

 École polytechnique fédérale de Lausanne

Content (6 weeks)

- W1 General concepts of image classification / segmentation
 Traditional supervised classification methods (RF)
- W2 Traditional supervised classification methods (SVM)
 Best practices
- W3 Elements of neural networks
- W4 Convolutional neural networks
- W5 Convolutional neural networks for semantic segmentation
- W6 Sequence modeling, change detection

And now: deep learning!

 In the rest of the classes, we will focus on a new type of model: neural networks

And now: deep learning!

- In the rest of the classes, we will focus on a new type of model: neural networks
- Neural networks are the base of what we call today deep learning

AI can figure out a place's politics by analyzing cars on Google Street View

A neural network combed through 50 million images

By Rob Verger November 30, 2017







IPEO course - 4-Neural netwo 17 October 2024

And now: deep learning!

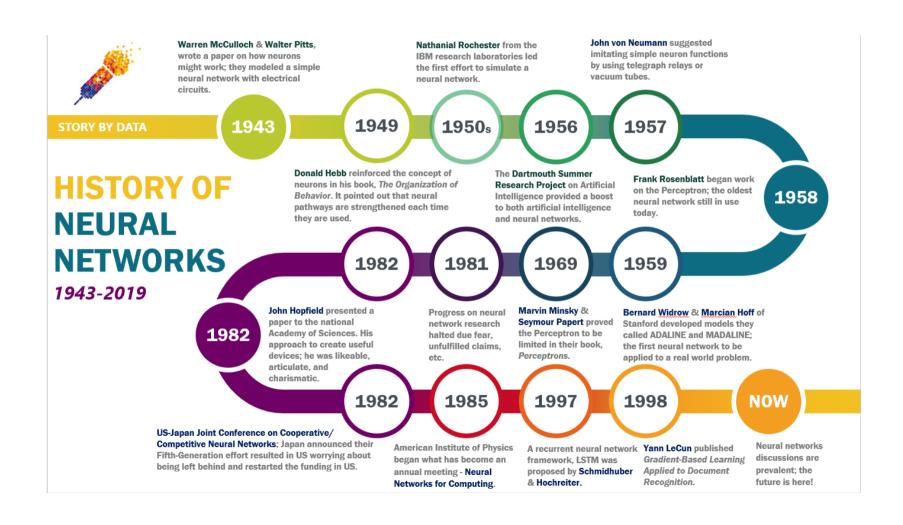
- In the rest of the classes, we will focus on a new type of model: neural networks
- Neural networks are the base of what we call today deep learning (DL)
- In the next 3 weeks we will learn about
 - neural networks in general
 - convolutional neural networks (CNNs)
 - how to use CNNs for semantic segmentation and sequence prediction
- Disclaimer: we will cover the basics only. For an in-depth study of DL, please consider more in depth courses in IC, STI.

Neural networks

History and myths

Basic principles of neural networks

A bit of history

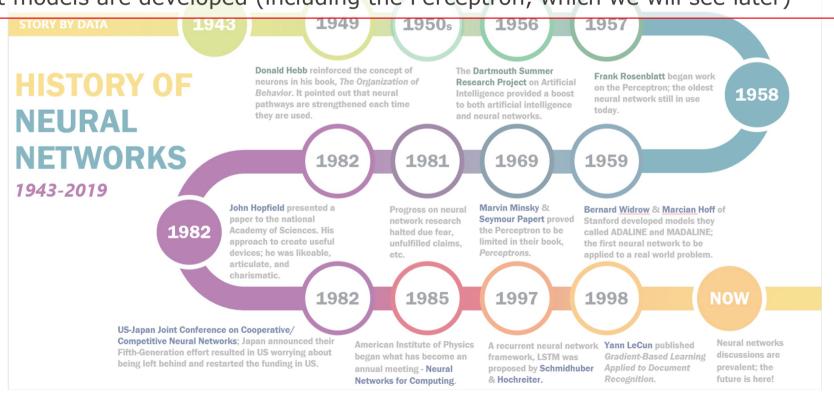


Phase 1: the birth

IPEO course – 4-Neural networks 17 October 2024 Warren McCulloch & Walter Pitts, wrote a paper on how neurons

- Basic concepts are proposed
- First models are developed (including the Perceptron, which we will see later)

John von Neumann suggested imitating simple neuron function by using telegraph relays or vacuum tubes.



A bit of history

Phase 1: the birth

EO course – 4-Neural networks October 2024 Basic concepts are proposed

First models are developed (including the Perceptron we will see later)

thanial Rochester from the
A research laboratories led imitating simple neuron function
first effort to simulate a by using telegraph relays or vacuum tubes.

Phase 2: the adolescence

- Strong research until the 80s of pointed out that requal

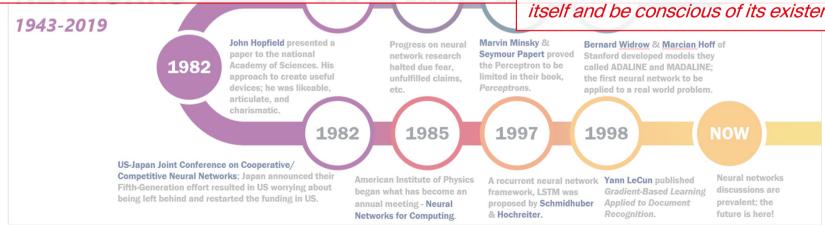
The National New York The Organization of The National New York The

Frustration due to unfulfilled claims

Considered dead until around 2000.

New York Times, July 8th 1958:

The Navy revealed the embryo of an electronic computer today that it expects will be able to walk, talk, see, write, reproduce itself and be conscious of its existence.



EPFL A bit of history

Phase 1: the birth

Basic concepts are proposed

First models are developed (including the Perceptron we will see later)

1957

Phase 2: the adolescence

- Strong research until the 80s
- Frustration due to unfulfilled claims
- Considered dead until around 2000.

Phase 3: the adult age

IPEO course - 4-Neural networks 17 October 2024

Neural nets become fashionable again!

Mainly due to improvements in hardware and size of datasets

Nowadays most research in machine learning is neural network-based

Competitive Neural Networks; Japan announced their Fifth-Generation effort resulted in US worrying about being left behind and restarted the funding in US.

American Institute of Physics began what has become an annual meeting - Neural **Networks for Computing**

framework, LSTM was proposed by Schmidhuber & Hochreiter.

Gradient-Based Learning Applied to Document Recognition.

Neural networks discussions are prevalent; the future is here!

US-Japan Joint Conference on Cooperative/

A recurrent neural network Yann LeCun published

Key moments

1943 : McCulloch and Pitts: model of neuron

1958 : Rosenblatt: perceptrons

• 1968 : Minsky and Papert point out limitations:

perceptrons are linear

■ 1982 : Hopfield network (associative memory), → Nobel Prize in Physics 2024!

Kohonen's self-organising map (clustering),

Fukushima's Neocognitron (vision)

• 1986 : Rumelhart, Hinton and Williams:

training of multilayer perceptrons

1989 : LeCun introduces convolutional networks

2012 : Hinton (re)introduces these: deep learning

So it's not that new?

- No, neural nets have been around for a while, and also DL
- What is new is that nowadays we can learn these models efficiently
- Since we have

Data



ImageNet

Natural images
1000 classes (person, car, cat, ...)

>14 Million images





Large computational resources







And this also in EO!

EPFL Ohh, but that's so far fetched!

- Well ... I am sure you can name a few apps that you use that make use of neural nets.
 - Google translate
 - Alexa
 - Photo upload in Whatsapp/Facebook
 - Amazon web search



Okok, but that's super complicated!



- True, if you want to master all the complexity of it
- Research in fundamentals of DL is just starting
- For using these algorithms, there is A LOT of resources.
- In particular python libraries that allow building DL models "easily".
- The one you will use: Pytorch.
- There are many others (Keras, Tensorflow, ...)
- But before diving into Pytorch, let's learn how NN and DL work.

Neural networks

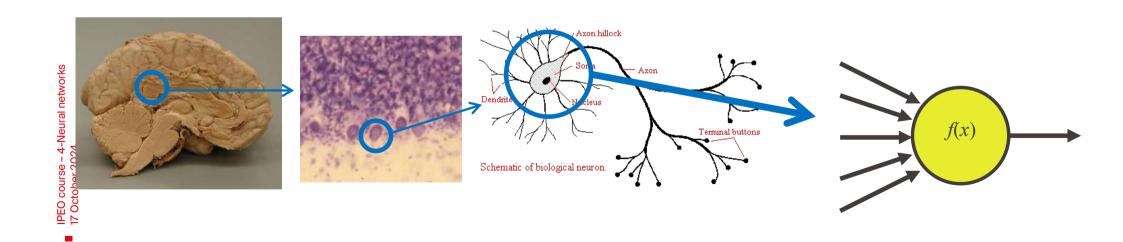
History and myths

Basic principles

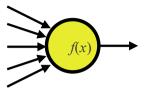
- Artificial neurons

The mathematical model of a neuron

- Some (not all!) networks are originally inspired by the brain
- Neural nets are large, densely interconnected networks of simple processing units, a.k.a neurons
- It remains just a parallel, the artificial neuron is just an approximation!



The neuron model

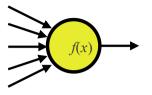


The neuron was first proposed by Warren McCulloch and Walter Pitts in 1943. It was modeled with electrical circuits.

The output y was modeled as a weighted linear combination of inputs X_i (i are the variables describing each one of your samples)

$$y = \sum_{j} \beta_{j} x_{j} + \beta_{0}$$

The neuron model



The neuron was first proposed by Warren McCulloch and Walter Pitts in 1943. It was modeled with electrical circuits.

The output y was modeled as a weighted linear combination of inputs x_j

Moreover, a transfer function, or activation function was used to make a decision:

$$y = f\left(\sum_{j} \beta_{j} x_{j} + \beta_{0}\right)$$

Activation functions

$$y = \int \left(\sum_{j} \beta_{j} x_{j} + \beta_{0}\right)$$

Let's call v the result of the parenthesis. Examples of functions:

$$f(v) = \begin{cases} 1 & v \ge 0 \\ 0 & v < 0 \end{cases}$$

Classifies into two groups

$$f(v) = \frac{1}{1 + \exp(-v)} \qquad \int$$

Outputs numbers between 0 and 1

$$f(v) = v$$



Outputs the linear combination itself

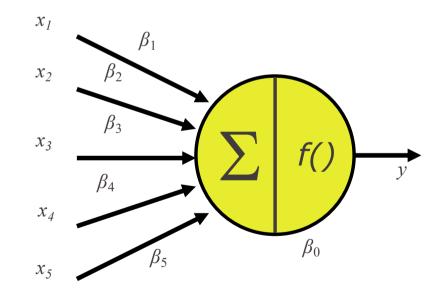
The artificial neuron

 A network made of many McCulloch-Pitts neurons can approximate any function, given the right weights β.

But how to find the right weights?

We will use gradient descent

$$y = f\left(\sum_{j} \beta_j x_j + \beta_0\right)$$



Neural networks (NN)

History and myths

Basic principles of NN

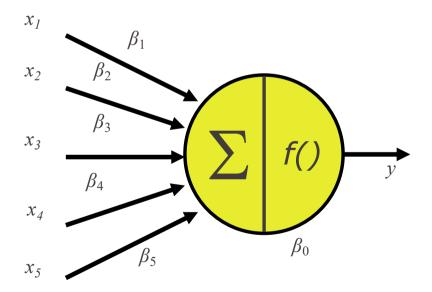
- Artificial neurons
- Gradient descent

IPEO course – 4-Neural networks 17 October 2024

Gradient descent in a nutshell

$$y = f\left(\sum_{j} \beta_j x_j + \beta_0\right)$$

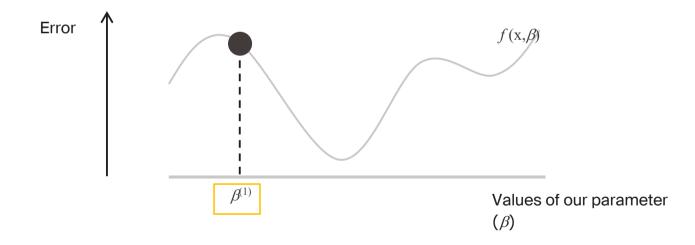
- We want to "find the β_s ".
- Roughly
 - We initialise
 - We calculate error of the current set of weights
 - We modify the β_s to make the error decrease
 - · We repeat until satisfaction



Gradient descent

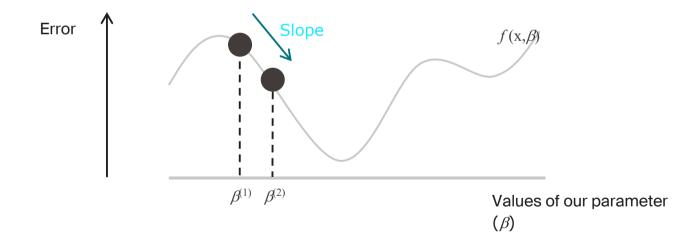
We want to find the minimum of an error function f(x,β), which you don't know

We start with an initial guess of the parameter β



Gradient descent

- We want to find the minimum of an error function f(x,β), which you don't know
- We start with an initial guess of the parameter β
- We change its value in the direction of maximal slope



Gradient descent - how

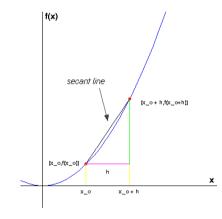
- We want to find the minimum of an error function
- We start with an initial guess of the parameter β
- We change its value in the direction of maximal slope

The derivative measures the steepness of the graph of a function at some particular point on the graph.

Thus, the derivative is a slope.

$$\frac{\partial f(x)}{\partial x} | x_0 = f'(x_0) = \lim_{h \to 0} \frac{f(x_0 + h) - f(x_0)}{h}$$

It is also a function, it varies according to x



Source:

https://www.ugrad.math.ubc.ca/coursedoc/math100/notes/derivs/deriv5.html

Gradient descent - how

To update a parameter β , we substract the value of its derivative from it

$$\beta(r+1) = \beta(r) - \alpha \frac{\partial f(x)}{\partial x}|_{\beta(r)}$$

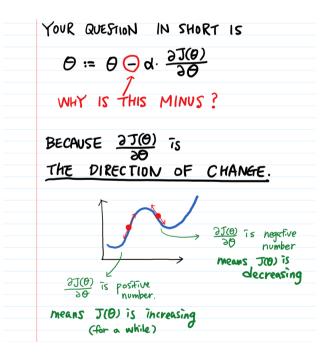
- α is called a learning rate.
- r is the iteration

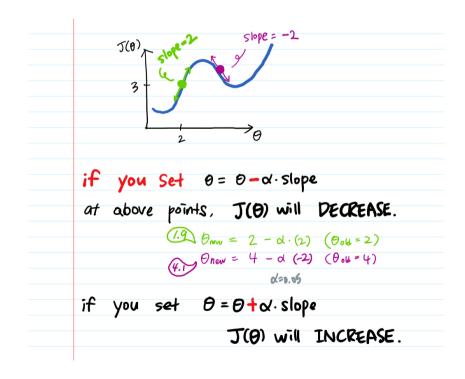
Remember, the derivative is a slope

Gradient descent - how

Source: https://medium.com/@aerinykim/why-do-we-subtract-the-slope-a-in-gradient-descent-73c7368644fa

So why the negative sign?

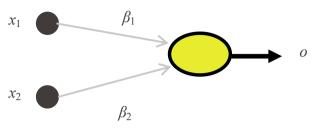


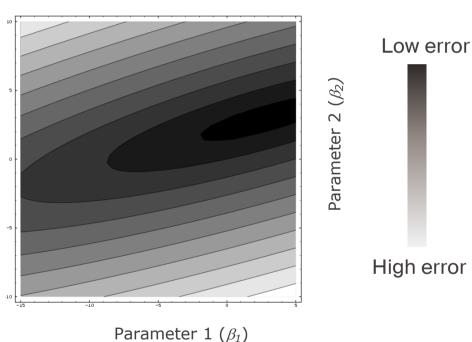


IPEO course – 4-Neural networks 17 October 2024

Some intuition

- In the 2D case this is simpler to visualize
- We consider a model with 2 parameters, β_1 and β_2
- (the grey tone of the plot is the error, the darker the smaller)

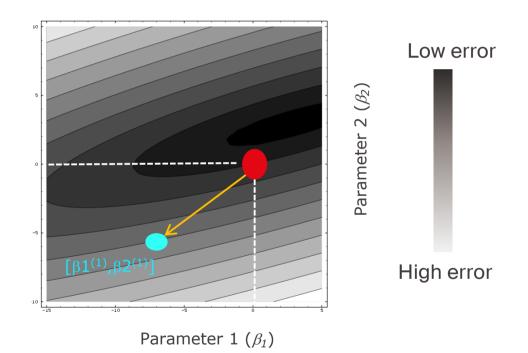




Some intuition

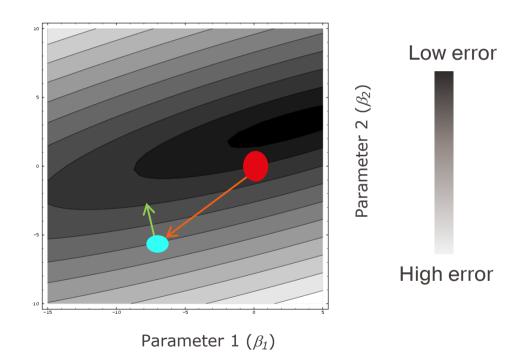
- In the 2D case this is simpler to visualize
- Your initial parameter $[\beta_1^{(1)}, \beta_2^{(1)}]$ is a 2D vector from the origin [0,0]

• Where would the next gradient step go?



Some intuition

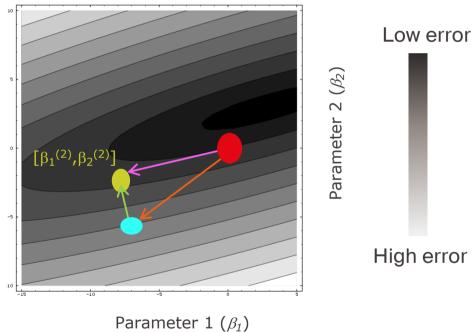
- In the 2D case this is simpler to visualize
- Your initial parameter $[\beta_1^{(1)}, \beta_2^{(1)}]$ is a 2D vector from the origin [0,0]
- We compute the direction of maximal slope



D. Tuia. ECEO

Some intuition

- In the 2D case this is simpler to visualize
- Your initial parameter $[\beta_1^{(1)}, \beta_2^{(1)}]$ is a 2D vector from the origin [0,0]
- We compute the direction of maximal slope
- The new parameter $[\beta_1^{(2)}, \beta_2^{(2)}]$ is the sum of the two vectors

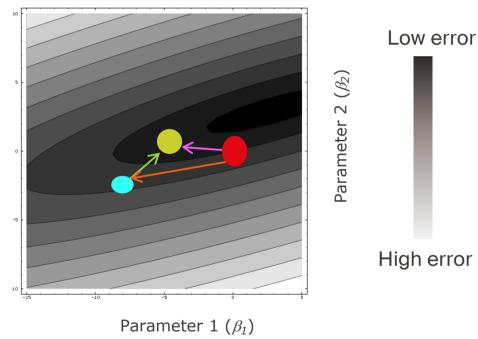


D. Tuia. ECEO

EPFL

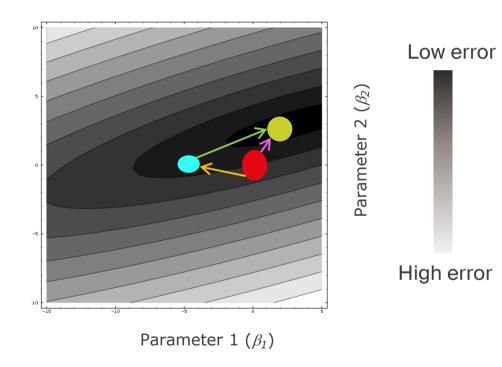
... and then again

- In the 2D case this is simpler to visualize
- Your initial parameter $[\beta_1^{(2)}, \beta_2^{(2)}]$ is a 2D vector from the origin [0,0]
- We compute the direction of maximal slope
- The new parameter $[\beta_1^{(3)}, \beta_2^{(3)}]$ is the sum of the two vectors



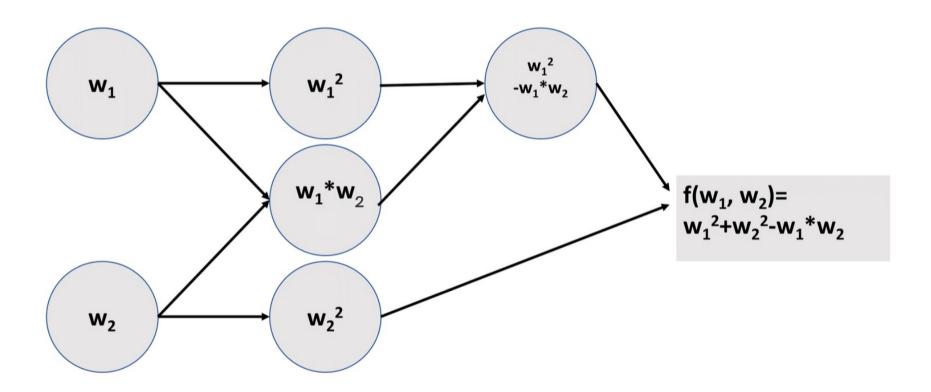
... and again

- In the 2D case this is simpler to visualize
- Your initial parameter $[\beta_1^{(3)}, \beta_2^{(3)}]$ is a 2D vector from the origin [0,0]
- We compute the direction of maximal slope
- The new parameter $[\beta_1^{(4)}, \beta_2^{(4)}]$ is the sum of the two vectors



A little exercice

Let's take this neural network



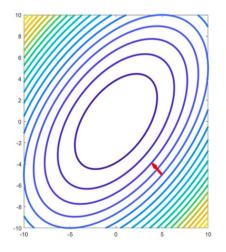
IPEO course – 4-Neural networks 17 October 2024

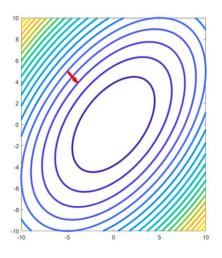
A little exercice

$$f(w_1, w_2) = w_1^2 + w_2^2 - w_1 * w_2$$

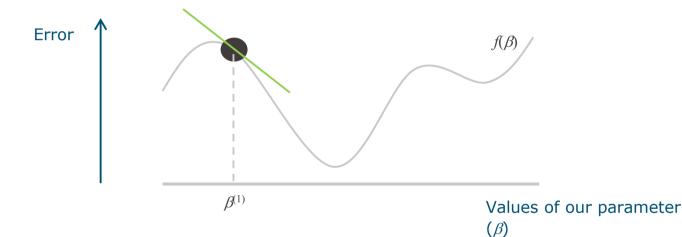
$$\nabla f(w_1, w_2) = \left[????,???\right]^{\top}$$

Calculate the gradient at $(w_1, w_2) = (5, -5)$ and at $(w_1, w_2) = (-5, 5)$



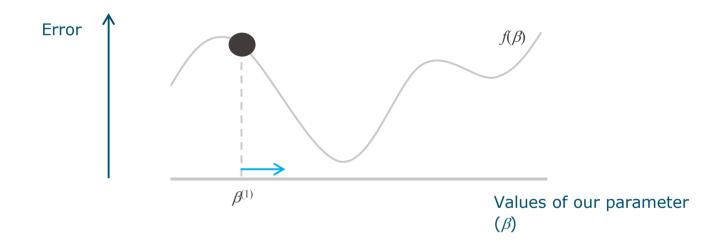


Remember 1: calculate the derivative



$$\beta(r+1) = \beta(r) - \alpha \frac{\partial f(x)}{\partial x}|_{\beta(r)}$$

Remember 2: α gives you the step vector (how far you go)



$$\beta(r+1) = \beta(r) - \boxed{\alpha} \frac{\partial f(x)}{\partial x}|_{\beta(r)}$$

Remember 3: this gives you the new value of β and the new error

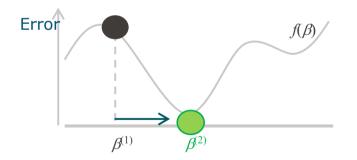
Error Values of our parameter (β)

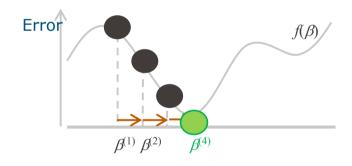
$$\beta(r+1) = \beta(r) - \alpha \frac{\partial f(x)}{\partial x}|_{\beta(r)}$$

The step value (or learning rate)

$$\beta(r+1) = \beta(r) - \alpha \frac{\partial f(x)}{\partial x}|_{\beta(r)}$$

It decides by how much you multiply the step vector Large α can lead to faster convergence than small

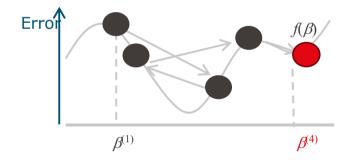




The learning rate

$\beta(r+1) = \beta(r) - \alpha \frac{\partial f(x)}{\partial x}|_{\beta(r)}$

• It decides by how much you multiply the step vector Too large α can lead to disaster



EPFL Momentum

- The last trick is called momentum
- Momentum discourages sharp changes of direction in the descent.
- It basically adds a term going in the direction of the previous step

$$\Delta \beta = \beta(r) - \beta(r-1)$$

Gradient descent like before

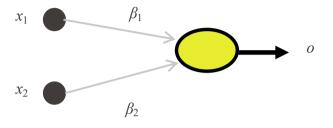
momentum

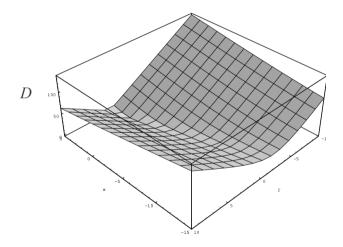
$$\beta(r+1) = \beta(r) - \alpha \frac{\partial f(x)}{\partial x}|_{\beta(r)} + \eta \Delta \beta$$

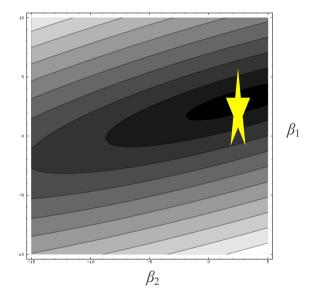


Numerical example

Perceptron with two weights $[\beta_1, \beta_2]$

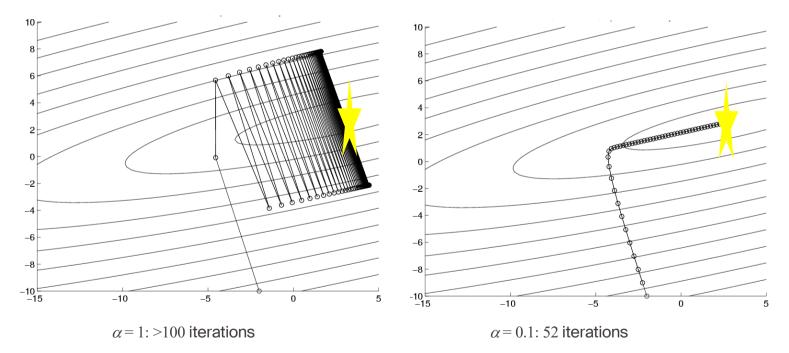






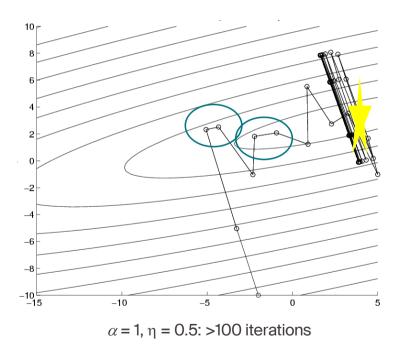
EPFL Numerical example

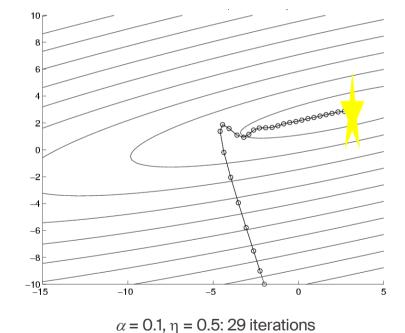
Learning rate controls oscillation and speed



EPFL Numerical example

Momentum uses a bit of the previous step





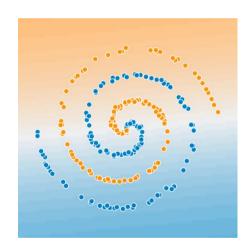
Neural networks (NN)

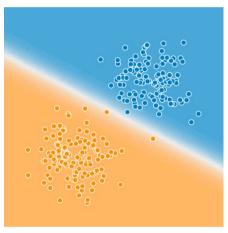
History and myths

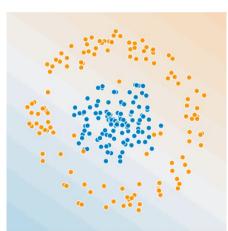
Basic principles of NN

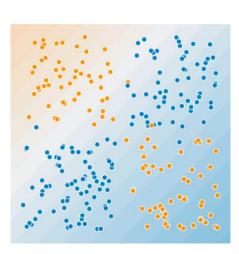
- Artificial neurons
- Gradient descent
- Multi-layer perceptror

Which data sets can the perceptron handle?







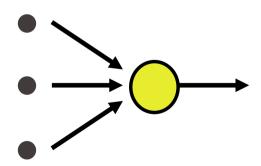




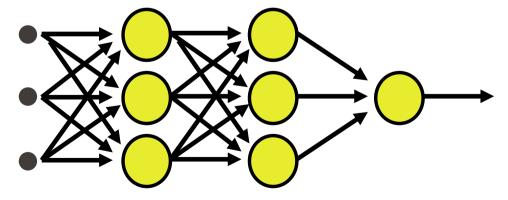


- The perceptron remained a linear classifier.
- It needed a nonlinear version.
- The Multilayer perceptron (MLP) was formalized in 1986 by Rumelhart and colleagues (including Geoffrey Hinton, considered to be one of the "fathers" of deep learning and Nobel prize for physics in 2024).

FROM the perceptron



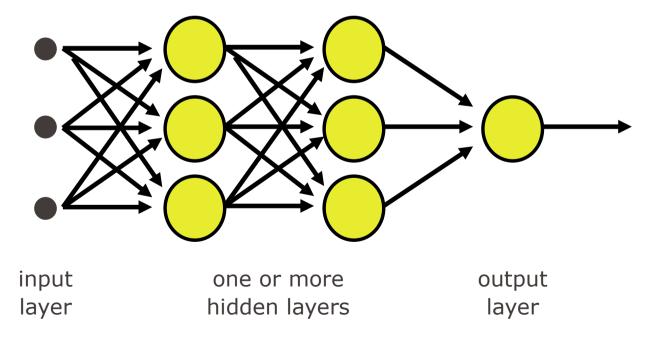
TO the MLP

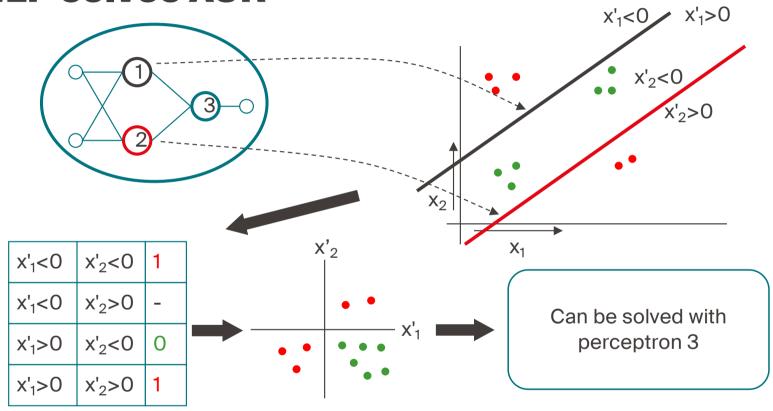


IPEO course – 4-Neural networks 17 October 2024

The multilayer perceptron (MLP)

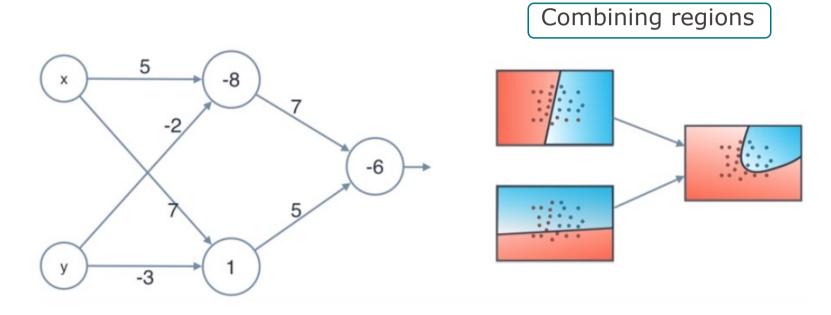
- It's a feed forward network: it goes from inputs to outputs, without loops
- Every neuron includes an activation function (e.g. sigmoid, see earlier slides 18-19)





IPEO course - 4-Neural networks 17 October 2024

Multi-Layer Perceptron



Courtesy Luis Serrano

IPEO course - 4-Neural networ

Neural networks (NN)

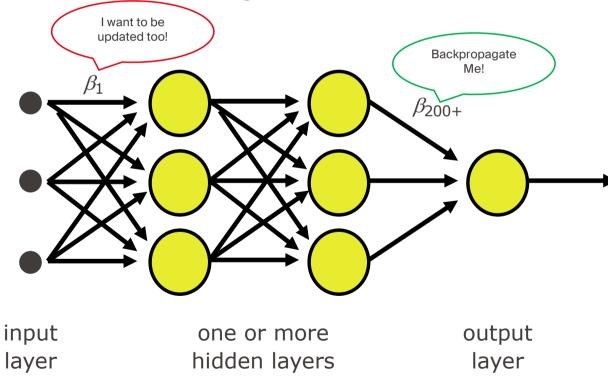
History and myths

Basic principles of NN

- Artificial neurons
- Gradient descent
- Multi-layer perceptron
- Training the MLP

The multilayer perceptron (MLP)

- How to compute gradient descent on a cascade of operations?
- How to reach all trainable weights?



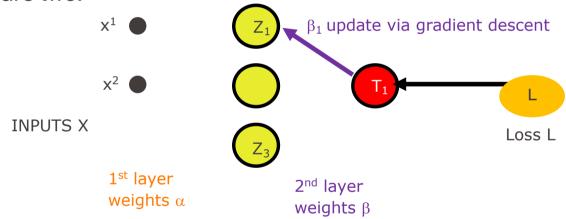
INPUTS X

• There are two: $x^1 \quad \bullet \quad X_1 \quad X_2 \quad \bullet \quad X_1 \quad \beta_1 \quad X_2 \quad \bullet \quad X_2 \quad \bullet \quad X_3 \quad X_4 \quad \bullet \quad X_4 \quad \bullet \quad X_5 \quad \bullet \quad X_6 \quad \bullet \quad X_7 \quad \bullet \quad X_8 \quad \times \quad X_8 \quad \bullet \quad X_8$

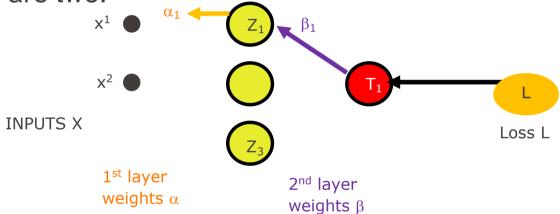
 1^{st} layer weights α

2nd layer weights β **OUTPUT Y**

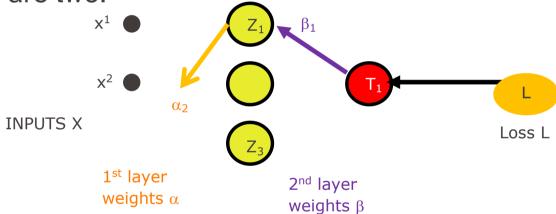




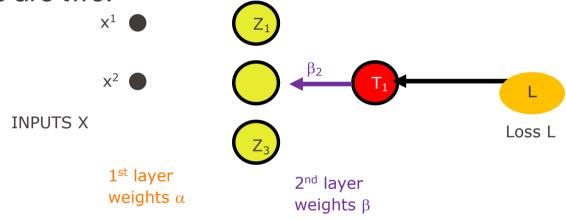
- 1. The forward pass: you put a data point in and obtain the prediction
- 2. The Backward pass: you update parameters by back-propagating the loss



- 1. The forward pass: you put a data point in and obtain the prediction
- 2. The Backward pass: you update parameters by back-propagating the loss



- 1. The forward pass: you put a data point in and obtain the prediction
- 2. The Backward pass: you update parameters by back-propagating the loss



- 1. The forward pass: you put a data point in and obtain the prediction
- 2. The Backward pass: you update parameters by back-propagating the loss

And one to rule them all: the chain rule



- But how to reach weights that are earlier in the network?
- We use the chain rule
- It's a classical rule of derivatives:

"the derivative of a function applied to another function is the product of their derivatives"

$$\frac{\partial f(g(x))}{\partial x} = f'(g(x)) \cdot g'(x)$$



And one to rule them all: the chain rule

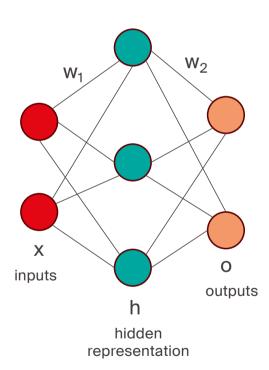


- But how to reach weights that are earlier in the network?
- We use the chain rule
- It's a classical rule of derivatives:

"the derivative of a function applied to another function is the product of their derivatives"

$$\frac{\partial f(g(h(x)))}{\partial x} = f'(g(h(x))) \cdot g'(h(x)) \cdot h'(x)$$

How to apply the chain rule: the network



Forward propagation:

$$z = w_1 x$$

$$h = \phi(z)$$

$$o = w_2 h$$

is an activation function (e.g. sigmoid)

Loss:

$$L = \ell(o, y)$$

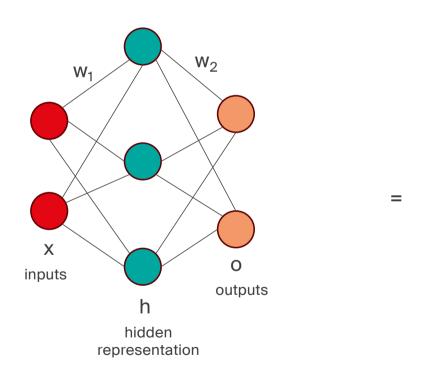
y is the ground truth

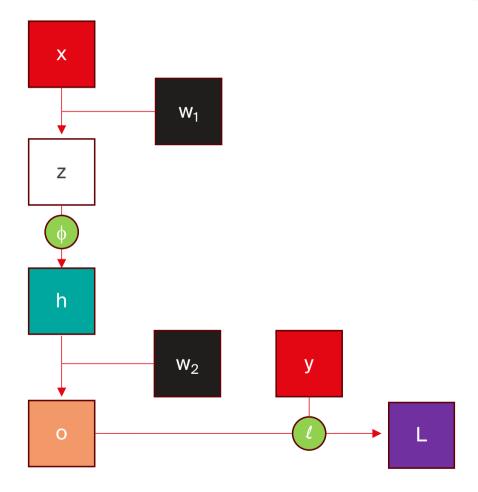
IPEO course – 4-Neural networks 17 October 2024

IPEO course - 4-Neural networks 17 October 2024

EPFL

How to apply the chain rule: the network is a computational graph



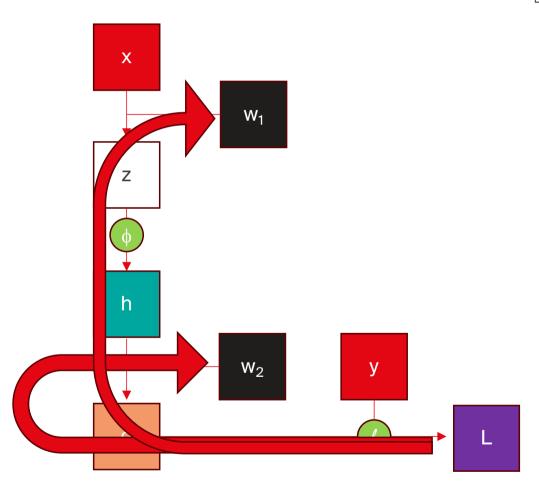


How to apply the chain rule: goal

Find the best weights w₁
 and w₂ to minimise Loss L

 $\frac{\partial L}{\partial w_1}$

 $\frac{\partial L}{\partial w_2}$

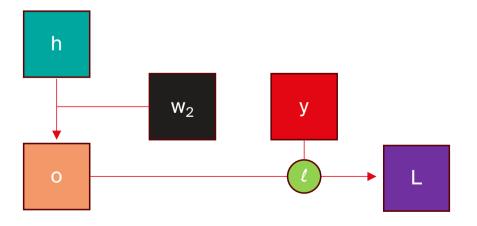


IPEO course - 4-Neural netw

$$\frac{\partial L}{\partial w_2} = \frac{\partial L}{\partial o} \cdot \frac{\partial o}{\partial w_2}$$
$$= \frac{\partial L}{\partial o} \cdot h^{\top}$$

since
$$o = w_2 \cdot h$$

$$\frac{\partial o}{\partial w_2} = h^{\top}$$



How to calculate $\frac{O}{2}$

- s = softmax = scales output between 0 and 1 (slide 19)
- $\hat{y}_i = \frac{e^{o_i}}{\sum_j e^{o_j}}$
- = classification loss = 0 if correct answer,> 0 otherwisewe use the cross entropy

$$L = \ell(\hat{y}, y) = -\log p(y|x)$$

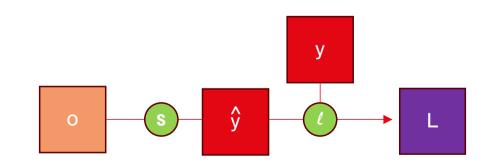
$$= -\sum_{i} y_{i} \log \hat{y}_{i}$$

$$= -\sum_{i} y_{i} \log \left[\frac{e^{o_{i}}}{\sum_{j} e^{o_{j}}}\right]$$

$$= -\sum_{i} y_{i} [o_{i} - \log \sum_{j} e^{o_{j}}]$$

$$= \sum_{i} y_{i} \log \sum_{j} e^{o_{j}} - \sum_{i} y_{i} o_{i}$$

$$= \log \sum_{j} e^{o_{j}} - \sum_{i} y_{i} o_{i}$$



How to calculate

= classification loss = 0 if correct answer, > 0 otherwise

we use the cross entropy

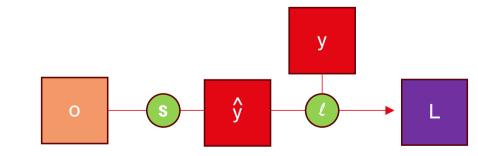
$$L = \ell(\hat{y}, y) = \log \sum_{i} e^{o_i} - \sum_{i} y_i o_i$$

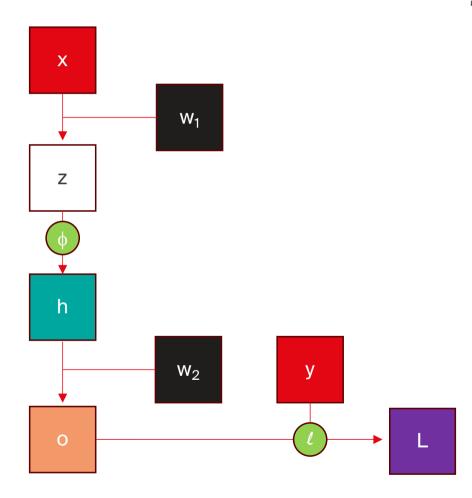
since $\frac{\partial f}{\partial x} \log f(x) = \frac{1}{f(x)} f'(x)$

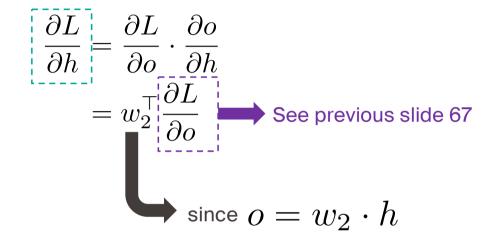
$$\frac{\partial L}{\partial o_i} = \frac{e^{o_i}}{\sum_j e^{o_j}} - y_i$$
$$= \hat{y}_i - y_i$$

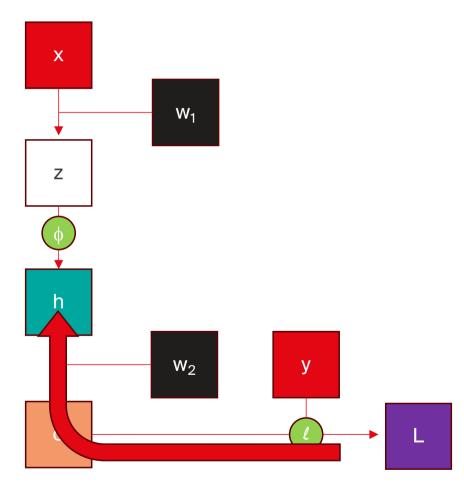
Prediction = between 0 and 1

True class = 1 or 0









$$\begin{vmatrix} \frac{\partial L}{\partial h} \end{vmatrix} = \frac{\partial L}{\partial o} \cdot \frac{\partial o}{\partial h}$$

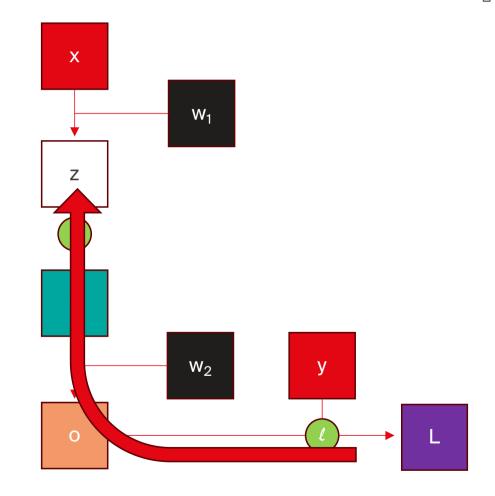
$$= w_2^{\top} \frac{\partial L}{\partial o}$$

$$\frac{\partial L}{\partial z} = \begin{vmatrix} \frac{\partial L}{\partial h} \end{vmatrix} \frac{\partial h}{\partial z}$$

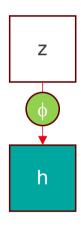
$$= \frac{\partial L}{\partial h} \odot \phi'(z)$$
?????

since
$$h=\phi(z)$$

$$\frac{\partial h}{\partial z}=\phi'(z)$$

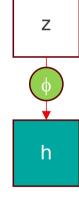


How to calculate $\phi'(z)$



How to calculate $\phi'(z)$

- Activation functions add nonlinearity to the network
- Otherwise it's a chain of linear operations
- We saw several functions at slide 19, e.g. the softmax a few slides ago (s)





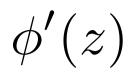
$$y = f\left(\sum_{j} \beta_{j} x_{j} + \beta_{0}\right)$$

Let's call \(\nu\) the result of the parenthesis. Examples of functions

$$f(v) = rac{1}{1 + \exp(-v)}$$
 Outputs numbers between 0 and

IPEO course – 4-Neural networks 17 October 2024

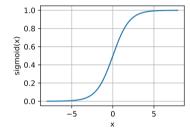
How to calculate $\phi'(z)$ **EPFL**



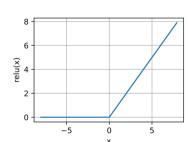


Sigmoid (or softmax)

- Output in [0,1]
- Close to a threshold function, with a smooth transition from 0 to 1
- In contrast to real threshold function, sigmoid is differentiable
- Often used at output neurons to get the probability for binary classification
- Usually not used in hidden layers: slow updating, vanishing gradient



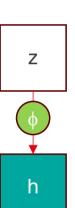






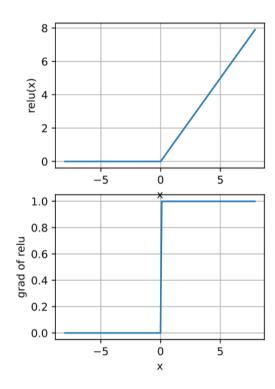
ReLU (Rectified Linear Unit)

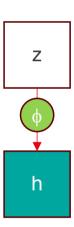
- ReLU(x) = max(x, 0)
- Its derivatives are simple: either they vanish (0) or they just let the argument through without modifications (1).
- This makes optimization better behaved and it mitigated well-documented problem of *vanishing gradients*
- So it is often used within intermediate neurons



How to calculate $\phi'(z)$ when using ReLU

$$\blacksquare \operatorname{ReLU}(x) = \max(x, 0)$$



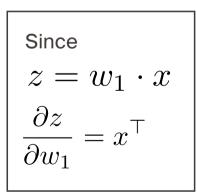


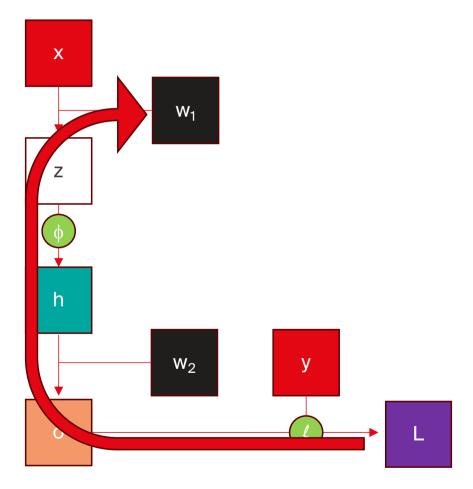
$$\frac{\partial L}{\partial h} = \frac{\partial L}{\partial o} \cdot \frac{\partial o}{\partial h}$$

$$= w_2^{\top} \frac{\partial L}{\partial o}$$

$$\frac{\partial L}{\partial z} = \frac{\partial L}{\partial h} \cdot \frac{\partial h}{\partial z}$$
$$= \frac{\partial L}{\partial h} \odot \phi'(z)$$

$$\frac{\partial L}{\partial w_1} = \frac{\partial L}{\partial z} \cdot \frac{\partial z}{\partial w_1}$$
$$= \frac{\partial L}{\partial z} x^{\top}$$





EPFL In summary

- NN update their parameters via gradient descent
- They have many parameters, recombined in a feed forward way
- They add nonlinearities at each step
- They can approximate very complex functions (so they overfit easily!)
- Remember to cross-validate!
 - Number of neurons
 - Number of layers
 - Learning rate
 - Momentum
 - •
- In the next course we will see how to use this to build convolutional neural networks