru