

Imperial College
London

April 28th, 2021

Feed-Forward Neural Networks

Olivier Dubrule

1

Imperial College
London

Objectives of the Day

- Focus on Supervised Learning
- Present Logistic Regression as a Single Neuron Operation
- Generalize to Feed-Forward Neural Networks
- Illustrate the Back-Propagation Algorithm
- Introduce Stochastic Gradient Descent
- Show Examples

2

Imperial College
London

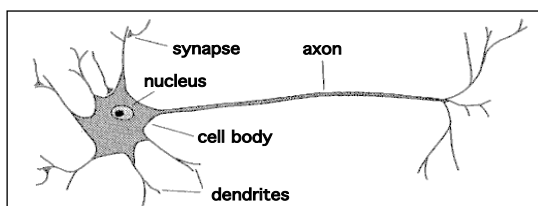
Feed-Forward Neural Networks

1. From Logistic Regression to Single Neuron Representation
2. Feed-Forward Neural Networks
3. Back-Propagation
4. Batch, Stochastic and Mini-Batch Gradient Descent
5. Conclusion

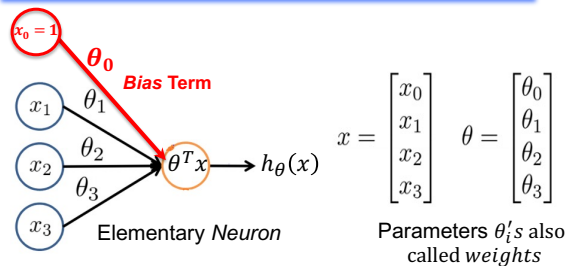
3

Imperial College
London

Logistic Regression: Analogy with Human Neuron



- A neuron has
 - Branching input (dendrites)
 - Branching output (the axon)



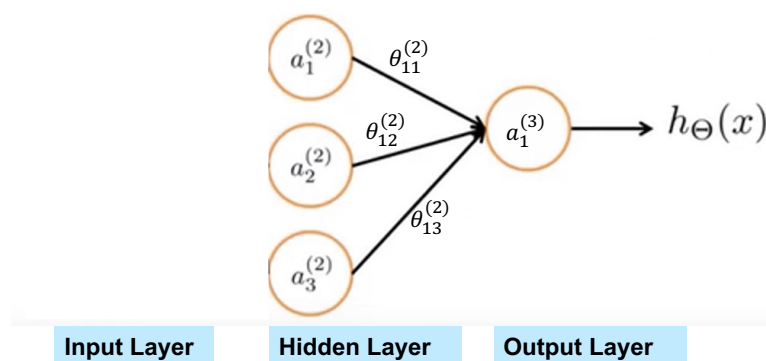
$$h_{\theta}(x) = \sigma(\theta^T x) = \frac{1}{1 + e^{-\theta^T x}}$$

σ also called *activation function* g

4

Imperial College
London

From Single Neuron to Feed-forward Neural Network



Also called “**Fully-Connected Network**”. Historically this is a generalization of the Multi-Layer Perceptron or MLP, which was the first network capable of learning its weights (Rosenblatt, 1961).

5

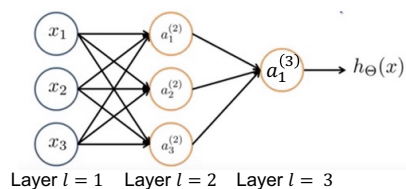
Imperial College
London

Feed-Forward Neural Networks

1. From Logistic Regression to Single Neuron Representation
2. Feed-Forward Neural Networks
3. Back-Propagation
4. Batch, Stochastic and Mini-Batch Gradient Descent
5. Conclusion

6

A Simple Neural Network in Matrix Form



$a^{(l)}$ = Activation vector of layer l

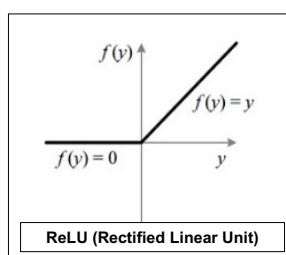
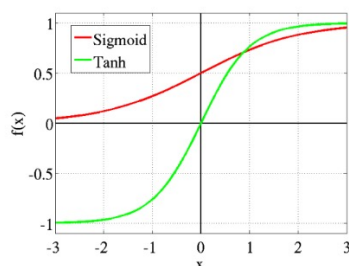
$\theta^{(l)}$ = Matrix of weights controlling mapping from layer l to layer $l + 1$
 $\theta_{jk}^{(l)}$ = weight from neuron k in layer l to neuron j in layer $l + 1$

For example (if we have a bias term):

$$a^{(2)} = \begin{pmatrix} a_1^{(2)} \\ a_2^{(2)} \\ a_3^{(2)} \end{pmatrix} = g(\theta^{(1)} a^{(1)}) = g(\theta^{(1)} x) = g \left(\begin{pmatrix} \theta_{10}^{(1)} & \theta_{11}^{(1)} & \theta_{12}^{(1)} & \theta_{13}^{(1)} \\ \theta_{20}^{(1)} & \theta_{21}^{(1)} & \theta_{22}^{(1)} & \theta_{23}^{(1)} \\ \theta_{30}^{(1)} & \theta_{31}^{(1)} & \theta_{32}^{(1)} & \theta_{33}^{(1)} \end{pmatrix} \begin{pmatrix} x_0 \\ x_1 \\ x_2 \\ x_3 \end{pmatrix} \right)$$

\uparrow
 g is the activation function
 (generalization of σ)

7

Possible Choices for the Non-Linear Activation Function g 

Sigmoid σ between 0 and 1, mostly used for output layer of binary classification problems, in order to allow a probabilistic interpretation of the result.

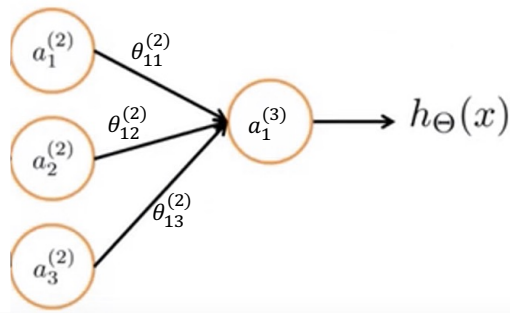
Tanh is between -1 and +1 with mean zero (instead of mean 0.5 for sigmoid function). Almost identity function close to the origin.

ReLU is such that gradient does not vanish for non-zero values, most often used activation function in hidden layers..

8

Imperial College
London

The last layer of a NN is a Logistic Regression!!








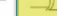

9

Imperial College
London

Why Logistic Regression is not Enough?

- There are some complex mathematical functions (such as the “Exclusive OR”) that Logistic Regression cannot represent

Logic Gates

Name	NOT	AND	NAND	OR	NOR	XOR	XNOR																																																																																																
Alg. Expr.	\overline{A}	AB	\overline{AB}	$A + B$	$\overline{A + B}$	$A \oplus B$	$\overline{A \oplus B}$																																																																																																
Symbol																																																																																																							
Truth Table	<table><tr><th>A</th><th>X</th></tr><tr><td>0</td><td>1</td></tr><tr><td>1</td><td>0</td></tr></table>	A	X	0	1	1	0	<table><tr><th>B</th><th>A</th><th>X</th></tr><tr><td>0</td><td>0</td><td>0</td></tr><tr><td>0</td><td>1</td><td>0</td></tr><tr><td>1</td><td>0</td><td>0</td></tr><tr><td>1</td><td>1</td><td>1</td></tr></table>	B	A	X	0	0	0	0	1	0	1	0	0	1	1	1	<table><tr><th>B</th><th>A</th><th>X</th></tr><tr><td>0</td><td>0</td><td>1</td></tr><tr><td>0</td><td>1</td><td>1</td></tr><tr><td>1</td><td>0</td><td>1</td></tr><tr><td>1</td><td>1</td><td>0</td></tr></table>	B	A	X	0	0	1	0	1	1	1	0	1	1	1	0	<table><tr><th>B</th><th>A</th><th>X</th></tr><tr><td>0</td><td>0</td><td>0</td></tr><tr><td>0</td><td>1</td><td>1</td></tr><tr><td>1</td><td>0</td><td>1</td></tr><tr><td>1</td><td>1</td><td>0</td></tr></table>	B	A	X	0	0	0	0	1	1	1	0	1	1	1	0	<table><tr><th>B</th><th>A</th><th>X</th></tr><tr><td>0</td><td>0</td><td>1</td></tr><tr><td>0</td><td>1</td><td>0</td></tr><tr><td>1</td><td>0</td><td>0</td></tr><tr><td>1</td><td>1</td><td>1</td></tr></table>	B	A	X	0	0	1	0	1	0	1	0	0	1	1	1	<table><tr><th>B</th><th>A</th><th>X</th></tr><tr><td>0</td><td>0</td><td>0</td></tr><tr><td>0</td><td>1</td><td>1</td></tr><tr><td>1</td><td>0</td><td>1</td></tr><tr><td>1</td><td>1</td><td>0</td></tr></table>	B	A	X	0	0	0	0	1	1	1	0	1	1	1	0	<table><tr><th>B</th><th>A</th><th>X</th></tr><tr><td>0</td><td>0</td><td>1</td></tr><tr><td>0</td><td>1</td><td>0</td></tr><tr><td>1</td><td>0</td><td>0</td></tr><tr><td>1</td><td>1</td><td>1</td></tr></table>	B	A	X	0	0	1	0	1	0	1	0	0	1	1	1
A	X																																																																																																						
0	1																																																																																																						
1	0																																																																																																						
B	A	X																																																																																																					
0	0	0																																																																																																					
0	1	0																																																																																																					
1	0	0																																																																																																					
1	1	1																																																																																																					
B	A	X																																																																																																					
0	0	1																																																																																																					
0	1	1																																																																																																					
1	0	1																																																																																																					
1	1	0																																																																																																					
B	A	X																																																																																																					
0	0	0																																																																																																					
0	1	1																																																																																																					
1	0	1																																																																																																					
1	1	0																																																																																																					
B	A	X																																																																																																					
0	0	1																																																																																																					
0	1	0																																																																																																					
1	0	0																																																																																																					
1	1	1																																																																																																					
B	A	X																																																																																																					
0	0	0																																																																																																					
0	1	1																																																																																																					
1	0	1																																																																																																					
1	1	0																																																																																																					
B	A	X																																																																																																					
0	0	1																																																																																																					
0	1	0																																																																																																					
1	0	0																																																																																																					
1	1	1																																																																																																					

B and A can be seen as two input features (x_1, x_2) and X as the label

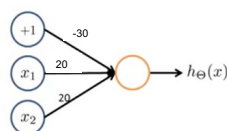
<https://medium.com/autonomous-agents/how-to-teach-logic-to-your-neuralnetworks-116215c71a49>

10

Imperial College
London

Single Neurons as Logical Functions

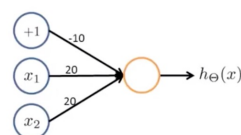
The activation function g used is the sigmoid σ



Show that $h_{\theta}(x)$ models the « AND » Logical Function

$$x_2 \begin{array}{|c|c|} \hline 0 & 1 \\ \hline 0 & 0 \\ \hline \end{array} x_1$$

x_1	x_2	$h_{\theta}(x)$
0	0	
0	1	
1	0	
1	1	



Show that $h_{\theta}(x)$ models the « OR » Logical Function

$$x_2 \begin{array}{|c|c|} \hline 1 & 1 \\ \hline 0 & 1 \\ \hline \end{array} x_1$$

x_1	x_2	$h_{\theta}(x)$
0	0	
0	1	
1	0	
1	1	

11

11

Example of the "AND"

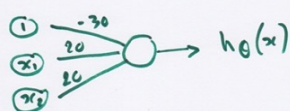


image of $(0,0) = 0$ because $\sigma(-30) < 0.5$
 image of $(1,0) = 0$ because $\sigma(-10) < 0.5$
 image of $(0,1) = 0$ because $\sigma(-10) < 0.5$
 image of $(1,1) = 1$ because $\sigma(+10) > 0.5$

$$x_2 \begin{array}{|c|c|} \hline 0 & 1 \\ \hline 0 & 0 \\ \hline \end{array} x_1$$

Example of the "OR"

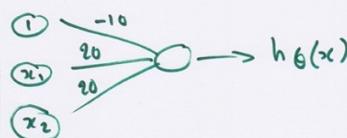


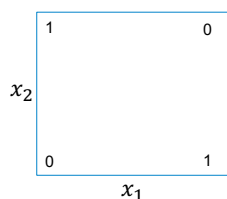
image of $(0,0) = 0$ because $\sigma(-10) < 0.5$
 image of $(1,0) = 1$ because $\sigma(10) > 0.5$
 image of $(0,1) = 1$ because $\sigma(10) > 0.5$
 image of $(1,1) = 1$ because $\sigma(30) > 0.5$

$$x_2 \begin{array}{|c|c|} \hline 1 & 1 \\ \hline 0 & 1 \\ \hline \end{array} x_1$$

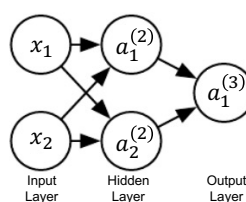
12

Modelling the XOR Function with a Neural Network (1)

The "Exclusive OR" or XOR function: when one of the input x_1, x_2 values is equal to 1 and the other to 0, XOR returns the value 1, otherwise 0.



This cannot be approximated by a single neuron, as we saw that the Logistic Regression Decision Boundary was a line. We need a hidden layer! Let us try the simplest possible neural network:



Inspired from Goodfellow et al, 2016

13

Modelling the XOR Function with a Neural Network (2)

Assume that we have a Bias term and the parameters of the hidden layer are:

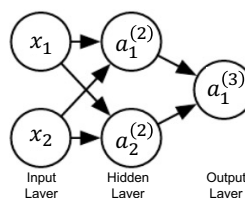
$$\theta^{(1)} = \begin{pmatrix} 0 & 1 & 1 \\ -1 & 1 & 1 \end{pmatrix}$$

The activation function is the ReLU Function.

And that the parameters of the output layer are:

$$\theta^{(2)T} = \begin{pmatrix} -\frac{1}{2} & \frac{5}{8} & -1 \end{pmatrix}$$

The activation function is the Sigmoid Function.



Exercise: Calculate the network output for each of the four input points $\begin{pmatrix} x_1 \\ x_2 \end{pmatrix}$:

$$\begin{pmatrix} 0 \\ 0 \end{pmatrix}, \begin{pmatrix} 0 \\ 1 \end{pmatrix}, \begin{pmatrix} 1 \\ 0 \end{pmatrix}, \begin{pmatrix} 1 \\ 1 \end{pmatrix}$$

Inspired from Goodfellow et al, 2016

14

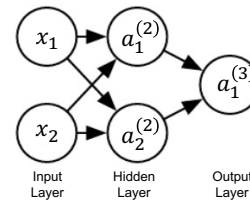
Modelling the XOR Function with a Neural Network (3)

$$\text{We have } \begin{pmatrix} a_1^{(2)} \\ a_2^{(2)} \end{pmatrix} = \text{ReLU} \left(\theta^{(1)} \begin{pmatrix} x_0 \\ x_1 \\ x_2 \end{pmatrix} \right)$$

$$\text{with } \theta^{(1)} \begin{pmatrix} x_0 \\ x_1 \\ x_2 \end{pmatrix} = \begin{pmatrix} 0 & 1 & 1 \\ -1 & 1 & 1 \end{pmatrix} \begin{pmatrix} x_0 \\ x_1 \\ x_2 \end{pmatrix} = \begin{pmatrix} 0 + x_1 + x_2 \\ -1 + x_1 + x_2 \end{pmatrix}$$

$$\text{And we have } a_1^{(3)} = \sigma \left(\theta^{(2)T} \begin{pmatrix} a_0^{(2)} \\ a_1^{(2)} \\ a_2^{(2)} \end{pmatrix} \right)$$

$$\text{with } \theta^{(2)T} \begin{pmatrix} a_0^{(2)} \\ a_1^{(2)} \\ a_2^{(2)} \end{pmatrix} = \begin{pmatrix} -\frac{1}{2} & \frac{5}{8} & -1 \end{pmatrix} \begin{pmatrix} a_0^{(2)} \\ a_1^{(2)} \\ a_2^{(2)} \end{pmatrix} = \left(-\frac{1}{2} + \frac{5}{8} a_1^{(2)} - a_2^{(2)} \right)$$



Inspired from Goodfellow et al, 2016

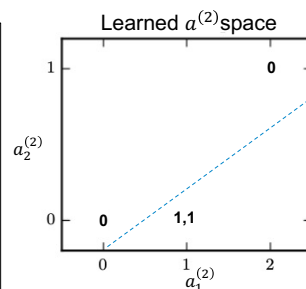
15

Modelling the XOR Function with a Neural Network (4)

It is easy to see that :

$$\begin{aligned} \begin{pmatrix} x_1 \\ x_2 \end{pmatrix} = \begin{pmatrix} 0 \\ 0 \end{pmatrix} &\Rightarrow \begin{pmatrix} a_1^{(2)} \\ a_2^{(2)} \end{pmatrix} = \text{ReLU} \begin{pmatrix} 0 \\ -1 \end{pmatrix} = \begin{pmatrix} 0 \\ 0 \end{pmatrix} \\ \begin{pmatrix} x_1 \\ x_2 \end{pmatrix} = \begin{pmatrix} 0 \\ 1 \end{pmatrix} &\Rightarrow \begin{pmatrix} a_1^{(2)} \\ a_2^{(2)} \end{pmatrix} = \text{ReLU} \begin{pmatrix} 1 \\ 0 \end{pmatrix} = \begin{pmatrix} 1 \\ 0 \end{pmatrix} \\ \begin{pmatrix} x_1 \\ x_2 \end{pmatrix} = \begin{pmatrix} 1 \\ 0 \end{pmatrix} &\Rightarrow \begin{pmatrix} a_1^{(2)} \\ a_2^{(2)} \end{pmatrix} = \text{ReLU} \begin{pmatrix} 1 \\ 0 \end{pmatrix} = \begin{pmatrix} 1 \\ 0 \end{pmatrix} \\ \begin{pmatrix} x_1 \\ x_2 \end{pmatrix} = \begin{pmatrix} 1 \\ 1 \end{pmatrix} &\Rightarrow \begin{pmatrix} a_1^{(2)} \\ a_2^{(2)} \end{pmatrix} = \text{ReLU} \begin{pmatrix} 2 \\ 1 \end{pmatrix} = \begin{pmatrix} 2 \\ 1 \end{pmatrix} \end{aligned}$$

The hidden layer changes the position of the four points such that they can now be linearly separated by the output layer!



..and the output layer $a_1^{(3)} = \sigma \left(\frac{5}{8} a_1^{(2)} - a_2^{(2)} - \frac{1}{2} \right)$ gives the probabilities that lead to the XOR function for the four points!

Inspired from Goodfellow et al, 2016

16

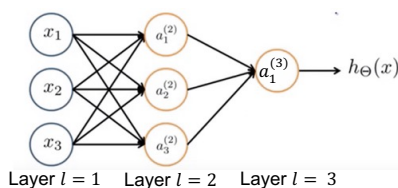
Neural Networks and Logistic Regression

- There are some complex mathematical functions (such as the “Exclusive OR”) that Logistic Regression cannot represent.
- Linear Logistic Regression represents a single neuron and is the shallowest form of neural network. The input features are given and the Decision Boundary is linear. The output layer of a neural network can be seen as a Logistic Regression module. But you may need to “massage” the input features to help this linear Decision Boundary separate the data properly in the output layer. The Neural Network, by transforming the initial variable at each hidden layer, does this “massaging”.
- The “Universal Approximation Theorem” (Hornik et al, 1989) states that a neural network with at least one hidden layer can approximate arbitrarily well any function from any finite dimensional space to another, provided that the network is given enough neurons.

17

17

Inside a Feed-Forward Neural Network



$a_i^{(l)}$ = “activation” of unit i in layer l
 $\theta^{(l)}$ = matrix of weights controlling mapping from layer l to layer $l + 1$
 $\theta_{jk}^{(l)}$ = weight from neuron k in layer l to neuron j in layer $l + 1$

$$a_1^{(2)} = g \left(\theta_{10}^{(1)} x_0 + \theta_{11}^{(1)} x_1 + \theta_{12}^{(1)} x_2 + \theta_{13}^{(1)} x_3 \right)$$

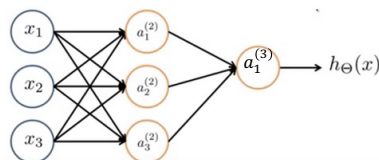
$$a_2^{(2)} = g \left(\theta_{20}^{(1)} x_0 + \theta_{21}^{(1)} x_1 + \theta_{22}^{(1)} x_2 + \theta_{23}^{(1)} x_3 \right)$$

$$a_3^{(2)} = g \left(\theta_{30}^{(1)} x_0 + \theta_{31}^{(1)} x_1 + \theta_{32}^{(1)} x_2 + \theta_{33}^{(1)} x_3 \right)$$

$$h_{\theta}(x) = a_1^{(3)} = g' \left(\theta_{10}^{(2)} a_0^{(2)} + \theta_{11}^{(2)} a_1^{(2)} + \theta_{12}^{(2)} a_2^{(2)} + \theta_{13}^{(2)} a_3^{(2)} \right)$$

18

The Number of Parameters of a Neural Network



Number of parameters (or weights) $\theta_{jk}^{(l)}$, assuming bias term for each layer?

On this example:

(Number of nodes in input layer + 1) \times Number of nodes in hidden layer +
(Number of nodes in hidden layer + 1) \times Number of nodes in output layer

$$= (3+1) \times 3 + (3+1) \times 1 = 16$$

19

Notations for Neural Network With Two Hidden Layers

Forward propagation

$$a^{(1)} = x$$

$$z^{(2)} = \theta^{(1)} a^{(1)}$$

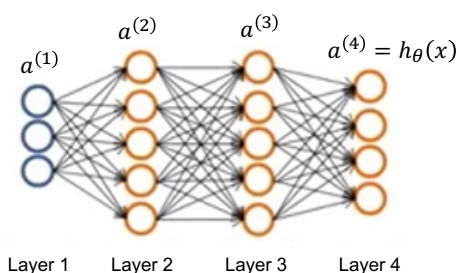
$$a^{(2)} = g(z^{(2)})$$

$$z^{(3)} = \theta^{(2)} a^{(2)}$$

$$a^{(3)} = g(z^{(3)})$$

$$z^{(4)} = \theta^{(3)} a^{(3)}$$

$$a^{(4)} = g'(z^{(4)}) = h_{\theta}(x)$$



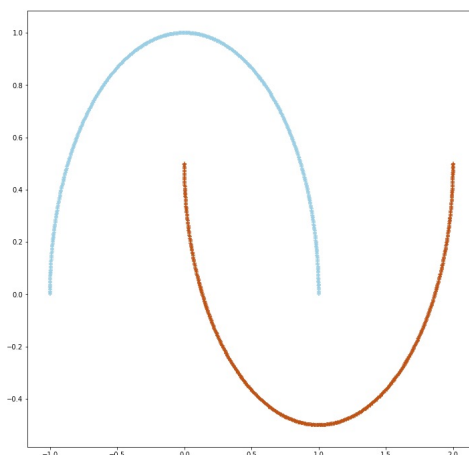
Number of parameters (assuming bias terms)?

$$(3+1) \times 5 + (5+1) \times 5 + (5+1) \times 4 = 74$$

20

Imperial College
London

Neural Network on Logistic Regression Example



The $m = 1000$ data points $(x_1^i, x_2^i)_{i=1,1000}$ are in the plane.

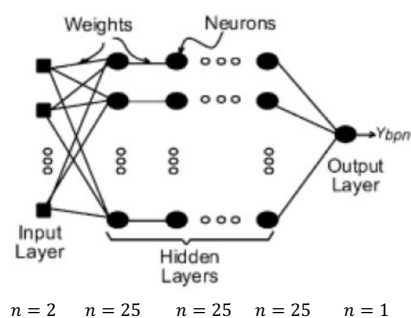
A group of data are one color,
the other group another color.

**Question: predict the color at
each location of the plane.**

21

Imperial College
London

A Simple Neural Network Example (2)



Neural Network Architecture

Input Layer: Two Neurons (x_1 and x_2)

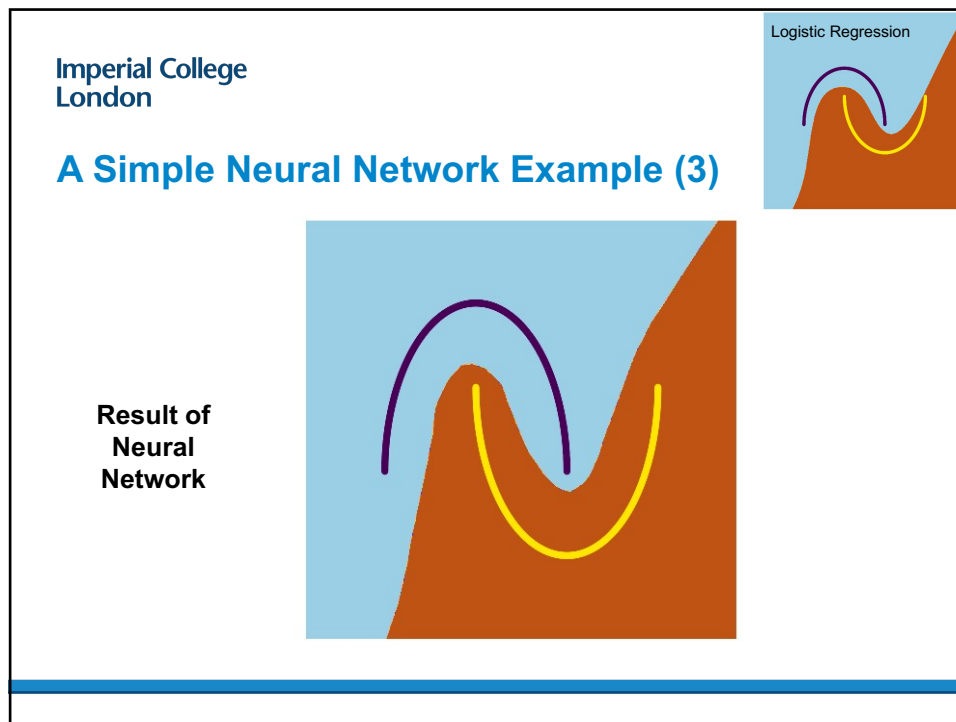
Three Hidden Layers Each 25 Neurons

Output Layer: One neuron: Probability of
Being Blue

Total Number of Weights (if bias terms):

$$(2+1) \times 25 + (25+1) \times 25 + (25+1) \times 25 + (25+1) \times 1 = 1401$$

22



23

Imperial College
London

More Simple Neural Network Examples

<https://playground.tensorflow.org>

24

Imperial College
London

More Simple Neural Network Classification Examples

<https://playground.tensorflow.org>

Four classification examples are discussed.

The third one can be addressed by logistic regression, that is with no hidden layer at all and just one output neuron. This is true also with not linearly separated data points.

The second one cannot be addressed by logistic regression, because it is clearly non linear and it is similar to the XOR function. We need at least one hidden layer.

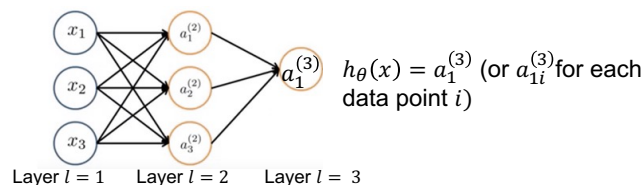
The first one is the circle. With just one hidden layer and 4 neurons in it we reproduce the circle (<https://www.youtube.com/watch?v=ru9dXF04iSE>). But of course if we use x_1^2 and x_2^2 as input we do not even need a hidden layer, meaning that logistic regression does the job!!

25

Imperial College
London

Cost Function for Binary Classification

The class of each data point i is: $y^{(i)} = 0$ or 1



The loss function is the cross-entropy already used in the Logistic Regression case:

$$J(\theta) = -\frac{1}{m} \sum_{i=1}^m \left[y^{(i)} \log a_{1i}^{(3)} + (1 - y^{(i)}) \log (1 - a_{1i}^{(3)}) \right]$$

We will see in the third week that cross-entropy is equal to (minus) the log-likelihood of a Bernoulli distribution.

26

Imperial College
London

Multi-Class Classification

The “One vs Rest” Approach. Example of 3 Classes.

- Train the 3 Neural Networks associated with the three following binary problems:
 - Merge classes 2 and 3 as class “zero” and calculate parameters θ_1 of function:

$$h_{\theta_1}^1(x) = P(y = 1|x, \theta_1)$$
 - Merge classes 1 and 3 as class “zero” and calculate parameters θ_2 of function:

$$h_{\theta_2}^2(x) = P(y = 2|x, \theta_2)$$
 - Merge classes 1 and 2 as class “zero” and calculate parameters θ_3 of function:

$$h_{\theta_3}^3(x) = P(y = 3|x, \theta_3)$$
- For each new unlabeled point x , calculate the three probabilities above and pick as the predicted class the one corresponding to the highest of the three probabilities.

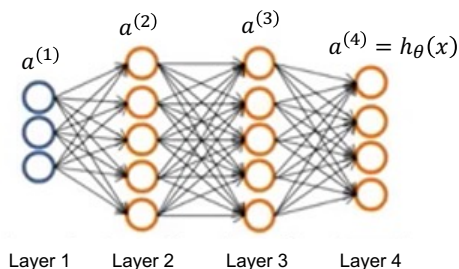
27

Imperial College
London

Multi-Class Classification with Softmax

Forward propagation

$$\begin{aligned}
 a^{(1)} &= x \\
 z^{(2)} &= \theta^{(1)} a^{(1)} \\
 a^{(2)} &= g(z^{(2)}) \\
 z^{(3)} &= \theta^{(2)} a^{(2)} \\
 a^{(3)} &= g(z^{(3)}) \\
 z^{(4)} &= \theta^{(3)} a^{(3)} \\
 a^{(4)} &= \text{Softmax}(z^{(4)}) = h_{\theta}(x)
 \end{aligned}$$



28

Imperial College
London

Multi-Class Classification with Softmax

Suppose the output before applying the activation ($z^{(4)}$ in the previous slide) of

the last layer of the neural network is $\begin{pmatrix} z_1 \\ z_2 \\ \vdots \\ z_i \\ \vdots \\ z_{n-1} \\ z_n \end{pmatrix}$ of size n

The **Softmax** function transforms it into an output probability vector:

$$\begin{pmatrix} p_1 \\ p_2 \\ p_3 \\ \vdots \\ p_i \\ \vdots \\ p_{n-1} \\ p_n \end{pmatrix} = \begin{pmatrix} \frac{e^{z_1}}{\sum_{j=1}^n e^{z_j}} \\ \frac{e^{z_2}}{\sum_{j=1}^n e^{z_j}} \\ \vdots \\ \frac{e^{z_i}}{\sum_{j=1}^n e^{z_j}} \\ \vdots \\ \frac{e^{z_{n-1}}}{\sum_{j=1}^n e^{z_j}} \\ \frac{e^{z_n}}{\sum_{j=1}^n e^{z_j}} \end{pmatrix}$$

29

Imperial College
London

A SoftMax Function Calculation Example

$$\begin{pmatrix} z_1 \\ z_2 \\ z_3 \\ z_4 \end{pmatrix} = \begin{pmatrix} 5 \\ -1 \\ 2 \\ -4 \end{pmatrix} \quad \begin{pmatrix} p_1 \\ p_2 \\ p_3 \\ p_4 \end{pmatrix} = \begin{pmatrix} \frac{148.41}{156.19} \\ \frac{0.37}{156.19} \\ \frac{7.39}{156.19} \\ \frac{0.02}{156.19} \end{pmatrix} = \begin{pmatrix} 0.950 \\ 0.003 \\ 0.047 \\ 0.000 \end{pmatrix}$$

SoftMax Function does two things:

- Transform the vector into probabilities that are between 0 and 1 and add up to 1
- Transform the largest value into a value close to 1 and the lowest into one close to 0

30

Imperial College
London

Which Cross-Entropy Function for Training with Softmax? One-Hot Encoding for Defining each Class Membership.

Take the example of predicting whether the colour at one pixel of an image is **brown**, **green** or **blue**. We have three classes.

First approach

Code each colour as a number: 1 for brown, 2 for green, 3 for blue.

But this may create an artificial distance between brown and blue larger than between brown and green or green and blue!

Second approach: one-hot encoding

Represent the class of each pixel by a vector c of dimension equal to the number of classes:

If pixel is brown: $c = \begin{pmatrix} 1 \\ 0 \\ 0 \end{pmatrix}$, If pixel is green: $c = \begin{pmatrix} 0 \\ 1 \\ 0 \end{pmatrix}$, If pixel is blue: $c = \begin{pmatrix} 0 \\ 0 \\ 1 \end{pmatrix}$.

31

Imperial College
London

Loss Function in Softmax Multi-Class Classification

For one hot-encoded yellow data point i of the Training Set: $y^{(i)} = \begin{pmatrix} 0 \\ 1 \\ 0 \end{pmatrix}$

If $\begin{pmatrix} p_1 \\ p_2 \\ p_3 \end{pmatrix}$ is the probability vector calculated by Softmax for this point, the cross-entropy is defined as:

$$J(\theta) = -0 \times \log p_1 - 1 \times \log p_2 - 0 \times \log p_3 = -\log p_2$$

So, in the multi-class Softmax case, the cross-entropy is minus the log of the output probability associated with the actual class of the data point.

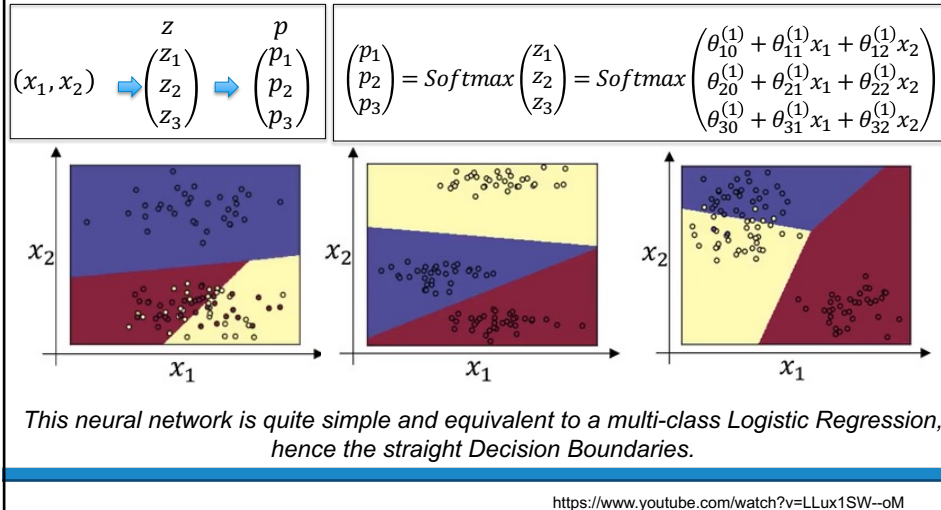
Logically, if $p_2 = 1$, $J(\theta) = 0$, and if $p_2 = 0$, $J(\theta) = +\infty$

For the binary case, it is easy to check that we get the formula already seen on Monday.

For m data points, the above is averaged for all the data points.

32

Examples of Classifying 2-D Data into 3 Classes



33

Once Trained, how does Softmax Predict a Test Point's Class?

Suppose the output of the last layer of the neural network is $\begin{pmatrix} z_1 \\ z_2 \\ \vdots \\ z_i \\ \vdots \\ z_{n-1} \\ z_n \end{pmatrix}$ of size n

The **Softmax** function transforms it into an output probability vector:

$$\begin{pmatrix} p_1 \\ p_2 \\ p_3 \\ \vdots \\ p_i \\ \vdots \\ p_{n-1} \\ p_n \end{pmatrix} = \begin{pmatrix} \frac{e^{z_1}}{\sum_{j=1}^n e^{z_j}} \\ \frac{e^{z_2}}{\sum_{j=1}^n e^{z_j}} \\ \vdots \\ \frac{e^{z_i}}{\sum_{j=1}^n e^{z_j}} \\ \vdots \\ \frac{e^{z_{n-1}}}{\sum_{j=1}^n e^{z_j}} \\ \frac{e^{z_n}}{\sum_{j=1}^n e^{z_j}} \end{pmatrix}$$

**The class with the highest
Softmax probability is selected!**

34

Imperial College
London

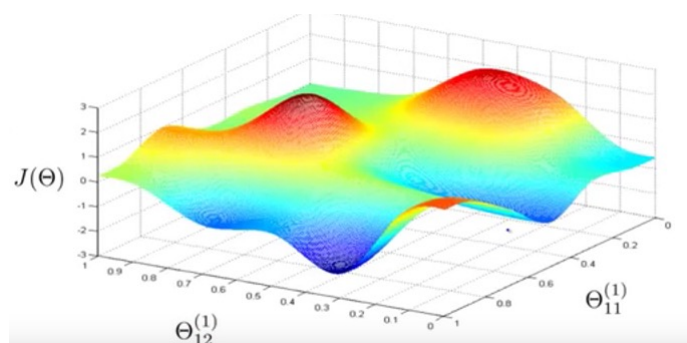
Feed-Forward Neural Networks

1. From Logistic Regression to Single Neuron Representation
2. Feed-Forward Neural Networks
3. Back-Propagation
4. Batch, Stochastic and Mini-Batch Gradient Descent
5. Conclusion

35

Imperial College
London

How to Minimize the Loss Function with Neural Nets?



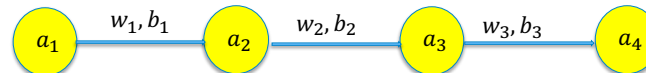
36

Imperial College
London

A Look at Back-Propagation Using a Simple Example (1)

Let us take a very simple network of $L = 4$ layers with one neuron in each.

We have six parameters, three weights w_i and three biases b_i (we could have used the θ notation as before but the (w_i, b_i) notation is quite common too!).



Suppose we have just one sample (x, y) (hence $a_1 = x$ and y is the real number target), and a regression neural network.

The L_2 cost function is $C = \frac{1}{2}(a_4 - y)^2$.

We have: $a_2 = g(z_2)$ with $z_2 = w_1 a_1 + b_1$
 $a_3 = g(z_3)$ with $z_3 = w_2 a_2 + b_2$
 $a_4 = g(z_4)$ with $z_4 = w_3 a_3 + b_3$

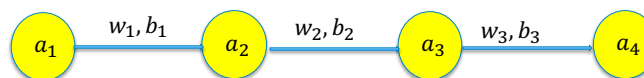
How can we calculate the derivative of the Cost Function C according to each of the six parameters?

<https://www.youtube.com/watch?v=tleHLnjs5U8>

37

Imperial College
London

A Look at Back-Propagation Using a Simple Example (2)



The L_2 cost function is $C = \frac{1}{2}(a_4 - y)^2$.

We have: $a_2 = g(z_2)$ with $z_2 = w_1 a_1 + b_1$
 $a_3 = g(z_3)$ with $z_3 = w_2 a_2 + b_2$
 $a_4 = g(z_4)$ with $z_4 = w_3 a_3 + b_3$

Using the chain rule, the derivatives according to w_3 and b_3 are:

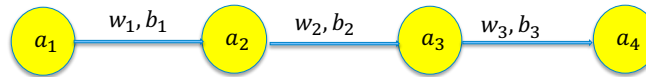
$$\frac{\partial C}{\partial w_3} = \frac{\partial C}{\partial a_4} \frac{\partial a_4}{\partial z_4} \frac{\partial z_4}{\partial w_3} = (a_4 - y)g'(z_4)a_3 \frac{\partial C}{\partial b_3} = \frac{\partial C}{\partial a_4} \frac{\partial a_4}{\partial z_4} \frac{\partial z_4}{\partial b_3} = (a_4 - y)g'(z_4)$$

<https://www.youtube.com/watch?v=tleHLnjs5U8>

38

Imperial College
London

A Look at Back-Propagation Using a Simple Example (3)



Thanks to the formula: $\frac{\partial C}{\partial a_3} = \frac{\partial C}{\partial a_4} \frac{\partial a_4}{\partial z_4} \frac{\partial z_4}{\partial a_3} = (a_4 - y)g'(z_4)w_3$

We can keep moving backwards and now obtain the derivatives of C in w_2 and b_2

$$\frac{\partial C}{\partial w_2} = \frac{\partial C}{\partial a_3} \frac{\partial a_3}{\partial z_3} \frac{\partial z_3}{\partial w_2} = (a_4 - y)g'(z_4)w_3g'(z_3)a_2$$

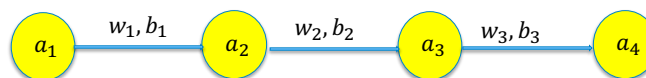
$$\frac{\partial C}{\partial b_2} = \frac{\partial C}{\partial a_3} \frac{\partial a_3}{\partial z_3} \frac{\partial z_3}{\partial b_2} = (a_4 - y)g'(z_4)w_3g'(z_3)$$

<https://www.youtube.com/watch?v=tleHLnjs5U8>

39

Imperial College
London

A Look at Back-Propagation Using a Simple Example (4)



Thanks to the formula: $\frac{\partial C}{\partial a_2} = \frac{\partial C}{\partial a_3} \frac{\partial a_3}{\partial z_3} \frac{\partial z_3}{\partial a_2} = (a_4 - y)g'(z_4)w_3g'(z_3)w_1$

We can keep moving backwards and obtain the derivatives of C in w_1 and b_1

$$\frac{\partial C}{\partial w_1} = \frac{\partial C}{\partial a_2} \frac{\partial a_2}{\partial z_2} \frac{\partial z_2}{\partial w_1} = (a_4 - y)g'(z_4)w_3g'(z_3)w_1g'(z_2)x$$

$$\frac{\partial C}{\partial b_1} = \frac{\partial C}{\partial a_2} \frac{\partial a_2}{\partial z_2} \frac{\partial z_2}{\partial b_1} = (a_4 - y)g'(z_4)w_3g'(z_3)w_1g'(z_2)$$

So we have obtained the six partial derivatives by back-propagation!

<https://www.youtube.com/watch?v=tleHLnjs5U8>

40

Imperial College
London

Feed-Forward Neural Networks

1. From Logistic Regression to Single Neuron Representation
2. Feed-Forward Neural Networks
3. Back-Propagation
4. Batch, Stochastic and Mini-Batch Gradient Descent
5. Conclusion

41

Imperial College
London

Back to Binary Logistic Regression Gradient Descent

$$J(\theta) = \frac{1}{m} \sum_{i=1}^m \text{Cost}(h_{\theta}(x^{(i)}), y^{(i)})$$

$$= -\frac{1}{m} \sum_{i=1}^m \left[y^{(i)} \log(h_{\theta}(x^{(i)})) + (1 - y^{(i)}) \log(1 - h_{\theta}(x^{(i)})) \right] = \text{Cross-Entropy}$$

To minimize, just calculate the derivatives $\frac{\partial J(\theta)}{\partial \theta_j}$ for $j = 1 \dots n$ and apply Gradient Descent

$$\frac{\partial J(\theta)}{\partial \theta_j} = \frac{1}{m} \sum_{i=1}^m [h_{\theta}(x^{(i)}) - y^{(i)}] x_j^{(i)}$$

42

Batch vs Stochastic Gradient Descent

"Batch" Gradient Descent uses all m training set data $(x^{(i)}, y^{(i)})_{i=1, \dots, m}$ at each gradient descent iteration. When m is very large (say m is in the 100,000's), the number of calculations is thus very large, and this is just to calculate one gradient on all the parameters $\theta_{ij}^{(l)}$!

Instead of summing over all the training set, then calculating the gradient and optimizing, we can iterate only using the gradient at individual data points. With *Stochastic Gradient Descent*: every iteration works on one data point at a time.

Stochastic gradient descent

1. Randomly shuffle (reorder) training examples
2. Repeat {
 - for $i := 1, \dots, m$ {
 - $\theta_j := \theta_j - \alpha(h_{\theta}(x^{(i)}) - y^{(i)})x_j^{(i)}$
(for every $j = 0, \dots, n$)



43

Mini-Batch Gradient Descent

Rather than the two extremes of *Batch* and *Stochastic Gradient Descent*, one often chooses the intermediate *Mini-Batch Gradient Descent*, where the gradient is calculated on a small subset of data.

If b is Batch size and m is total number of training data:

Mini-batch gradient descent

Say $b = 10, m = 1000$.

Repeat {

for $i = 1, 11, 21, 31, \dots, 991$ {

$$\theta_j := \theta_j - \alpha \frac{1}{10} \sum_{k=i}^{i+9} (h_{\theta}(x^{(k)}) - y^{(k)})x_j^{(k)}$$

(for every $j = 0, \dots, n$)

}

}

44

Imperial College
London

Gradient Descent: Different Ways to use the Data

Batch (also called Full-Batch) Gradient Descent:

Using all m training set data $(x^{(i)}, y^{(i)})_{i=1, \dots, m}$ at each gradient descent iteration

Stochastic Gradient Descent:

Use one single data $(x^{(i)}, y^{(i)})$ at each gradient descent iteration

Mini-Batch Gradient Descent:

Use small number (say a few tens or hundreds) of data $(x^{(i)}, y^{(i)})$ at each gradient descent iteration


An **epoch** is a training iteration over the whole training set. It is thus composed of one single gradient descent iteration in the Batch case, and as many gradient descent iterations as there are training data in the Stochastic Gradient Descent case.

<https://playground.tensorflow.org>

45

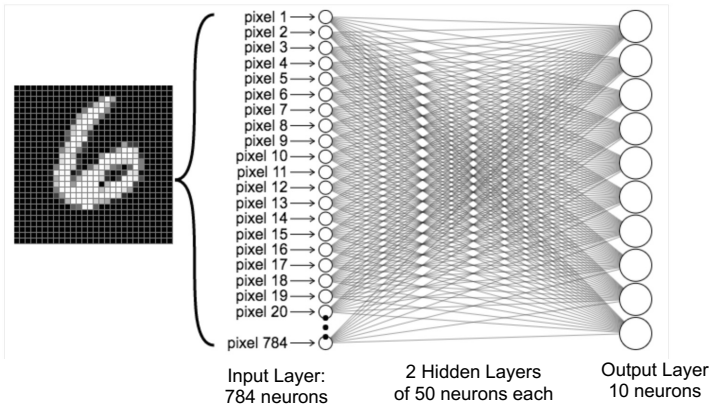
Imperial College
London

Basic Neural Network Training Algorithm (one Epoch)

1. Initialize the training with random parameters $\theta_{jk}^{(l)}$
 2. Calculate $h_{\theta}(x^{(i)})$ for new mini-batch of training set data $x^{(i)}$
 3. Calculate cost/loss function $J(\theta)$
 4. Calculate gradients by back-propagation
 5. Modify parameters $\theta_{jk}^{(l)}$ by gradient descent
- 

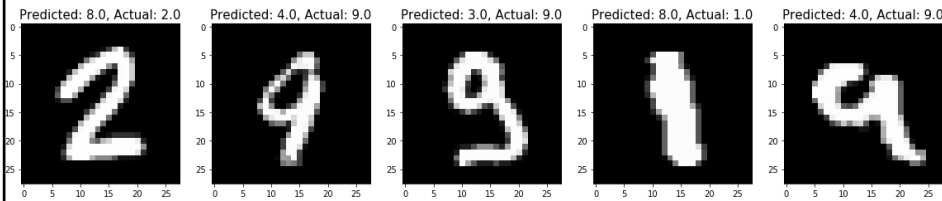
46

MNIST Example of Neural Network Architecture



MNIST Example of Neural Network Architecture

Examples of Misclassified Images



Imperial College
London

Feed-Forward Neural Networks on MNIST

The Results:

On the 60000 Training Images

Mean Accuracy: 0.962

Misclassified Images: 2301 (3.84%)

On the 10000 Test Images

Mean Accuracy: 0.942

Misclassified Images: 579 (5.79%)

49

Imperial College
London

Reminder: Result of Logistic Regression on MNIST

The Results:

On the 60000 Training Images

Mean Accuracy: 0.934

Misclassified Images: 3939 (6.57%)

On the 10000 Test Images

Mean Accuracy: 0.918

Misclassified Images: 818 (8.18%)

50

50

Imperial College
London

Feed-Forward Neural Networks

1. From Logistic Regression to Single Neuron Representation
2. Feed-Forward Neural Networks
3. Back-Propagation
4. Batch, Stochastic and Mini-Batch Gradient Descent
5. Conclusion

51

Imperial College
London

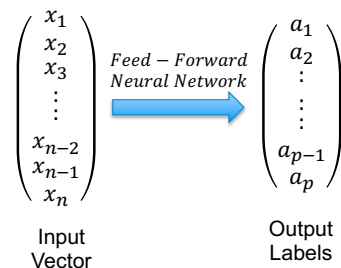
A Simple Way to See Supervised Neural Networks

Suppose we have m pairs of data. Each pair is composed of a vector of dimension n (the features) and a vector of dimension p (the labels).

A neural network is simply a function that maps any vector of dimension n into a (discrete or continuous) vector of dimension p .

In order to calculate the parameters of this function, we train the parameters of the neural network by back-propagation using the m pairs of data as Training Set.

If the n input features are the coordinates of a point in dimension n , the network is an interpolator in R^n .

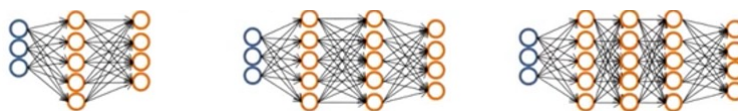


52

Imperial College
London

Which Architecture for my Feed-Forward Neural Network?

How many hidden layers, how many neurons in each hidden layer?



1. Number of input and output neurons is number of features and number of outputs.
2. One hidden layer is a good starting point.
3. If several hidden layers, good to use same number of neurons in each hidden layer.

How can the above choices be made more objective? See next lesson!

53

Imperial College
London

What have we Learnt?

- Logistic Regression as a Single Neuron Neural Network.
- In Feed-Forward Neural Networks, hidden layers provide a « messaging » of the input features in order to make them linearly separable by the output layer.
- The Back-Propagation algorithm allows calculation of all gradients of the loss function according to the trainable parameters so that Gradient Descent can be applied.
- Batch, Stochastic and Mini-Batch Gradient Descent are different ways to organize the data for speeding-up Gradient Calculation and Gradient Descent.
- In simple terms, Supervised Neural Networks just provide a mapping of one input vector into one output vector optimized using the data from the Training Set.

54

Imperial College
London

Exercise 1: Experiment with
<https://playground.tensorflow.org>

**Exercise 2: Gradient Calculations on “Half-Moon”
Dataset, with associated code**

