

# MSBigData

## Part: Neural Networks

Geoffroy Peeters

LTCI - Télécom ParisTech

2018



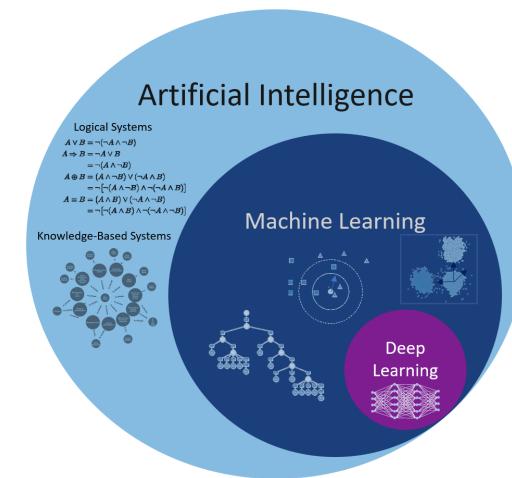
- 1. Introduction
- 1.1 Deep Learning and Neural Networks : History
- 1.2 Neural Networks get Deeper
- 1.3 What has changed ? Why now ?
- 1.4 Applications
- 1.5 Some of the famous Deep Learning people
- 1.6 Three main types of Nets
- 2. Multi-Layer-Perceptron (MLP)
- 2.1 Perceptron
- 2.2 Logistic Regression (0 hidden layers)
- 2.3 Logistic regression, Softmax regression, Multi-label, Multi-Class
- 2.4 Chain rule and Back-propagation
- 3. Backpropagation in matrix form
- 3.1 for a single example
- 3.2 for a whole mini-batch
- 3.3 Neural Networks (1 hidden layer)
- 3.4 Deep Neural Networks (> 2 hidden layers)
- 3.5 Computation Graph (from Alexandre Allauzen)
- 4. Deep Learning Frameworks
- 4.1 Playground
- 4.2 Example of an MLP
- 4.3 Example of an MLP in python
- 4.4 Example of an MLP in pytorch
- 4.5 Example of an MLP in tensorflow
- 4.6 Tensorboard
- 4.7 Example of an MLP in Keras
- 5. Evaluating the performances of a classification system
- 5.1 Train, Dev and Test sets
- 5.2 Measuring the performances
- 5.3 Various types of training
- 6. Making things work
- 6.1 Activation functions  $a = g(z)$
- 6.2 Weight initialization
- 6.3 Regularization
- 6.4 Alternative gradient descent algorithms
- 6.5 Normalization inputs
- 6.6 Batch Normalization (BN)

Geoffroy Peeters ~ LTCI / Télécom ParisTech - 1/20

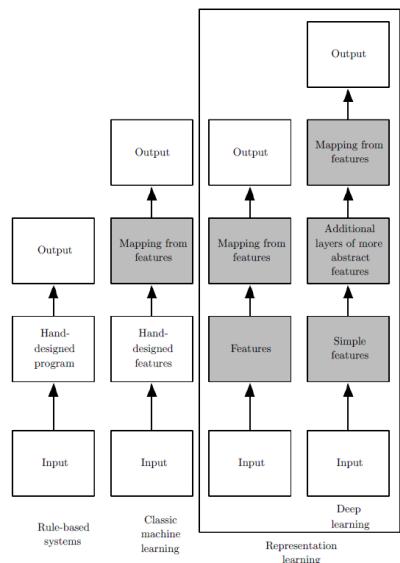
Geoffroy Peeters ~ LTCI / Télécom ParisTech - 2/20

## Deep Learning and Neural Networks : History

## Deep learning (a subset of machine learning)

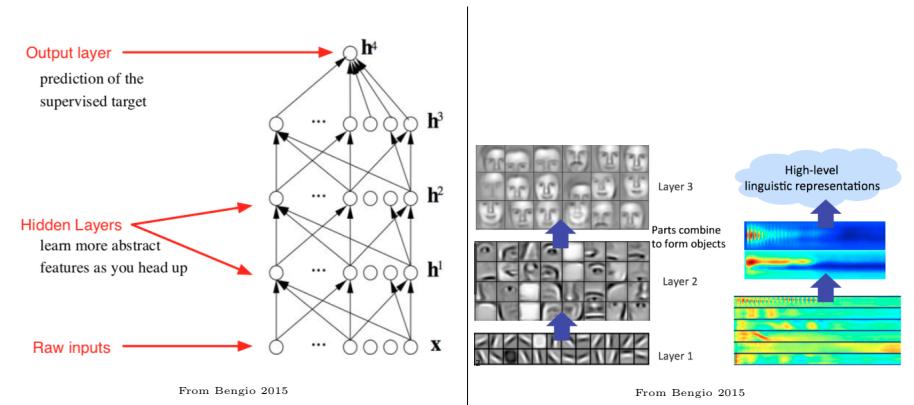


## Deep learning : learning hierarchical representations



From Deep Learning by Ian Goodfellow and Yoshua Bengio [Geoffroy Peeters - LIP6 / Télécom ParisTech - 6](#)

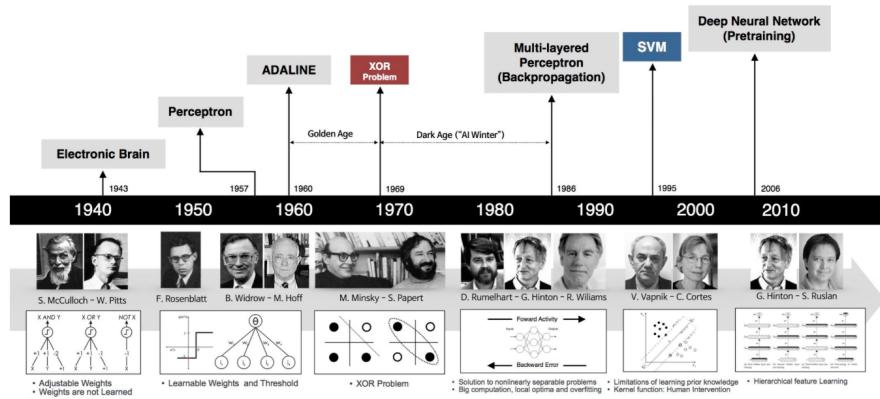
## Deep learning : learning hierarchical representations



From Bengio 2015

Geoffroy Peeters - LIP6 / Télécom ParisTech - 7

## Deep Learning and Neural Networks : History



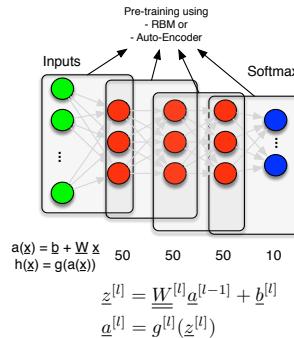
Geoffroy Peeters - LIP6 / Télécom ParisTech - 8

## Three main types of Nets

## Three main types of Nets

### Multi Layers Perceptron (MLP) or Fully-Connected (FC)

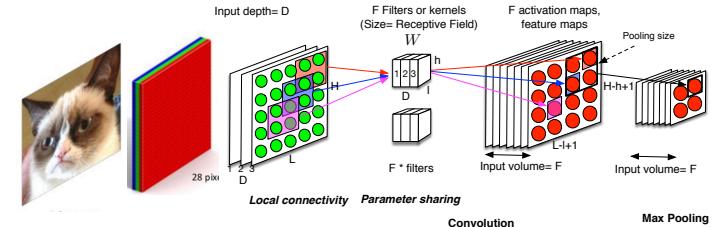
Multi Layers Perceptron (Fully Connected)



## Three main types of Nets

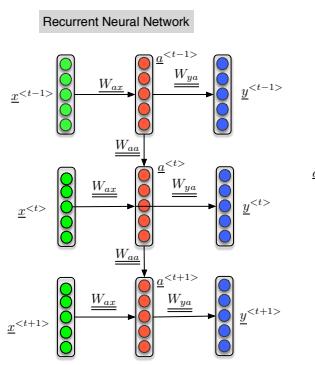
### Convolutional Neural Networks (CNN)

Convolutional Neural Network



## Three main types of Nets

### Recurrent Neural Networks (RNN)



## Perceptron

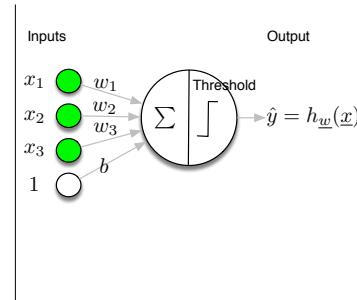
## Perceptron

- **Perceptron ?**

- A binary **classification algorithm** :
  - output is 0 or 1
- Models the decision using a linear function, followed by a threshold
 
$$\hat{y} = h_{\underline{w}}(\underline{x}) = \text{Threshold}(\underline{w} \cdot \underline{x}) \quad (1)$$
  - where
    - $\text{Threshold}(z) = 1$  if  $z \geq 0$
    - $\text{Threshold}(z) = 0$  if  $z < 0$

- **Parameters :**

- weights  $\underline{w}$  and bias  $b$



Geoffroy Peeters - LTCI / Télécom ParisTech - 34) Q C

## Perceptron

### Training

- **Training ?**

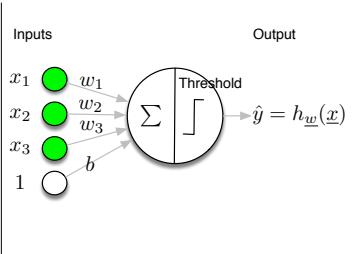
- Adapt  $\underline{w}$  and  $b$  so that  $\hat{y} = h_{\underline{w}}(\underline{x})$  gives the correct answer  $y$  on training data

- **Algorithm :**

- for each pair  $(\underline{x}^{(i)}, y^{(i)})$ 
  - compute  $h_{\underline{w}}(\underline{x}^{(i)}) = \text{Threshold}(\underline{w} \cdot \underline{x}^{(i)})$
  - if  $\hat{y}^{(i)} \neq y^{(i)}$  update
    - $w_d \leftarrow w_d + \alpha(y^{(i)} - \hat{y}^{(i)})x_d^{(i)} \quad \forall d$
    - where  $\alpha$  is the learning rate

- **Three cases :**

- $(y^{(i)} - \hat{y}^{(i)}) = 0$ 
  - $\Rightarrow$  no update
- $(y^{(i)} - \hat{y}^{(i)}) = 1$ 
  - $\Rightarrow$  the weights are too low
    - $\Rightarrow$  increase  $w_d$  for positive  $x_d^{(i)}$
- $(y^{(i)} - \hat{y}^{(i)}) = -1$ 
  - $\Rightarrow$  the weights are too high
    - $\Rightarrow$  decrease  $w_d$  for positive  $x_d^{(i)}$



Geoffroy Peeters - LTCI / Télécom ParisTech - 35) Q C

## Perceptron

### Formulating the training problem as an optimization problem

- In the case of the **Perceptron**, we want to **minimize the loss**

$$\mathcal{L}(y^{(i)}, h_{\underline{w}}(\underline{x}^{(i)})) = -(y^{(i)} - h_{\underline{w}}(\underline{x}^{(i)}))\underline{w} \cdot \underline{x}^{(i)}$$

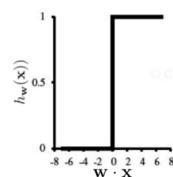
- if the prediction is good

- $\Rightarrow$  the cost is 0

- if the prediction is wrong

- $\Rightarrow$  the cost is the distance between  $\underline{w} \cdot \underline{x}^{(i)}$  and the necessary threshold for the prediction to be good

- Not exactly, since the derivative of  $h_{\underline{w}}(\underline{x}^{(i)})$  does not exist at  $\underline{w} \cdot \underline{x} = 0$
- $\Rightarrow$  Using the Threshold function can lead to instability during training



- Gradient descent ?

$$w_d \leftarrow w_d - \alpha \frac{\partial \mathcal{L}(y^{(i)}, h_{\underline{w}}(\underline{x}^{(i)}))}{\partial w_d} \quad \forall d$$

- Gradient ?

$$\begin{aligned} \frac{\partial \mathcal{L}(y^{(i)}, h_{\underline{w}}(\underline{x}^{(i)}))}{\partial w_d} &= - \left( \frac{\partial (y^{(i)} - h_{\underline{w}}(\underline{x}^{(i)}))}{\partial w_d} \underline{w} \cdot \underline{x}^{(i)} + (y^{(i)} - h_{\underline{w}}(\underline{x}^{(i)})) \cdot x_d^{(i)} \right) \\ &\simeq - \left( 0 + (y^{(i)} - h_{\underline{w}}(\underline{x}^{(i)})) \cdot x_d^{(i)} \right) \\ &= -(y^{(i)} - h_{\underline{w}}(\underline{x}^{(i)})) \cdot x_d^{(i)} \end{aligned}$$

Geoffroy Peeters - LTCI / Télécom ParisTech - 37) Q C

## Perceptron

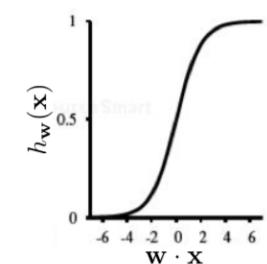
### Making things correct

- In **Logistic Regression** instead of predicting a class, we will **predict a probability** to belong to this class

$$p(y=1|\underline{x}) = h_{\underline{w}}(\underline{x}) = \text{Logistic}(\underline{w} \cdot \underline{x}) = \frac{1}{1 + e^{-\underline{w} \cdot \underline{x}}}$$

- We choose the class (0 or 1) with the largest probability

- If  $h_{\underline{w}}(\underline{x}) \geq 0.5$  choose 1
- If  $h_{\underline{w}}(\underline{x}) < 0.5$  choose 0



- We define the Loss as (binary cross-entropy)

$$\mathcal{L}(y^{(i)}, h_{\underline{w}}(\underline{x}^{(i)})) = -y^{(i)} \log(h_{\underline{w}}(\underline{x}^{(i)})) - (1 - y^{(i)}) \log(1 - h_{\underline{w}}(\underline{x}^{(i)}))$$

- Gradient descent ?

$$w_d \leftarrow w_d - \alpha \frac{\partial \mathcal{L}(y^{(i)}, h_{\underline{w}}(\underline{x}^{(i)}))}{\partial w_d} \quad \forall d$$

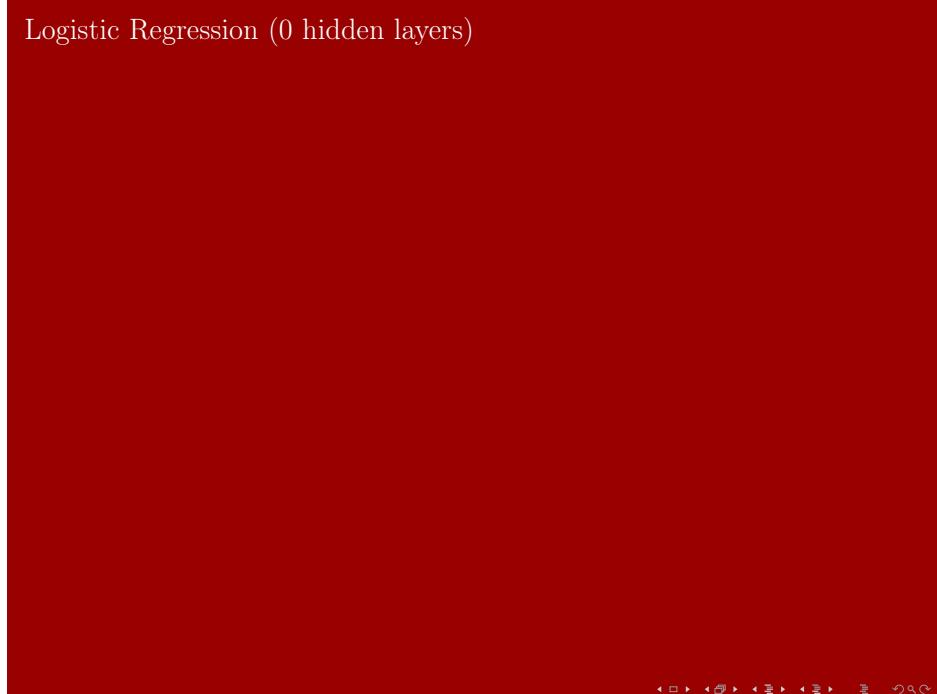
- Gradient ?

$$\frac{\partial \mathcal{L}(y^{(i)}, h_{\underline{w}}(\underline{x}^{(i)}))}{\partial w_d} = -(y^{(i)} - h_{\underline{w}}(\underline{x}^{(i)})) \cdot x_d^{(i)} \quad \forall d$$

- The update rule is therefore the same as for the Perceptron but the definition of the Loss and of  $h_{\underline{w}}(\underline{x})$  is different

Geoffroy Peeters - LTCI / Télécom ParisTech - 38) Q C

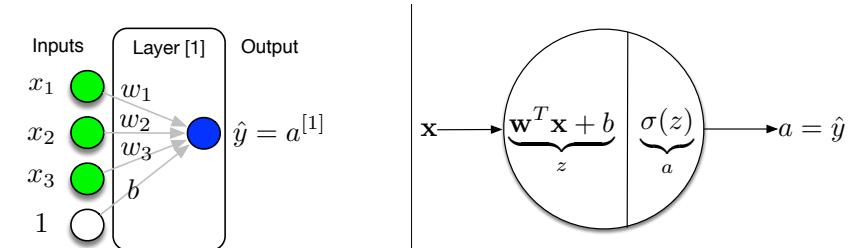
## Logistic Regression (0 hidden layers)



## Logistic Regression (0 hidden layers)

$$x \rightarrow [w_1 x_1 + w_2 x_2 + b] \xrightarrow{z} [\sigma(z)] \xrightarrow{\hat{y}=a} [\mathcal{L}(a, y)]$$

- Problem :
  - Given  $\underline{x} \in \mathbb{R}^{n_x}$ , we want to predict  $\hat{y} = P(y = 1 | \underline{x})$  with  $0 \leq y \leq 1$
- Model :
  - $\hat{y} = \sigma(\underline{w}^T \underline{x} + b)$
  - Parameters :  $\underline{w} \in \mathbb{R}^{n_x}$  and  $b \in \mathbb{R}$
- Sigmoid function :  $\sigma(z) = \frac{1}{1+e^{-z}}$ 
  - If  $z$  is very large then  $\sigma(z) \approx \frac{1}{1+0} = 1$
  - If  $z$  is very (negative) small then  $\sigma(z) = 0$



Geoffroy Peeters - LTCI / Télécom ParisTech - 40

## Logistic Regression (0 hidden layers)

### Cost Function (Empirical risk minimization (ERM))

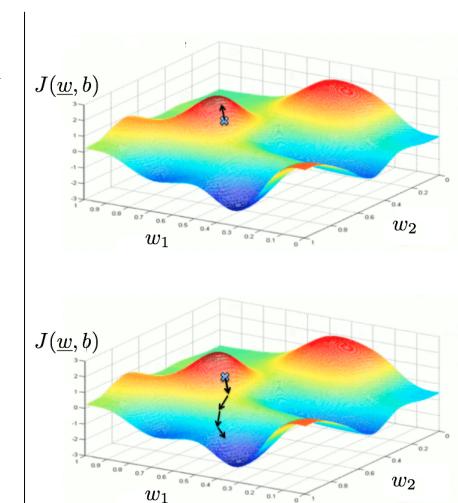
- **Training :**
  - Given  $\{(\underline{x}^{(1)}, y^{(1)}), (\underline{x}^{(2)}, y^{(2)}), \dots, (\underline{x}^{(m)}, y^{(m)})\}$ ,
  - We want  $\hat{y}^{(i)} \simeq y^{(i)}$
- How to measure ? **Loss  $\mathcal{L}$  (error) function**
  - Mean Square Error (MSE)
    - $\mathcal{L}(\hat{y}, y) = \frac{1}{2}(\hat{y} - y)^2 \Rightarrow$  not convex, many local minima
  - Binary Cross-Entropy
    - $\mathcal{L}(\hat{y}, y) = -(y \log(\hat{y}) + (1 - y) \log(1 - \hat{y}))$ 
      - If  $y = 1 \Rightarrow \mathcal{L}(\hat{y}, y) = -\log(\hat{y}) \Rightarrow$  want  $\log(\hat{y})$  large  $\Rightarrow$  want  $\hat{y}$  large
      - If  $y = 0 \Rightarrow \mathcal{L}(\hat{y}, y) = -\log(1 - \hat{y}) \Rightarrow$  want  $\log(1 - \hat{y})$  large  $\Rightarrow$  want  $\hat{y}$  small
- For all training example : **Cost  $J$  function**

$$\begin{aligned} J(\underline{w}, b) &= \frac{1}{m} \sum_{i=1}^m \mathcal{L}(\hat{y}^{(i)}, y^{(i)}) \\ &= -\frac{1}{m} \sum_{i=1}^m \left( y^{(i)} \log(\hat{y}^{(i)}) + (1 - y^{(i)}) \log(1 - \hat{y}^{(i)}) \right) \end{aligned}$$
- We want to find the values  $\underline{w}$  and  $b$  that minimize  $J(\underline{w}, b)$

## Logistic Regression (0 hidden layers)

### Gradient Descent

- **How to minimize  $J(\underline{w}, b)$  ?**
- The gradient  $\frac{\partial J(\underline{w}, b)}{\partial \underline{w}}$  points in the direction of the greatest rate of increase of the function
- We will go in the opposite direction :  $-\frac{\partial J(\underline{w}, b)}{\partial \underline{w}}$ 
  - We move down-hill in the steepest direction
- **Gradient descent :**
  - Repeat
    - $\underline{w} \leftarrow \underline{w} - \alpha \cdot d\underline{w}$
    - $b \leftarrow b - \alpha \cdot db$
  - where  $\alpha$  is the "learning rate"
- Notation :
  - we will use  $d\underline{w}$  for  $\frac{\partial J(\underline{w}, b)}{\partial \underline{w}}$
  - we will use  $db$  for  $\frac{\partial J(\underline{w}, b)}{\partial b}$



Geoffroy Peeters - LTCI / Télécom ParisTech - 42

## Logistic Regression (0 hidden layers)

$$x \rightarrow [w_1x_1 + w_2x_2 + b] \xrightarrow{z} [\sigma(z)] \xrightarrow{\hat{y}=a} [\mathcal{L}(a, y)]$$

### Forward propagation

$$\begin{aligned} z &= \underline{w}^T \underline{x} + b \\ \hat{y} &= a = \sigma(z) \end{aligned}$$

with  $\sigma(z) = \frac{1}{1+e^{-z}}$ .

## Logistic Regression (0 hidden layers)

$$x \rightarrow [w_1x_1 + w_2x_2 + b] \xrightarrow{z} [\sigma(z)] \xrightarrow{\hat{y}=a} [\mathcal{L}(a, y)]$$

### Backward propagation

$$\begin{aligned} da &= \frac{\partial \mathcal{L}}{\partial a} = \frac{-y}{a} + \frac{(1-y)}{(1-a)} \\ dz &= \frac{\partial \mathcal{L}}{\partial z} = \frac{\partial \mathcal{L}}{\partial a} \frac{\partial a}{\partial z} = \dots = a - y \\ dw_1 &= \frac{\partial \mathcal{L}}{\partial w_1} = \frac{\partial \mathcal{L}}{\partial z} \frac{\partial z}{\partial w_1} = x_1 \cdot dz \\ dw_2 &= \frac{\partial \mathcal{L}}{\partial w_2} = \frac{\partial \mathcal{L}}{\partial z} \frac{\partial z}{\partial w_2} = x_2 \cdot dz \\ db &= \frac{\partial \mathcal{L}}{\partial b} = \frac{\partial \mathcal{L}}{\partial z} \frac{\partial z}{\partial b} = dz \end{aligned}$$

## Logistic Regression (0 hidden layers)

$$x \rightarrow [w_1x_1 + w_2x_2 + b] \xrightarrow{z} [\sigma(z)] \xrightarrow{\hat{y}=a} [\mathcal{L}(a, y)]$$

### Parameters update

$$\begin{aligned} w_1 &:= w_1 - \alpha dw_1 \\ w_2 &:= w_2 - \alpha dw_2 \\ b &:= b - \alpha db \end{aligned}$$

where  $\alpha$  is the learning rate

## Logistic Regression (0 hidden layers)

### Computing the gradient descent on $m$ training examples

- For one training example  $i$  :  
 $\hat{y}^{(i)} = a^{(i)} = \sigma(z^{(i)}) = \sigma(\underline{w}^T \underline{x}^{(i)} + b)$
- Cost (sum of the Loss  $\mathcal{L}$  over all training examples  $m$ )

$$J(w, b) = \frac{1}{m} \sum_{i=1}^m \mathcal{L}(\hat{y}^{(i)}, y^{(i)})$$

- Minimizing the Cost  $\Rightarrow$  derivative of the Cost w.r.t. the parameters

$$\begin{aligned} \frac{\partial J(w, b)}{\partial w_1} &= \frac{1}{m} \sum_{i=1}^m \frac{\partial \mathcal{L}(\hat{y}^{(i)}, y^{(i)})}{\partial w_1} \\ &= \frac{1}{m} \sum_{i=1}^m dw_1^{(i)} \end{aligned}$$

- $\frac{\partial J(w, b)}{\partial w_1}$  is the average over the  $m$  training examples of the gradients  $dw_1^{(i)}$

## Logistic Regression (0 hidden layers)

Equivalence between minimizing the Cost and maximum likelihood

- We interpret  $\hat{y}$  as  $p(y = 1|x)$   
 If  $y = 1 \rightarrow p(y|x) = \hat{y}$   
 If  $y = 0 \rightarrow p(y|x) = 1 - \hat{y}$   

$$p(y|x) = \hat{y}^y \cdot (1 - \hat{y})^{(1-y)}$$

|            |   |
|------------|---|
| If $y = 1$ | $= \hat{y}^1 \cdot (1 - \hat{y})^0 = \hat{y}$       |
| If $y = 0$ | $= \hat{y}^0 \cdot (1 - \hat{y})^1 = (1 - \hat{y})$ |

- We want to maximize  $\log p(y|x) = \log(\hat{y}^y \cdot (1 - \hat{y})^{(1-y)})$ 

$$= y \log(\hat{y}) + (1 - y) \log(1 - \hat{y})$$

$$= -\mathcal{L}(\hat{y}, y)$$
  - $\Rightarrow$  Maximizing  $\log p(y|x)$  is equivalent to minimizing  $\mathcal{L}(\hat{y}, y)$

$$\begin{aligned} p(\text{labels}) &= \prod_{i=1}^m p(y^{(i)} | x^{(i)}) \\ \log p(\text{labels}) &= \sum_{i=1}^m \log p(y^{(i)} | x^{(i)}) \\ &= \sum_{i=1}^m -\mathcal{L}(y^{(i)}, y^{(i)}) \end{aligned}$$

- On the whole training set
$$p(\text{labels}) = \prod_{i=1}^m p(y^{(i)} | x^{(i)})$$

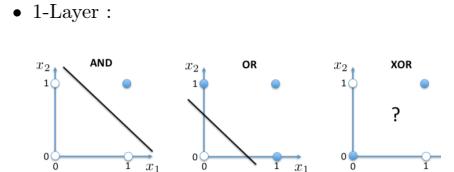
$$\log p(\text{labels}) = \sum_{i=1}^m \log p(y^{(i)} | x^{(i)})$$

$$= \sum_{i=1}^m -\mathcal{L}(y^{(i)}, y^{(i)})$$

## Logistic Regression (0 hidden layers)

## Limitation of linear classifiers

- Perceptron and Logistic Regression are linear classifiers
    - What if classes are not linearly separable?
    - What about the XOR function?
      - Need 2 layers
        - $x_1^{[1]} = AND(x_1^{[0]}, x_2^{[0]})$
        - $x_2^{[1]} = OR(x_1^{[0]}, x_2^{[0]})$



- The diagram shows the internal structure of an AND gate. It consists of two XOR gates connected in series. The first XOR gate takes inputs  $x_1$  and  $x_2$ , with output  $y_1 = \text{OR}(x_1, x_2)$ . This output  $y_1$  and the second input  $x_2$  are fed into a second XOR gate, which produces the final output  $y = \text{OR}(y_1, x_2) = \text{AND}(x_1, x_2)$ . The question mark indicates that the intermediate output  $y_1$  is not used directly.

Logistic regression, Softmax regression, Multi-label, Multi-Class

- Logistic regression (binary classification)

- **Logistic Regression** does binary classification (not regression)

$$P(y=1|\underline{x}) = h_w(x) = Logistic(\underline{w} \cdot \underline{x}) = \underbrace{\frac{1}{1+e^{-\underline{w} \cdot \underline{x}}}}_{\text{sigmoid function}} \quad (2)$$

$$P(y = 0 | \underline{x}) = 1 - P(y = 1 | \underline{x}) \quad (3)$$

- We can show that it models the posterior probability of the class  $y = 1$  using linear functions of the inputs  $x$

- ### • Proof

$$P(y = 1 | \underline{x}) = \frac{1}{1 + e^{-f(\underline{x})}} = \frac{e^{f(\underline{x})}}{1 + e^{f(\underline{x})}} \quad (4)$$

$$\frac{1}{P(y=1|x)} = 1 + e^{-f(x)} \quad (5)$$

$$\frac{1 - P(y = 1|x)}{P(y = 1|x)} = e^{-f(x)} \quad (6)$$

$$\log \left( \frac{P(y = 1 | \underline{x})}{P(y = 0 | \underline{x})} \right) = f(\underline{x}) \quad (7)$$

$$\log \left( \frac{P(y=1|x)}{\hat{P}(y=0|x)} \right) = w \cdot x \quad (8)$$

## Softmax regression (multi-class classification)

- **Softmax Regression :**

- Model the posterior probabilities of the K classes using linear functions of the inputs  $\underline{x}$ , according to :

$$\log \left( \frac{P(y = c_1 | \underline{x})}{P(y = c_k | \underline{x})} \right) = \underline{w}_1 \underline{x} + w_{10} \quad (9)$$

$$\log \left( \frac{P(y = c_2 | \underline{x})}{P(y = c_k | \underline{x})} \right) = \underline{w}_2 \underline{x} + w_{20} \quad (10)$$

$$\dots \quad (11)$$

$$\log \left( \frac{P(y = c_k - 1 | \underline{x})}{P(y = c_k | \underline{x})} \right) = \underline{w}_{k-1} \underline{x} + w_{(k-1)0} \quad (12)$$

$$(13)$$

- It is easy to deduce that

$$P(c_k | \underline{x}) = \frac{e^{\underline{w}_k \cdot \underline{x} + w_{k0}}}{1 + \sum_{l=1}^{K-1} e^{\underline{w}_l \cdot \underline{x} + w_{l0}}} \quad k = 1 \dots K-1 \quad (14)$$

$$P(c_K | \underline{x}) = \frac{1}{1 + \sum_{l=1}^{K-1} e^{\underline{w}_l \cdot \underline{x} + w_{l0}}} \quad (15)$$

- We of course have

$$P(c_1 | \underline{x}) + P(c_2 | \underline{x}) + \dots + P(c_K | \underline{x}) = 1 \quad (16)$$

- **Softmax function** (generalization of logistic function)

$$P(y = j | \underline{x}) = \frac{e^{\underline{w}_j \cdot \underline{x}}}{\sum_{k=1}^K e^{\underline{w}_k \cdot \underline{x}}} \quad (17)$$

using  $\underline{w}_K = 0$

## Chain rule and Back-propagation

- **Chain rule**

- formula for computing the derivative of the composition of two or more functions

$$\begin{aligned} \frac{d}{dt} f(x(t)) &= \frac{df}{dx} \frac{dx}{dt} \\ \frac{d}{dt} f(x(t), y(t)) &= \frac{\partial f}{\partial x} \frac{dx}{dt} + \frac{\partial f}{\partial y} \frac{dy}{dt} \end{aligned}$$

- Example

$$z_1 = z_1(x_1, x_2)$$

$$z_2 = z_2(x_1, x_2)$$

$$p = p(z_1, z_2)$$

$$\frac{\partial p}{\partial x_1} = \frac{\partial p}{\partial z_1} \frac{\partial z_1}{\partial x_1} + \frac{\partial p}{\partial z_2} \frac{\partial z_2}{\partial x_1}$$

- Example :

$$h(x) = f(x)g(x)$$

$$\frac{\partial h}{\partial x} = \frac{\partial h}{\partial f} \frac{\partial f}{\partial x} + \frac{\partial h}{\partial g} \frac{\partial g}{\partial x} = f'g + fg'$$

- **Back-propagation**

- efficient algorithm to compute the chain-rule by storing intermediate (and re-use derivatives)

## Chain rule and Back-propagation

## Chain rule and Back-propagation

- **Univariate logistic least squares model**

$$z = wx + b$$

$$\hat{y} = \sigma(z)$$

$$\mathcal{L} = \frac{1}{2}(\hat{y} - y)^2$$

- **Computing derivative as in calculus class**

$$\begin{aligned} \mathcal{L} &= \frac{1}{2} (\sigma(wx + b) - y)^2 \\ \frac{\partial \mathcal{L}}{\partial w} &= \frac{\partial}{\partial w} \left[ \frac{1}{2} (\sigma(wx + b) - y)^2 \right] \\ &= \frac{1}{2} \frac{\partial}{\partial w} (\sigma(wx + b) - y)^2 \\ &= (\sigma(wx + b) - y) \frac{\partial}{\partial w} (\sigma(wx + b) - y) \\ &= (\sigma(wx + b) - y) \sigma'(wx + b) \frac{\partial}{\partial w} (wx + b) \\ &= (\sigma(wx + b) - y) \sigma'(wx + b) x \end{aligned}$$

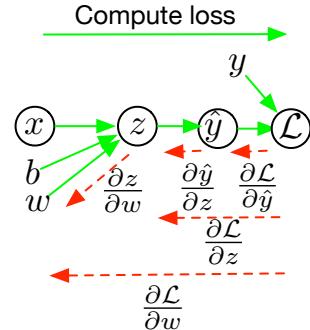
## Chain rule and Back-propagation

- Univariate logistic least squares model

$$z = wx + b$$

$$\hat{y} = \sigma(z)$$

$$\mathcal{L} = \frac{1}{2}(\hat{y} - y)^2$$



- Computing derivative using back-propagation

$$\frac{\partial \mathcal{L}}{\partial \hat{y}} = \hat{y} - y$$

$$\frac{\partial \mathcal{L}}{\partial z} = \frac{\partial \mathcal{L}}{\partial \hat{y}} \sigma'(z)$$

$$\frac{\partial \mathcal{L}}{\partial w} = \frac{\partial \mathcal{L}}{\partial z} x$$

$$\frac{\partial \mathcal{L}}{\partial b} = \frac{\partial \mathcal{L}}{\partial z}$$

- We can diagram out the computations using a **computation graph**

- nodes represent all the inputs and computed quantities,
- edges represent which nodes are computed directly as a function of which other nodes

## Chain rule and Back-propagation

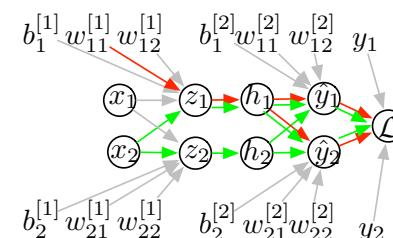
- Multilayer Perceptron (multiple outputs)

$$z_i = \sum_j w_{ij}^{[1]} x_j + b_i^{[1]}$$

$$h_i = \sigma(z_i)$$

$$\hat{y}_k = \sum_i w_{ki}^{[2]} h_i + b_k^{[2]}$$

$$\mathcal{L} = \frac{1}{2} \sum_k (\hat{y}_k - y_k)^2$$



- Computing derivative using back-propagation

- How much changing  $w_{11}/x_2$  affect  $\mathcal{L}$
- need to take into account all possible paths

$$\frac{\partial \mathcal{L}}{\partial \hat{y}_k} = \hat{y}_k - y_k$$

$$\frac{\partial \mathcal{L}}{\partial w_{ki}^{[2]}} = \frac{\partial \mathcal{L}}{\partial \hat{y}_k} h_i \quad \frac{\partial \mathcal{L}}{\partial b_k^{[2]}} = \frac{\partial \mathcal{L}}{\partial \hat{y}_k}$$

$$\frac{\partial \mathcal{L}}{\partial h_i} = \sum_k \frac{\partial \mathcal{L}}{\partial \hat{y}_k} w_{ki}^{[2]}$$

$$\frac{\partial \mathcal{L}}{\partial z_i} = \frac{\partial \mathcal{L}}{\partial h_i} \sigma'(z_i)$$

$$\frac{\partial \mathcal{L}}{\partial w_{ij}^{[1]}} = \frac{\partial \mathcal{L}}{\partial z_i} x_j \quad \frac{\partial \mathcal{L}}{\partial b_i^{[1]}} = \frac{\partial \mathcal{L}}{\partial z_i}$$

## Chain rule and Back-propagation

$$z_1 = z_1(x_1, x_2)$$

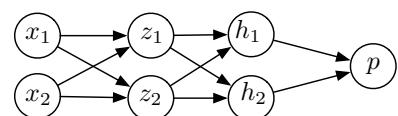
$$z_2 = z_2(x_1, x_2)$$

$$h_1 = z_1(z_1, z_2)$$

$$h_2 = z_2(z_1, z_2)$$

$$p = p(h_1, h_2)$$

$$\begin{aligned} \frac{\partial p}{\partial x_1} &= \frac{\partial p}{\partial h_1} \underbrace{\frac{\partial h_1}{\partial z_1}}_{\frac{\partial h_1}{\partial z_1} + \frac{\partial h_1}{\partial z_2}} + \underbrace{\frac{\partial p}{\partial h_2}}_{\frac{\partial h_2}{\partial z_1}} \underbrace{\frac{\partial z_1}{\partial x_1}}_{\frac{\partial h_1}{\partial z_1} + \frac{\partial h_1}{\partial z_2}} \\ &\quad + \frac{\partial p}{\partial h_1} \frac{\partial h_1}{\partial z_1} \frac{\partial z_1}{\partial x_1} + \frac{\partial p}{\partial h_1} \frac{\partial h_1}{\partial z_2} \frac{\partial z_2}{\partial x_1} + \frac{\partial p}{\partial h_2} \frac{\partial h_2}{\partial z_1} \frac{\partial z_1}{\partial x_1} + \frac{\partial p}{\partial h_2} \frac{\partial h_2}{\partial z_2} \frac{\partial z_2}{\partial x_1} \end{aligned}$$



- Each edge is assigned to derivative of origin w.r.t. destination

## Neural Networks (1 hidden layer)

## Neural Networks (1 hidden layer)

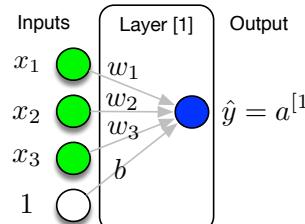
- Logistic Regression** (1 layer, **0 hidden layer**)

$$x \rightarrow [w_1 x_1 + w_2 x_2 + b] \xrightarrow{z} [\sigma(z)] \xrightarrow{\hat{y}=a} [\mathcal{L}(a, y)]$$

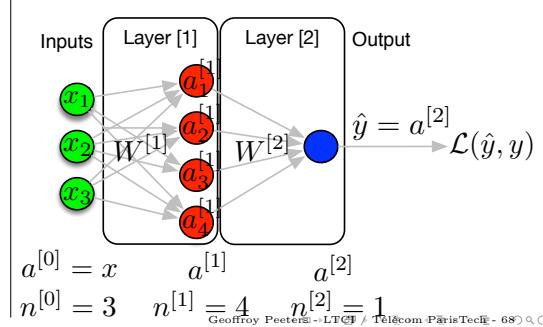
- Neural Network** (2 layers, **1 hidden layer**)

$$\underline{x} \xrightarrow{a^{[0]}} \underbrace{[\underline{W}^{[1]} \underline{x} + \underline{b}^{[1]}]}_{Layer[1]} \xrightarrow{\underline{z}^{[1]}} \underbrace{[g^{[1]}(\underline{z}^{[1]})]}_{Layer[2]} \xrightarrow{a^{[1]}} \underbrace{[\underline{W}^{[2]} \underline{a}^{[1]} + \underline{b}^{[2]}]}_{Layer[2]} \xrightarrow{z^{[2]}} \underbrace{[g^{[2]}(z^{[2]})]}_{Layer[2]} \xrightarrow{a^{[2]}} \mathcal{L}(a^{[2]}, y)$$

### Logistic Regression



### Neural Network

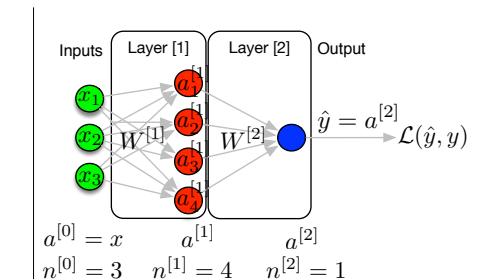


## Neural Networks (1 hidden layer)

$$\underline{x} \xrightarrow{a^{[0]}} \underbrace{[\underline{W}^{[1]} \underline{x} + \underline{b}^{[1]}]}_{Layer[1]} \xrightarrow{\underline{z}^{[1]}} \underbrace{[g^{[1]}(\underline{z}^{[1]})]}_{Layer[2]} \xrightarrow{a^{[1]}} \underbrace{[\underline{W}^{[2]} \underline{a}^{[1]} + \underline{b}^{[2]}]}_{Layer[2]} \xrightarrow{z^{[2]}} \underbrace{[g^{[2]}(z^{[2]})]}_{Layer[2]} \xrightarrow{a^{[2]}} \mathcal{L}(a^{[2]}, y)$$

### Notations :

- variables of **layer**  $l \Rightarrow *^{[l]}$
- dimension**  $d$  of variables  $\Rightarrow *_d^{[l]}$
- training **example**  $i \Rightarrow *_d^{[l](i)}$
- input  $x \Rightarrow a^{[0]}$
- dimensions of layer  $l \Rightarrow n^{[l]}$
- $g^{[l]}$  : activation function
  - (possibly a sigmoid  $\sigma$  function)



Geoffroy Peeters - LTCI / Télécom ParisTech - 68

## Neural Networks (1 hidden layer)

$$\underline{x} \xrightarrow{a^{[0]}} \underbrace{[\underline{W}^{[1]} \underline{x} + \underline{b}^{[1]}]}_{Layer[1]} \xrightarrow{\underline{z}^{[1]}} \underbrace{[g^{[1]}(\underline{z}^{[1]})]}_{Layer[2]} \xrightarrow{a^{[1]}} \underbrace{[\underline{W}^{[2]} \underline{a}^{[1]} + \underline{b}^{[2]}]}_{Layer[2]} \xrightarrow{z^{[2]}} \underbrace{[g^{[2]}(z^{[2]})]}_{Layer[2]} \xrightarrow{a^{[2]}} \mathcal{L}(a^{[2]}, y)$$

### Gradient Descent

- Parameters** : to update

$$\underbrace{\underline{W}^{[1]}}_{(n^{[1]}, n^{[0]})}, \underbrace{\underline{b}^{[1]}}_{(n^{[1]}, 1)}, \underbrace{\underline{W}^{[2]}}_{(n^{[2]}, n^{[1]})}, \underbrace{\underline{b}^{[2]}}_{(n^{[2]}, 1)}$$

- Gradient Descent** :

- Initialize the parameters
- Repeat
  - Forward propagation** : compute prediction  $\hat{y}^{(i)} \quad \forall i = 1, \dots, m$
  - Compute the loss**
    - Cost function (binary classification) :
    - $J(w, b) = \frac{1}{m} \sum_{i=1}^m \mathcal{L}(\hat{y}^{(i)}, y^{(i)})$
  - Backward propagation** : compute gradients  $dW^{[l]}, db^{[l]}$
  - Update the parameters** using the learning rate  $\alpha$

## Neural Networks (1 hidden layer)

$$\underline{x} \xrightarrow{a^{[0]}} \underbrace{[\underline{W}^{[1]} \underline{x} + \underline{b}^{[1]}]}_{Layer[1]} \xrightarrow{\underline{z}^{[1]}} \underbrace{[g^{[1]}(\underline{z}^{[1]})]}_{Layer[2]} \xrightarrow{a^{[1]}} \underbrace{[\underline{W}^{[2]} \underline{a}^{[1]} + \underline{b}^{[2]}]}_{Layer[2]} \xrightarrow{z^{[2]}} \underbrace{[g^{[2]}(z^{[2]})]}_{Layer[2]} \xrightarrow{a^{[2]}} \mathcal{L}(a^{[2]}, y)$$

### Forward propagation (all dimensions, all $m$ training examples)

$$\begin{aligned} \underbrace{\underline{Z}^{[1]}}_{(n^{[1]}, m)} &= \underbrace{\underline{W}^{[1]}}_{(n^{[1]}, n^{[0]})} \underbrace{\underline{X}}_{(n^{[0]}, m)} + \underbrace{\underline{b}^{[1]}}_{(n^{[1]}, 1)} \\ \underbrace{\underline{A}^{[1]}}_{(n^{[1]}, m)} &= g^{[1]} \left( \underline{Z}^{[1]} \right) \\ \underbrace{\underline{Z}^{[2]}}_{(n^{[2]}, m)} &= \underbrace{\underline{W}^{[2]}}_{(n^{[2]}, n^{[1]})} \underbrace{\underline{A}^{[1]}}_{(n^{[1]}, m)} + \underbrace{\underline{b}^{[2]}}_{(n^{[2]}, 1)} \\ \underbrace{\underline{A}^{[2]}}_{(n^{[2]}, m)} &= g^{[2]} \left( \underline{Z}^{[2]} \right) \end{aligned}$$

Geoffroy Peeters - LTCI / Télécom ParisTech - 74

## Neural Networks (1 hidden layer)

### Backward propagation (each training example)

$$da^{[2]} = \frac{\partial \mathcal{L}}{\partial a^{[2]}} = \text{derivative of the loss} = - \left( \frac{y}{a^{[2]}} + \frac{(1-y)}{(1-a^{[2]})} \right)$$

$$dz^{[2]} = \frac{\partial \mathcal{L}}{\partial z^{[2]}} = \frac{\partial \mathcal{L}}{\partial a^{[2]}} \frac{\partial a^{[2]}}{\partial z^{[2]}} = da^{[2]} \odot g^{[2]}'(z^{[2]}) = da^{[2]} \odot a^{[2]}(1-a^{[2]}) = a^{[2]} - y$$

$$dW^{[2]} = \frac{\partial \mathcal{L}}{\partial W^{[2]}} = \frac{\partial \mathcal{L}}{\partial z^{[2]}} \frac{\partial z^{[2]}}{\partial W^{[2]}} = \underbrace{dz^{[2]}}_{(n^{[2]},1)} \underbrace{a^{[1]}}_{(n^{[1]},1)}^T$$

$$db^{[2]} = \frac{\partial \mathcal{L}}{\partial b^{[2]}} = \frac{\partial \mathcal{L}}{\partial z^{[2]}} \frac{\partial z^{[2]}}{\partial b^{[2]}} = dz^{[2]}$$

$$da^{[1]} = \frac{\partial \mathcal{L}}{\partial a^{[1]}} = \frac{\partial \mathcal{L}}{\partial z^{[2]}} \frac{\partial z^{[2]}}{\partial a^{[1]}} = \underbrace{W^{[2]}}_{(n^{[2]}, n^{[1]})}^T \underbrace{dz^{[2]}}_{(n^{[2]},1)}$$

$$dz^{[1]} = \frac{\partial \mathcal{L}}{\partial z^{[1]}} = \frac{\partial \mathcal{L}}{\partial a^{[1]}} \frac{\partial a^{[1]}}{\partial z^{[1]}} = da^{[1]} \odot g^{[1]}'(z^{[1]})$$

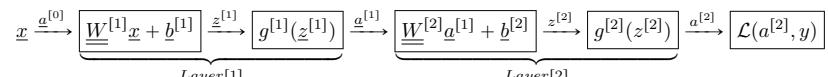
$$dW^{[1]} = \frac{\partial \mathcal{L}}{\partial W^{[1]}} = \frac{\partial \mathcal{L}}{\partial z^{[1]}} \frac{\partial z^{[1]}}{\partial W^{[1]}} = \underbrace{dz^{[1]}}_{(n^{[1]},1)} \underbrace{a^{[0]}}_{(n^{[0]},1)}^T$$

$$db^{[1]} = \frac{\partial \mathcal{L}}{\partial b^{[1]}} = \frac{\partial \mathcal{L}}{\partial z^{[1]}} \frac{\partial z^{[1]}}{\partial b^{[1]}} = dz^{[1]}$$

where  $\odot$  is the element-wise (also named Hadamard) product.

Geoffroy Peeters - LTCI / Télécom ParisTech - 75) Q. C.

## Neural Networks (1 hidden layer)



### Parameters update

- Parameters update :**

$$W^{[l]} = W^{[l]} - \alpha \cdot dW^{[l]}$$

$$b^{[l]} = b^{[l]} - \alpha \cdot db^{[l]}$$

where  $\alpha$  is the learning rate

## Neural Networks (1 hidden layer)

### Backward propagation (all $m$ training examples)

$$\underbrace{dZ^{[2]}}_{(n^{[2]},m)} = \underbrace{\underline{A}^{[2]}}_{(n^{[2]},m)} - \underbrace{\underline{Y}}_{(n^{[2]},m)}$$

$$\underbrace{dW^{[2]}}_{(n^{[2]},n^{[1]})} = \frac{1}{m} \underbrace{dZ^{[2]}}_{(n^{[2]},m)} \underbrace{\underline{A}^{[1]}}_{(n^{[1]},m)}^T$$

$$\underbrace{db^{[2]}}_{(n^{[2]},1)} = \frac{1}{m} \sum_{i=1}^m \underbrace{d\underline{Z}^{[2]}}_{(n^{[2]},m)}$$

$$\underbrace{d\underline{A}^{[1]}}_{(n^{[1]},m)} = \underbrace{W^{[2]}}_{(n^{[2]},n^{[1]})}^T \underbrace{d\underline{Z}^{[2]}}_{(n^{[2]},m)}$$

$$\underbrace{d\underline{Z}^{[1]}}_{(n^{[1]},m)} = \underbrace{d\underline{A}^{[1]}}_{(n^{[1]},m)} \odot g^{[1]}'(\underbrace{\underline{Z}^{[1]}}_{(n^{[1]},m)})$$

$$\underbrace{dW^{[1]}}_{(n^{[1]},n^{[0]})} = \frac{1}{m} \underbrace{d\underline{Z}^{[1]}}_{(n^{[1]},m)} \underbrace{\underline{A}^{[0]}}_{(n^{[0]},m)}^T$$

$$\underbrace{db^{[1]}}_{(n^{[1]},1)} = \frac{1}{m} \sum_{i=1}^m \underbrace{d\underline{Z}^{[1]}}_{(n^{[1]},m)}$$

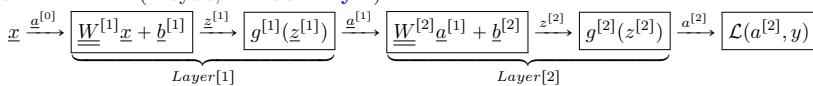
Geoffroy Peeters - LTCI / Télécom ParisTech - 76) Q. C.

## Deep Neural Networks ( $> 2$ hidden layers)

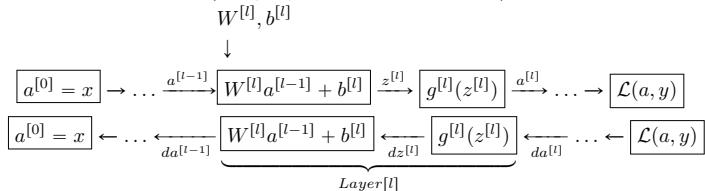
Geoffroy Peeters - LTCI / Télécom ParisTech - 77) Q. C.

Deep Neural Networks ( $> 2$  hidden layers)

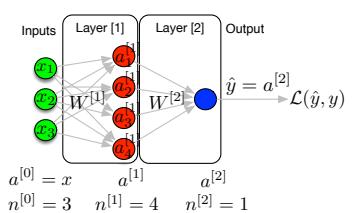
- Neural Network (2 layers, 1 hidden layer)



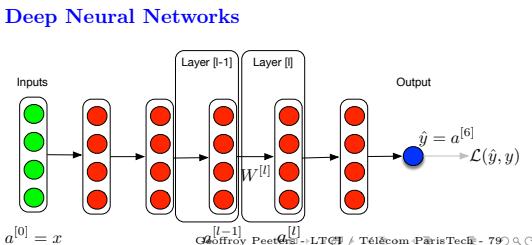
- Deep Neural Networks (many layers, **>2** hidden layer)



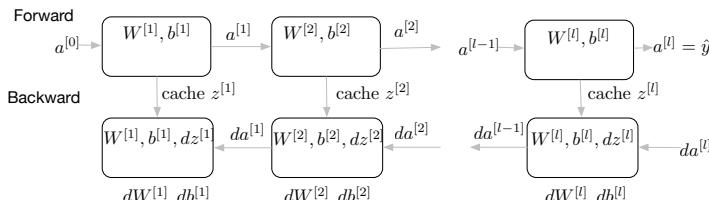
## Neural Network



## Deep Neural Networks



Deep Neural Networks (> 2 hidden layers)



Backward (general formulation for layer  $l$ , all  $m$  training examples)

Input :  $da^{[l]}$

Cache (passed from Forward) :  $z^{[l]}$

$$dZ^{[l]} = dA^{[l]} \odot {g^{[l]}}'(Z^{[l]})$$

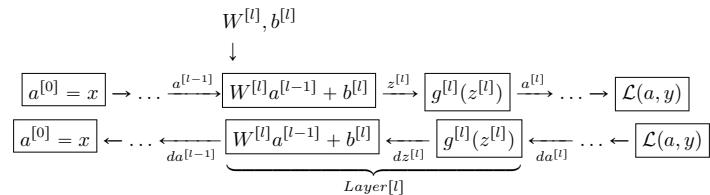
$$d\hat{W}^{[l]} = \frac{1}{m} d\hat{Z}^{[l]} \cdot \hat{A}^{[l-1]}{}^T$$

$$db^{[l]} = \frac{1}{m} \sum^m dZ^{[l]}$$

$$dA^{[l-1]} = W^{[l]}^T \cdot dZ^{[l]}$$

**Output :**  $da^{[l-1]}$   $dW^{[l]}$   $db^{[l]}$

Deep Neural Networks (> 2 hidden layers)



Forward (general formulation for layer  $l$ , all  $m$  training examples)

**Input :**  $A^{[l-1]}$

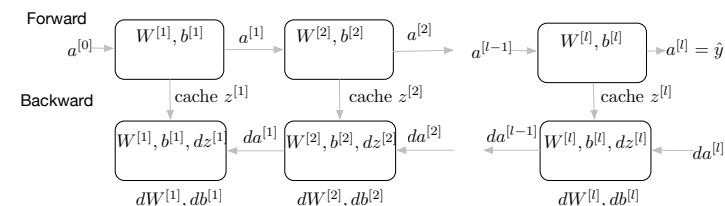
$$\underbrace{\underline{Z}^{[l]}}_{(n^{[l]}, m)} = \underbrace{\underline{W}^{[l]}}_{(n^{[l]}, n^{[l-1]})} \underbrace{\underline{A}^{[l-1]}}_{(n^{[l-1]}, m)} + \underbrace{\underline{b}^{[l]}}_{(n^{[l]}, 1)}$$

$$\underbrace{\underline{A}^{[l]}}_{(n^{[l]}, m)} = g^{[l]}(\underline{Z}^{[l]})$$

**Output :**  $\underline{A}^{[l]}$ , cache( $\underline{Z}^{[l]}$ ,  $\underline{W}^{[l]}$ ,  $\underline{b}^{[l]}$ )

Geoffroy Peeters - LTCI / Télécom ParisTech - 83

Deep Neural Networks (> 2 hidden layers)



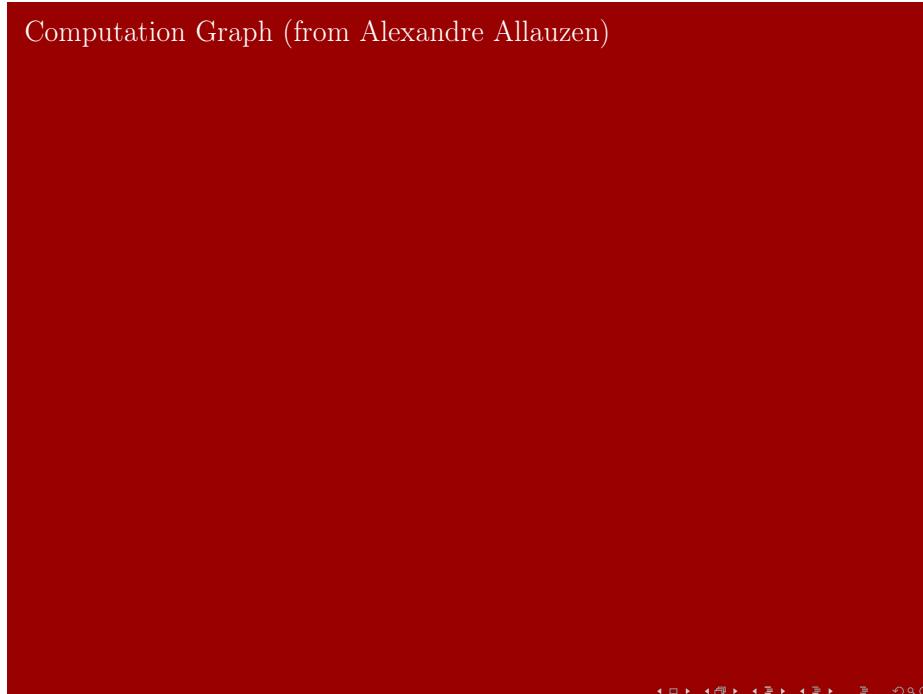
## Parameters update

$$\begin{aligned}\underline{\underline{W}}^{[l]} &= \underline{\underline{W}}^{[l]} - \alpha d\underline{\underline{W}}^{[l]} \\ \underline{b}^{[l]} &\equiv \underline{b}^{[l]} - \alpha db^{[l]}\end{aligned}$$

where  $\alpha$  is the learning rate

Geoffroy Peeters - LTCI / Télécom ParisTech - 860 & C

## Computation Graph (from Alexandre Allauzen)

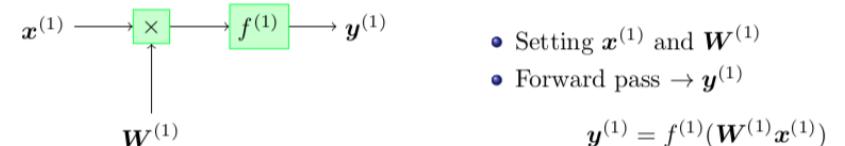


## Computation graph

A convenient way to represent a complex mathematical expressions :

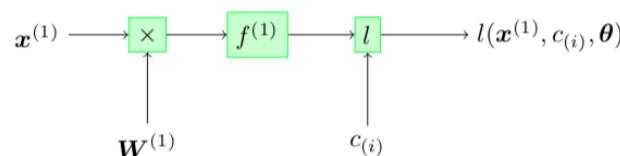
- each node is an operation or a variable
- an operation has some inputs / outputs made of variables

### Example 1 : A single layer network



Geoffroy Peeters - LTCI / Télécom ParisTech - 88/90

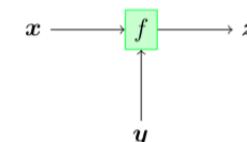
## Computation graph for training



- A variable node encodes the label
- To compute the output for a given input  
→ forward pass
- To compute the gradient of the loss wrt the parameters ( $\mathbf{W}^{(1)}$ )  
→ backward pass

## A function node

### Forward pass



This node implements :

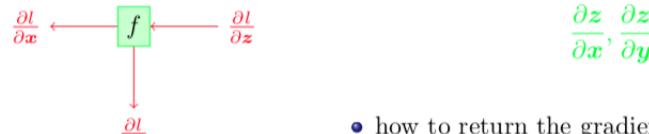
$$z = f(x, y)$$

## A function node - 2

### Backward pass

A function node knows :

- the "local gradients" computation



- how to return the gradient to the inputs :

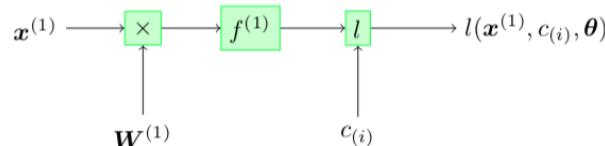
$$\left( \frac{\partial l}{\partial \mathbf{z}} \frac{\partial \mathbf{z}}{\partial \mathbf{x}}, \left( \frac{\partial l}{\partial \mathbf{z}} \frac{\partial \mathbf{z}}{\partial \mathbf{y}} \right) \right)$$

## Summary of a function node

$f :$

$$\begin{aligned} \mathbf{x}, \mathbf{y}, \mathbf{z} & \quad \# \text{ store the values} \\ \mathbf{z} = f(\mathbf{x}, \mathbf{y}) & \quad \# \text{ forward} \\ \frac{\partial \mathbf{z}}{\partial \mathbf{x}} \rightarrow \frac{\partial f}{\partial \mathbf{x}} & \quad \# \text{ local gradients} \\ \frac{\partial \mathbf{z}}{\partial \mathbf{y}} \rightarrow \frac{\partial f}{\partial \mathbf{y}} & \\ \left( \frac{\partial l}{\partial \mathbf{z}} \frac{\partial \mathbf{z}}{\partial \mathbf{x}}, \left( \frac{\partial l}{\partial \mathbf{z}} \frac{\partial \mathbf{z}}{\partial \mathbf{y}} \right) \right) & \quad \# \text{ backward} \end{aligned}$$

## Example of a single layer network



### Forward

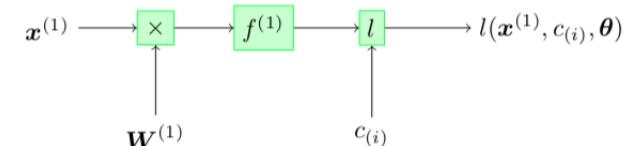
For each function node in topological order

- forward propagation

Which means :

- ➊  $\mathbf{a}^{(1)} = \mathbf{W}^{(1)} \mathbf{x}^{(1)}$
- ➋  $\mathbf{y}^{(1)} = f^{(1)}(\mathbf{a}^{(1)})$
- ➌  $l(\mathbf{y}^{(1)}, c_{(i)})$

## Example of a single layer network



### Backward

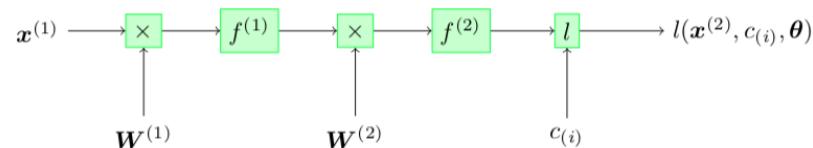
For each function node in reversed topological order

- backward propagation

Which means :

- ➊  $\nabla_{\mathbf{y}^{(1)}}$
- ➋  $\nabla_{\mathbf{a}^{(1)}}$
- ➌  $\nabla_{\mathbf{W}^{(1)}}$

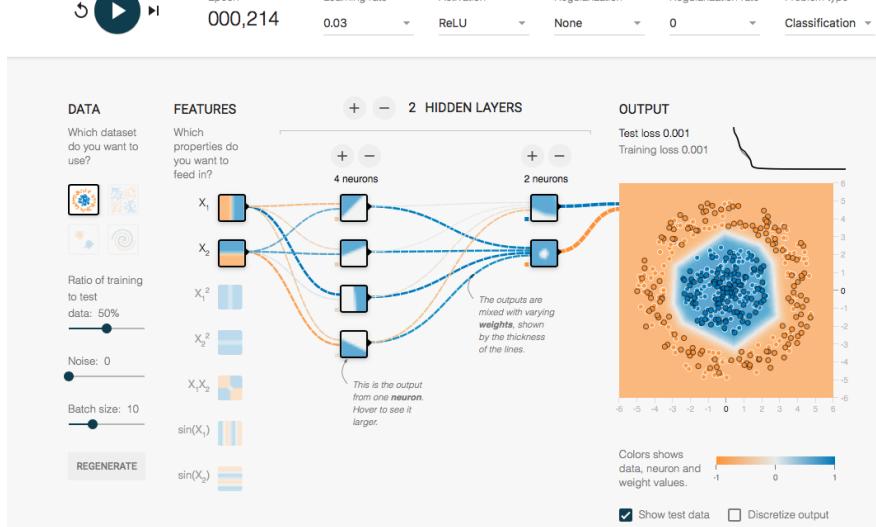
## Example of a two layers network



- The algorithms remain the same,
- even for more complex architectures
- Generalization by coding your own function node or by
- Wrapping a layer in a module

## Deep Learning Frameworks

## Playground



## Deep Learning Frameworks

- **Python**
  - coding forward, loss, backward, upgrading
    - painfull, prone to coding errors
- **(py)Torch**
  - The Facebook library with Lua/python API
  - <https://pytorch.org>
- **Tensorflow**
  - The Google library with python API
  - <https://www.tensorflow.org>
- **Keras**
  - A high-level API, in Python, running on top of TensorFlow
  - <https://keras.io>
- CPU/GPU
- Automatic differentiation based on computational graph

## Example of an MLP in python

```

def nn_model(X, Y, n_h, num_iterations=10000):
    for i in range(0, num_iterations):
        A2, cache = forward_propagation(X, parameters)
        cost = compute_cost(A2, Y, parameters)
        grads = backward_propagation(parameters, cache, X, Y)
        parameters = update_parameters(parameters, grads)
    return parameters

def forward_propagation(X, parameters):
    Z1 = np.dot(W1, X) + b1
    A1 = np.tanh(Z1)
    Z2 = np.dot(W2, A1) + b2
    A2 = sigmoid(Z2)
    return A2, cache

def compute_cost(A2, Y, parameters):
    logprobs = np.multiply(np.log(A2), Y) + np.multiply((1 - Y), np.log(1 - A2))
    cost = -np.sum(logprobs) / m
    return cost

def backward_propagation(parameters, cache, X, Y):
    dZ2 = A2 - Y
    dW2 = (1 / m) * np.dot(dZ2, A1.T)
    db2 = (1 / m) * np.sum(dZ2, axis=1, keepdims=True)
    dZ1 = np.multiply(np.dot(W2.T, dZ2), 1 - np.power(A1, 2))
    dW1 = (1 / m) * np.dot(dZ1, X.T)
    db1 = (1 / m) * np.sum(dZ1, axis=1, keepdims=True)
    return grads

def update_parameters(parameters, grads, learning_rate=1.2):
    W1 = W1 - learning_rate * dW1
    b1 = b1 - learning_rate * db1
    W2 = W2 - learning_rate * dW2
    b2 = b2 - learning_rate * db2
    return parameters

```

## Example of an MLP in pytorch

```

import torch
from torch.autograd import Variable
# N is batch size; D_in is input dimension;
# H is hidden dimension; D_out is output dimension. N, D_in, H, D_out = 64, 1000, 100, 10
# Create random Tensors to hold inputs and outputs,
# and wrap them in Variables.
x = torch.randn(N, D_in)
y = torch.randn(N, D_out)

# Use the nn package to define our model # as a sequence of layers.
model = torch.nn.Sequential(
    torch.nn.Linear(D_in, H),
    torch.nn.ReLU(),
    torch.nn.Linear(H, D_out),
)
loss_fn = torch.nn.MSELoss(size_average=False)
# Optimizer will update the weights of the model. lr0 = 1e-4
optimizer = torch.optim.Adam(model.parameters(), lr=learning_rate)
for t in range(10):
    # Forward pass: compute predicted y by passing x.
    y_pred = model(x)
    # Compute and print loss.
    loss = loss_fn(y_pred, y)
    print(t, loss.data[0])
    # Optim in two steps
    optimizer.zero_grad()
    # Backward pass: compute gradient of the loss wrt parameters
    loss.backward()
    # Calling the step function on an Optimizer makes an update
    optimizer.step()

```

## Example of an MLP in tensorflow

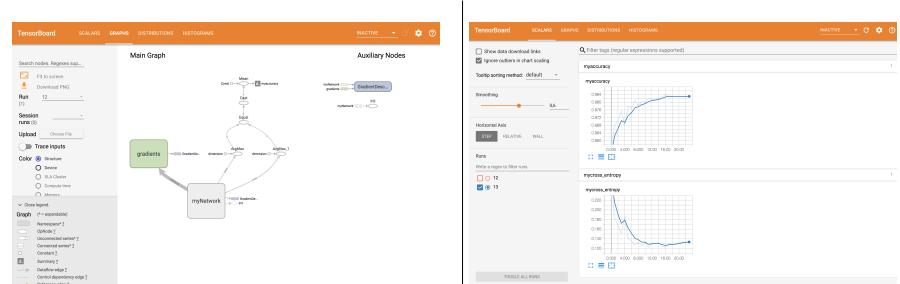
```

import tensorflow as tf

# declare the training data placeholders
# input x - for 28 x 28 pixels = 784
x = tf.placeholder(tf.float32, [None, D_in])
# now declare the output data placeholder - 10 digits
y = tf.placeholder(tf.float32, [None, D_out])
# now declare the weights connecting the input to the hidden layer
W1 = tf.Variable(tf.random_normal([D_in, H], stddev=0.03), name='W1')
b1 = tf.Variable(tf.random_normal([H]), name='b1')
# and the weights connecting the hidden layer to the output layer
W2 = tf.Variable(tf.random_normal([H, D_out], stddev=0.03), name='W2')
b2 = tf.Variable(tf.random_normal([D_out]), name='b2')
# calculate the output of the hidden layer
hidden_out = tf.add(tf.matmul(x, W1), b1)
hidden_out = tf.nn.relu(hidden_out)
# output layer: now calculate the hidden layer output - in this case, let's use a softmax activated
y_ = tf.nn.softmax(tf.add(tf.matmul(hidden_out, W2), b2))
y_clipped = tf.clip_by_value(y_, 1e-10, 0.9999999)
cross_entropy = -tf.reduce_mean(tf.reduce_sum(y * tf.log(y_clipped) + (1 - y) * tf.log(1 - y_clipped), axis=1))
# add an optimiser
optimiser = tf.train.GradientDescentOptimizer(learning_rate=learning_rate).minimize(cross_entropy)
# define an accuracy assessment operation
correct_prediction = tf.equal(tf.argmax(y, 1), tf.argmax(y_, 1))
accuracy = tf.reduce_mean(tf.cast(correct_prediction, tf.float32))
# finally setup the initialisation operator
init_op = tf.global_variables_initializer()
# start the session
with tf.Session() as sess:
    # initialise variables
    sess.run(init_op)
    total_batch = int(N / batch_size)
    for epoch in range(nbEpochs):
        avg_cost = 0
        for i in range(total_batch):
            batch_x, batch_y = x_train[i*batch_size:(i+1)*batch_size], y_train_ohv[i*batch_size:(i+1)*batch_size]
            _, c = sess.run([optimiser, cross_entropy], feed_dict={x: batch_x, y: batch_y})
            avg_cost += c / total_batch
        test_acc = sess.run(accuracy, feed_dict={x: x_test, y: y_test_ohv})
        print("Epoch:", (epoch + 1), "cost =", "{:.3f}".format(avg_cost), " test accuracy: {:.3f}".format(test_acc))
    print("\nTraining complete!")
    print(sess.run(accuracy, feed_dict={x: x_test, y: y_test_ohv}))

```

## Tensorboard



## Example of an MLP in Keras

```
import keras
from keras.models import Sequential
from keras.layers import Dense

# Convert labels to categorical one-hot encoding
y_train_ohv = keras.utils.to_categorical(y_train, num_classes=D_out)
y_test_ohv = keras.utils.to_categorical(Y_test, num_classes=D_out)

model = Sequential()
model.add(Dense(H, activation='relu', input_dim=D_in))
model.add(Dense(D_out, activation='softmax'))
model.compile(optimizer='rmsprop', loss='categorical_crossentropy', metrics=['accuracy'])
# Train the model, iterating on the data in batches of 32 samples
model.fit(x_train, y_train_ohv, epochs=nbEpochs, batch_size=batch_size, validation_data=(x_test, y_test_ohv))
score = model.evaluate(x_test, y_test_ohv, verbose=0)
print(score)
```

<https://colab.research.google.com/>

Bonjour Colaboratory

Colaboratory est un environnement de notebook Jupyter gratuit qui ne nécessite aucune configuration et qui s'exécute entièrement dans le cloud. Pour en savoir plus, consultez les [questions fréquentes](#).

Premiers pas

- Présentation de Colaboratory
- Chargement et sauvegarde des données : fichiers locaux, unité, feuilles, Google Cloud Storage
- Importation de bibliothèques et installation de dépendances
- Utilisation de Google Cloud BigQuery
- Formulaires, Gradiennes, Markdown et Widgets
- TensorFlow avec GPU
- Cours d'initiation au machine learning : [Introduction à Pandas](#) et [Premiers pas avec TensorFlow](#)

Caractéristiques principales

Exécution de TensorFlow

Colaboratory vous permet d'exécuter du code TensorFlow dans votre navigateur en un seul clic. L'exemple ci-dessous ajoute deux matrices.

$$\begin{bmatrix} 1. & 1. & 1. \end{bmatrix} + \begin{bmatrix} 1. & 2. & 3. \end{bmatrix} = \begin{bmatrix} 2. & 3. & 4. \end{bmatrix}$$

```
[ ] import tensorflow as tf
input1 = tf.ones((2, 3))
input2 = tf.reshape(tf.range(1, 7, dtype=tf.float32), (2, 3))
```

## Various types of training

## Various types of training

### Batch Gradient Descent

- Gradient computation**  $\frac{\partial J(w,b)}{\partial w_d}$ 
  - average of the gradient  $w_d^{(i)}$  over **all**  $m$  training examples
- Parameters update**
  - after the  $m$  training examples went to forward
- Pros and Cons**
  - (+) The gradient is accurate
  - (-) Can be very costly
    - if  $m = 5.000.000$ , we need to do  $m$  forward pass before doing any parameter update, can be very slow

## Various types of training

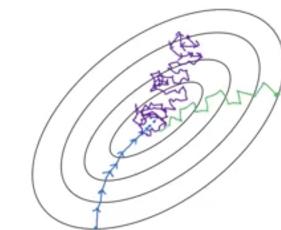
### Mini Batch Gradient Descent

- **Gradient computation**  $\frac{\partial J(w,b)}{\partial w_d}$ 
  - average of the gradient  $w_d^{(i)}$  over each **mini-batch**
- **Mini-batch ?**
  - Split training set into small training sets :  $X^{\{t\}}, Y^{\{t\}}$
  - Suppose mini-batch are of size 1000 :
$$X = [\underbrace{X^{(1)} \dots X^{(1000)}}_{(n_x, 1)} | \dots | \underbrace{X^{(2001)} \dots X^{(5.000.000)}}_{(n_x, 1000)}] = [\underbrace{X^{\{1\}}}_{(n_x, 1000)} \underbrace{X^{\{2\}}}_{(n_x, 1000)} \dots \underbrace{X^{\{5000\}}}_{(n_x, 1000)}]$$
- **Parameters update**
  - after each mini-batch
- **Pseudo-code**
  - for  $t=1, \dots, 5000 :$ 
    - Forward propagation on  $X^{\{t\}} : Z^{[1]} = W^{[1]}X^{[t]} + b^{[1]}, A^{[1]} = g^{[1]}(Z^{[1]}), \dots$
    - Compute cost :  $J^{[t]} = \frac{1}{1000} \sum_{i=1}^l \mathcal{L}(y^{(i)}, \hat{y}^{(i)}) + \frac{\lambda}{2 \cdot 1000} \sum_l \|W^{[l]}\|_F^2$
    - Backpropagation to compute gradients of  $J^{[t]}$  w.r.t the parameters
    - Update :  $W^{[l]} = W^{[l]} - \alpha dW^{[l]}, b^{[l]} = b^{[l]} - \alpha db^{[l]}$
  - One pass to the training set= one "**epoch**"
- **Pros and Cons**
  - (+) Allows to do 5000 gradient-descent steps
  - (-) Gradient is an approximation

## Various types of training

### Stochastic Gradient Descent (SGD)

- **Gradient computation**  $\frac{\partial J(w,b)}{\partial w_d}$ 
  - directly the gradient **on a single training examples**  $i$  (as mini-batch but with a mini-batch size of 1)
- **Parameters update**
  - after each training examples
- **Pros and Cons**
  - (+) Allows to do 5.000.000 gradient-descent steps
  - (-) Can lead to a very noisy gradient descent (see Figure)
  - (-) Lose speed up from vectorization



## Activation functions $a = g(z)$

## Activation functions $a = g(z)$

### Why non-linear activation functions?

- Consider the following network
$$\begin{aligned} \underline{z}^{[1]} &= \underline{W}^{[1]}\underline{x} + \underline{b}^{[1]} \\ \underline{a}^{[1]} &= g^{[1]}(\underline{z}^{[1]}) \\ \underline{z}^{[2]} &= \underline{W}^{[2]}\underline{a}^{[1]} + \underline{b}^{[2]} \\ \underline{a}^{[2]} &= g^{[2]}(\underline{z}^{[2]}) \end{aligned}$$
- If  $g^{[1]}$  and  $g^{[2]}$  are **linear activation functions** (identity function),
$$\begin{aligned} \underline{a}^{[1]} &= \underline{z}^{[1]} = \underline{W}^{[1]}\underline{x} + \underline{b}^{[1]} \\ \underline{a}^{[2]} &= \underline{z}^{[2]} = \underline{W}^{[2]}\underline{a}^{[1]} + \underline{b}^{[2]} \\ &= \underline{W}^{[2]} (\underline{W}^{[1]}\underline{x} + \underline{b}^{[1]}) + \underline{b}^{[2]} \\ &= \underline{W}^{[2]}\underline{W}^{[1]}\underline{x} + \underline{W}^{[2]}\underline{b}^{[1]} + \underline{b}^{[2]} \\ &= \underline{W}'\underline{x} + \underline{b}' \end{aligned}$$
  - then the network **the network reduces to a simple linear function**
- Linear activation? only interesting for regression problem :  $y \in \mathbb{R}$ 
  - linear function for last layer :  $g^{[L]}$

## Activation functions $a = g(z)$

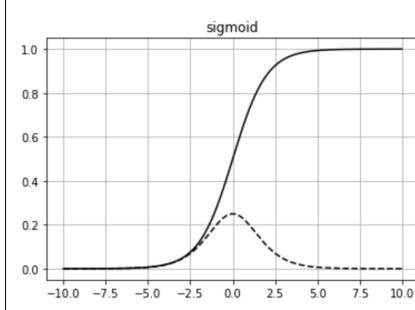
### Sigmoid $\sigma$

- Sigmoid function**

$$a = g(z) = \sigma(z) = \frac{1}{1 + e^{-z}}$$

- Derivative**

$$\begin{aligned} g'(z) &= -e^{-z} \frac{1}{(1 + e^{-z})^2} \\ &= \frac{1 + e^{-z} - 1}{(1 + e^{-z})^2} \\ &= \frac{1}{1 + e^{-z}} \left(1 - \frac{1}{1 + e^{-z}}\right) \\ &= \sigma(z)(1 - \sigma(z)) \\ g'(z) &= a(1 - a) \end{aligned}$$



## Activation functions $a = g(z)$

### Hyperbolic tangent $g$

- Hyperbolic tangent function**

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

- Derivative**

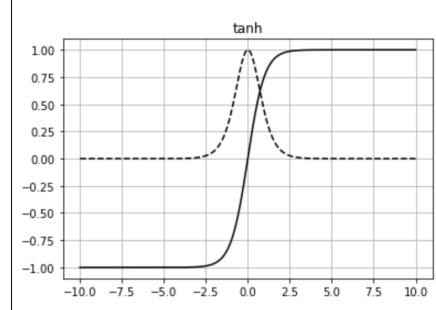
$$\begin{aligned} g'(x) &= 1 - (\tanh(z))^2 \\ g'(z) &= 1 - a^2 \end{aligned}$$

- Usage**

- $\tanh(z)$  better than  $\sigma(z)$  in middle hidden layers because its mean = zero ( $a \in [-1, 1]$ ).

- Problem** with  $\sigma$  and  $\tanh$  :

- if  $z$  is very small (negative) or very large (positive)
  - ⇒ slope becomes zero
  - ⇒ slow down Gradient Descent



## Activation functions $a = g(z)$

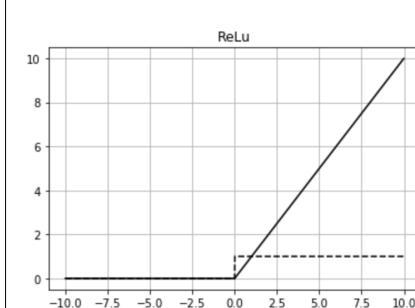
### ReLU (Rectified Linear Unit)

- ReLU function**

$$a = g(z) = \max(0, z)$$

- Derivative**

$$\begin{aligned} g'(x) &= 1 && \text{if } z > 0 \\ &= 0 && \text{if } z \leq 0 \end{aligned}$$



## Activation functions $a = g(z)$

### Leaky ReLU / Parametric ReLU

- Leaky ReLU function**

$$a = g(x) = \max(0.01z, z)$$

- allows avoiding the zero slope ("the neuron dies") for  $z < 0$

- Derivative**

$$\begin{aligned} g'(x) &= 1 && \text{if } z > 0 \\ &= 0.01 && \text{if } z \leq 0 \end{aligned}$$

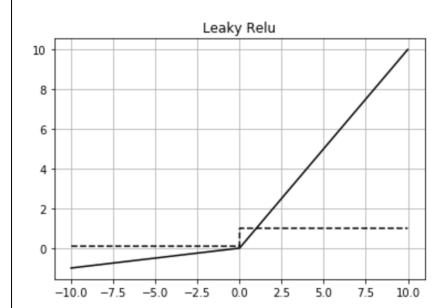
- PRelu function**

$$a = g(x) = \max(\alpha z, z)$$

- same as Leaky ReLU but  $\alpha$  is a parameter to be learnt

- Derivative**

$$\begin{aligned} g'(x) &= 1 && \text{if } z > 0 \\ &= \alpha && \text{if } z \leq 0 \end{aligned}$$



- Softplus function**

$$g(x) = \log(1 + e^x)$$

- continuous approximation of ReLU

- Derivative**

1

Activation functions  $a = g(z)$

## List of possible activation functions

| Name   | Plot | Equation   | Derivative  |
|--|------|--|---|
| Identity                                     |      | $f(x) = x$   | $f'(x) = 1$   |
| Binary step                                  |      | $f(x) = \begin{cases} 0 & \text{for } x < 0 \\ 1 & \text{for } x \geq 0 \end{cases}$               | $f'(x) = \begin{cases} 0 & \text{for } x \neq 0 \\ ? & \text{for } x = 0 \end{cases}$             |
| Logistic (a.k.a Soft step)                   |      | $f(x) = \frac{1}{1 + e^{-x}}$  | $f'(x) = f(x)(1 - f(x))$  |
| Tanh   |      | $f(x) = \tanh(x) = \frac{2}{1 + e^{-2x}} - 1$  | $f'(x) = 1 - f(x)^2$  |
| Arctan                                       |      | $f(x) = \tan^{-1}(x)$  | $f'(x) = \frac{1}{x^2 + 1}$   |
| Rectified Linear Unit (ReLU)                 |      | $f(x) = \begin{cases} 0 & \text{for } x < 0 \\ x & \text{for } x \geq 0 \end{cases}$               | $f'(x) = \begin{cases} 0 & \text{for } x < 0 \\ 1 & \text{for } x \geq 0 \end{cases}$             |
| Parametric Rectified Linear Unit (PReLU) [2] |      | $f(x) = \begin{cases} \alpha x & \text{for } x < 0 \\ x & \text{for } x \geq 0 \end{cases}$        | $f'(x) = \begin{cases} \alpha & \text{for } x < 0 \\ 1 & \text{for } x \geq 0 \end{cases}$        |
| Exponential Linear Unit (ELU) [3]            |      | $f(x) = \begin{cases} \alpha(e^x - 1) & \text{for } x < 0 \\ x & \text{for } x \geq 0 \end{cases}$ | $f'(x) = \begin{cases} f(x) + \alpha & \text{for } x < 0 \\ 1 & \text{for } x \geq 0 \end{cases}$ |
| SoftPlus                                     |      | $f(x) = \log_e(1 + e^x)$   | $f'(x) = \frac{1}{1 + e^{-x}}$  |

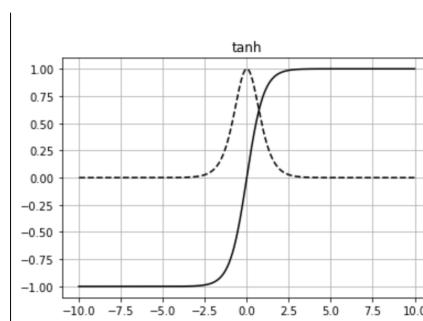
## Weight initialization



## Weight initialization

### Weight initialization with random values

- **Random initialization**
    - $W^{[1]} = \text{np.random.randn}(2, 2) * 0.01$
    - $b^{[1]} = 0$
    - $W^{[2]} = \text{np.random.randn}(1, 2) * 0.01$
    - $b^{[1]} = 0$
  - Remark : b doesn't have the symmetry problem
  - **Why 0.01 ?**
    - If  $W$  is big  $\Rightarrow Z$  is also big  
 $Z^{[1]} = W^{[1]}X + b^{[1]}$  (19)
    - $A^{[1]} = g^{[1]}(Z^{[1]})$  (20)
    - $\Rightarrow$  we are in the flat part of the sigmoid/tanh
      - $\Rightarrow$  slope is small
      - $\Rightarrow$  gradient descent slow
      - $\Rightarrow$  learning slow
    - Better to initialize to a very small value (valid for sigmoid and tanh)

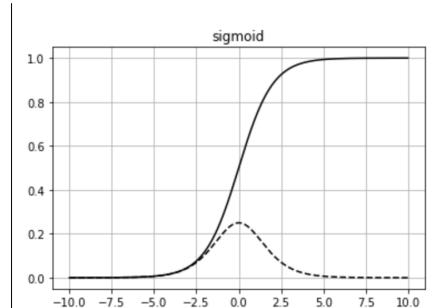


## Weight initialization

## Vanishing gradient (part 2)

- Reminder :
 
$$dz^{[l]} = W^{[l+1]} dz^{[l+1]} \odot g^{[l]}'(z^{[l]})$$

$$db^{[l]} = \frac{1}{m} \sum_m dZ^{[l]}$$
  - Hence for a deep network (supposing  $g^{[l]}(z) = \sigma(z)$ )
 
$$\frac{\partial \mathcal{L}}{\partial b^{[1]}} = \sigma'(z^{[1]})W^{[2]}\sigma'(z^{[2]})W^{[3]}\sigma'(z^{[3]})W^{[4]}\sigma'(z^{[4]})\frac{\partial \mathcal{L}}{\partial a^{[4]}}$$
  - But  $\max_z \sigma'(z) = \frac{1}{4}$  !
    - Therefore, the deeper the network, the fastest the gradient diminishes/vanishes during backpropagation
      - Consequence ?
      - The network stop learning
  - Solution ?
    - Use a ReLu activation



## Regularization

## Regularization

### L1 and L2 regularization

- **Goal ?**

- avoid over-overfitting (high variance)

- **How ?**

- reduce model complexity

- In logistic regression

$$J(\underline{w}, b) = \frac{1}{m} \sum_{i=1}^m \mathcal{L}(\hat{y}^{(i)}, y^{(i)}) + \frac{\lambda_1}{2m} \|w\|_1 + \frac{\lambda_2}{2m} \|w\|_2^2$$

with  $\|w\|_1 = \sum_{j=1}^{n_x} |w_j|$  and  $\|w\|_2^2 = \sum_{j=1}^{n_x} w_j^2 = \underline{w}^T \underline{w}$

- **L1 regularization** (Lasso) :

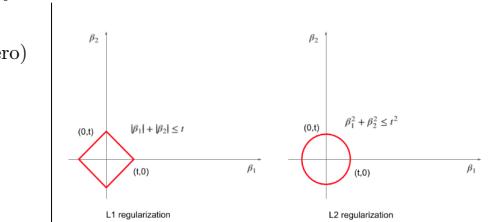
- will end up with sparse  $\underline{w}$  (many zero)

- **L2 regularization** (Ridge) :

- will end up with small values of  $\underline{w}$

- L1+L2 : ELastic Search

- $\lambda$  is the regularization parameter (hyperparameter)



## Regularization

### L1 and L2 regularization

- In neural network

$$J(\dots) = \frac{1}{m} \sum_{i=1}^m \mathcal{L}(\hat{y}^{(i)}, y^{(i)}) + \frac{\lambda}{2m} \sum_{l=1}^L \|\underline{W}^{[l]}\|^2$$

where  $\|\underline{W}\|_F^2 = \|\underline{W}\|_2^2$  is the "Frobenius norm"

$$\|\underline{W}\|^2 = \sum_{i=1}^{n^{[l-1]}} \sum_{j=1}^{n^{[l]}} (\underline{W}_{i,j}^{[l]})^2$$

- In gradient descent ?

$$d'W^{[l]} = dW^{[l]} + \frac{\lambda}{m} W^{[l]}$$

Therefore

$$\begin{aligned} W^{[l]} &\leftarrow W^{[l]} - \alpha d'W^{[l]} \\ &\leftarrow W^{[l]} - \alpha(dW^{[l]} + \frac{\lambda}{m} W^{[l]}) \\ &\leftarrow W^{[l]} - \frac{\alpha\lambda}{m} W^{[l]} - \alpha dW^{[l]} \\ &\leftarrow W^{[l]} \underbrace{\left(1 - \frac{\alpha\lambda}{m}\right)}_{\leq 1} - \alpha dW^{[l]} \end{aligned}$$

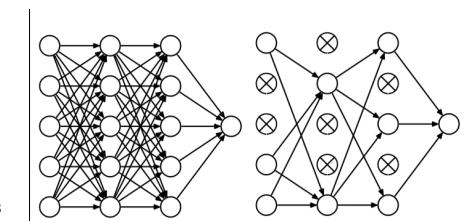
weight decay

## Regularization

### DropOut regularization

- **During training**

- for each training example **randomly turn-off** the neurons of hidden units (with  $p = 0.5$ )
  - this also removes the connections
- for different training examples, turn-off different units
- possible to vary the probability across layers
  - for large matrix  $\underline{W} \Rightarrow p$  is higher
  - for small matrix  $\underline{W} \Rightarrow p$  is lower



- **During testing**

- no drop out

- **Dropout effects :**

- prevents co-adaptation between units
- can be seen as averaging different models that share parameters.
- acts as a powerful regularization scheme.
- since the network is smaller, it is easier to train (as regularization)
- The network cannot rely on any feature, it has to spread out weights
  - Effect : shrinking the squared norm of the weights (similar to L2 regularization)
  - Can be shown to be an adaptive form of L2-regularization

## Alternative gradient descent algorithms

## Alternative gradient descent algorithms

### Momentum

- What ?
  - helps accelerating gradient descent in the relevant direction and dampens oscillations
- How ?
  - add a fraction  $\beta$  of the update vector of the past time step to the current step
- On iteration  $t$ 
  - Compute  $\frac{\partial \mathcal{L}(\dots W(t-1)\dots)}{\partial W}$ ,  $\frac{\partial \mathcal{L}(\dots b(t-1)\dots)}{\partial b}$  on current mini-batch
  - $V_{dw}(t) = \beta V_{dw}(t-1) + (1-\beta) \frac{\partial \mathcal{L}(\dots W(t-1)\dots)}{\partial W}$
  - $W(t) = W(t-1) - \alpha V_{dw}(t)$
- usual choice :  $\beta = 0.9$
- Explanation : the momentum term
  - increases for dimensions whose gradients point in the same directions
  - reduces updates for dimensions whose gradients change directions.
  - gain faster convergence and reduced oscillation.

Geoffroy Peeters LTCE/Télécom ParisTech 14/15

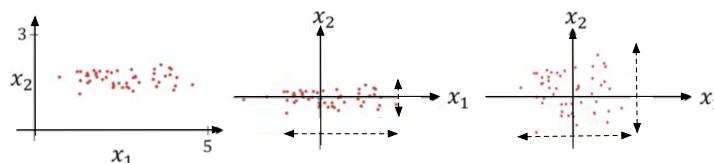
## Normalization inputs

- "Standardizing" :
  - subtracting a measure of location and dividing by a measure of scale
  - subtract the mean and divide by the standard deviation

$$\mu_d = \frac{1}{m} \sum_{i=1}^m x_d^{(i)} \rightarrow x_d = x_d - \mu_d$$

$$\sigma_d^2 = \frac{1}{m} \sum_{i=1}^m (x_d^{(i)} - \mu_d)^2 \rightarrow x_d = \frac{1}{\sigma_d} x_d$$

- We will use the same  $\mu$  and  $\sigma^2$  to normalize the test set



- "Normalizing" :
  - rescaling by the minimum and range of the vector
  - make all the elements lie between 0 and 1

## Batch Normalization (BN)

- Objective ?
  - Apply the same normalization for the input of each layer  $l$ 
    - allows to learn faster
  - Try to **reduce the "covariate shift"**
    - the inputs of a given layer  $l$  is the outputs of the previous layer  $l-1$ 
      - these outputs depends on the current parameters of the previous layer which change over training!
    - **normalize the output of the previous layer  $a^{[l-1]}$** 
      - in practice normalize the pre-activation  $z^{[l-1]}$
  - We don't want all units to always mean 0 and standard-deviation 1
    - We still allows to learn an appropriate bias  $\beta$  and scale  $\gamma$  to the input to the non-linear function  $g^{[l]}$

Geoffroy Peeters LTCE/Télécom ParisTech 15/15

Geoffroy Peeters LTCE/Télécom ParisTech 14/15

## Batch Normalization (BN)

- **Batch Normalization**

- Given some intermediate values in the network :  $z^{[l](1)}, z^{[l](2)}, \dots, z^{[l](m)}$
- Compute

$$\mu^{[l]} = \frac{1}{m} \sum_i z^{[l](i)}$$

$$\sigma^2{}^{[l]} = \frac{1}{m} \sum_i (z^{[l](i)} - \mu^{[l]})^2$$

$$z_{norm}^{[l](i)} = \frac{z^{[l](i)} - \mu^{[l]}}{\sqrt{\sigma^2{}^{[l]}} + \epsilon}$$

$$\tilde{z}^{[l](i)} = \gamma \cdot z_{norm}^{[l](i)} + \beta$$

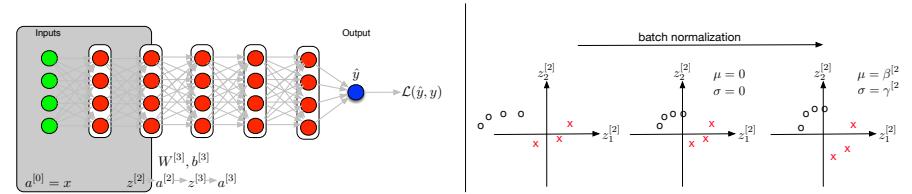
$$a^{[l](i)} = g^{[l]}(\tilde{z}^{[l](i)})$$

- **New parameters to be trained :  $\gamma$  and  $\beta$**

- $\gamma$  allows to set the variance of  $\tilde{z}^{[l]}$
- $\beta$  allows to set the mean of  $\tilde{z}^{[l]}$
- If  $\gamma = \sqrt{\sigma^2{}^{[l]}} + \epsilon$  and  $\beta = \mu$  then  $\tilde{z}^{[l](i)} = z^{[l](i)}$

## Batch Normalization (BN)

- **Why does it work ?**



- We consider the left part as the input ( $a^{[2]}$ ) of the current layer  $l = 3$
- **Problem :** when we train  $W^{[3]}, b^{[3]}$ , the input values  $a^{[2]}$  will change over time (since previous layers weight also change during training)
- **Solution :** BN allows to reduce these changes

- BN ensures that the mean and variance will remain the same ( $\beta$  and  $\gamma$ )
  - the values become more stable
  - it reduces the coupling between the layers
  - it speed up training

- **Regularization effect**

- each mini-batch is scaled using  $\mu$  and  $\sigma$  computed only on  $X^{\{t\}}$
- add some noise to the values  $z^{[l]}$  within that mini-batch
  - // to drop-out which add noise to each hidden layer activation
- to obtain regularization effect, using small mini-batch (64) is better than large ones (512)