


Neural Networks

Réda DEHAK
<http://ia.dehak.org>

1






Contents

- Origins/History
- Human Brain / Biological Neuron
- Perceptron
- Neural Networks Topology
- Learning:
 - Perceptron
 - Multi Layers Neural Networks: Backpropagation Algorithm
- Conclusions

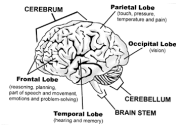
Réda DEHAK 2

2



Origins

- Many machine learning methods inspired by biology, e.g., the (human) brain
- Neural Networks:
 - Algorithms that try to mimic the brain
 - Produce artificial systems capable of complex calculation similar to the human brain



Réda DEHAK 3

3



History

1943 : McCulloch and Pitts proposed the McCulloch-Pitts neuron model

1949 : Hebb published his book The Organization of Behavior, in which the Hebbian learning rule was proposed.

1958 : Rosenblatt introduced the simple single layer networks now called Perceptrons.

1969 : Minsky and Papert's book Perceptrons demonstrated the limitation of single layer perceptrons, and almost the whole field went into hibernation.

1982 : Hopfield published a series of papers on Hopfield networks.

1982 : Kohonen developed the Self-Organising Maps that now bear his name.

1986 : The Back-Propagation learning algorithm for Multi-Layer Perceptrons was rediscovered and the whole field took off again.



1990s : The sub-field of Radial Basis Function Networks was developed.

2000s : The power of Neural Networks, Support Vector Machines, and Bayesian Techniques methods become apparent.

2009 : DNNs outperform other Machine Learning Methods in different competitions


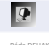
Réda DEHAENE
4

4





Human brain

- Our brain has $\sim 10^{11}$ neurons, each of which communicates (is connected) to $\sim 10^4$ other neurons
- The brain can fire all the neurons in a single step (parallelism).
- The human brain is extremely energy efficient, using approximately 10^{-16} joules per operation per second.
- Brains have been evolving for tens of millions of years

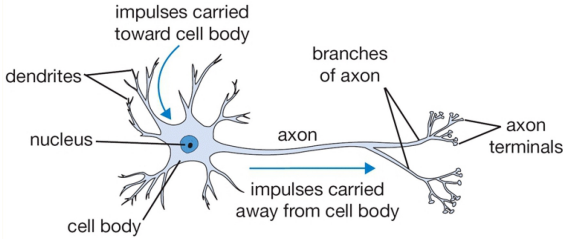
	Processing Elements	Element size	Energy use	Processing speed	Style of computation	Fault Tolerant	Learns	Intelligent, conscious
	10^{14} synapses	10^{-6} m	30 W	100 Hz	parallel, distributed	yes	yes	usually
	10^8 transistors	10^{-6} m	30 W (CPU)	10^8 Hz	serial, centralized	no	a little	not (yet)

Réda DEHAENE
5

5



Biological Neuron



impulses carried toward cell body

dendrites

nucleus

cell body

axon

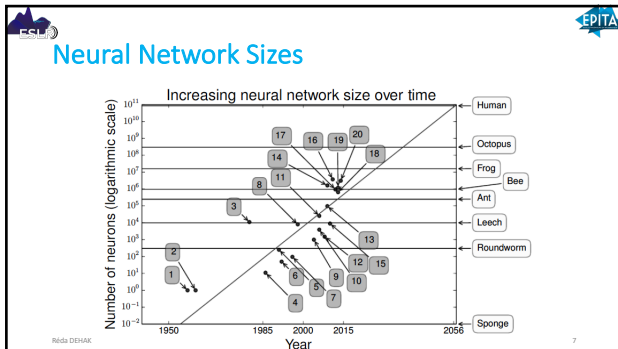
branches of axon

axon terminals

impulses carried away from cell body

Réda DEHAENE
6

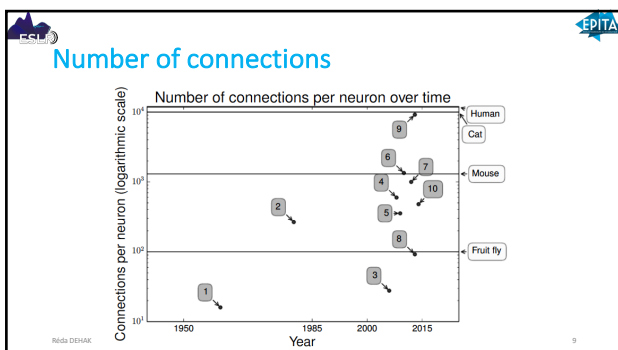
6





7

1. Perceptron (Rosenblatt, 1958, 1962)
 2. Adaptive Linear Element (Widrow and Hoff, 1960)
 3. Neocognitron (Fukushima, 1980)
 4. Early backpropagation network (Rumelhart *et al.*, 1986b)
 5. Recurrent neural network for speech recognition (Robinson and Fallside, 1991)
 6. Multilayer perceptron for speech recognition (Bengio *et al.*, 1991)
 7. Mean field sigmoid belief network (Saul *et al.*, 1996)
 8. LeNet-5 (LeCun *et al.*, 1998b)
 9. Echo state network (Jaeger and Haas, 2004)
 10. Deep belief network (Hinton *et al.*, 2006)
 11. GPU-accelerated convolutional network (Chellapilla *et al.*, 2006)
 12. Deep Boltzmann machines (Salakhutdinov and Hinton, 2009a)
 13. GPU-accelerated deep belief network (Raina *et al.*, 2009)
 14. Unsupervised convolutional network (Jarrett *et al.*, 2009b)
 15. GPU-accelerated multilayer perceptron (Ciresan *et al.*, 2010)
 16. OMP-1 network (Coates and Ng, 2011)
 17. Distributed autoencoder (Le *et al.*, 2012)
 18. Multi-GPU convolutional network (Krizhevsky *et al.*, 2012a)
 19. COTS HPC unsupervised convolutional network (Coates *et al.*, 2013)
 20. GoogleLeNet (Szegedy *et al.*, 2014a)
- Réda DENAE

8





9

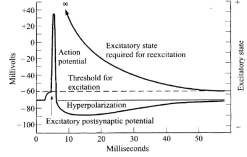
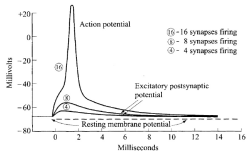
1. Adaptive Linear Element (Widrow and Hoff, 1960)
2. Neocognitron (Fukushima, 1980)
3. GPU-accelerated convolutional network (Chellapilla *et al.*, 2006)
4. Deep Boltzmann machines (Salakhutdinov and Hinton, 2009a)
5. Unsupervised convolutional network (Jarrett *et al.*, 2009b)
6. GPU-accelerated multilayer perceptron (Ciresan *et al.*, 2010)
7. Distributed autoencoder (Le *et al.*, 2012)
8. Multi-GPU convolutional network (Krizhevsky *et al.*, 2012a)
9. COTS HPC unsupervised convolutional network (Coates *et al.*, 2013)
10. GoogLeNet (Szegedy *et al.*, 2014a)

Réda DEHAË
10

10



Biological Neuron Activation

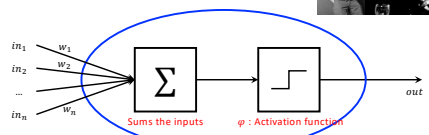
From: Principles of Neurocomputing for Science & Engineering, Hum & Karttunen, McGraw-Hill
From: Principles of Neurocomputing for Science & Engineering, Hum & Karttunen, McGraw-Hill

Réda DEHAË
11

11

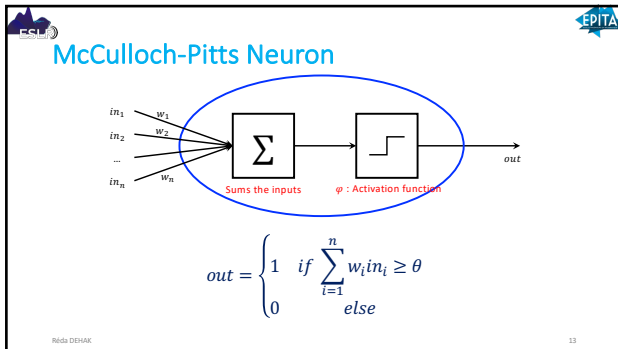
McCulloch-Pitts Neuron



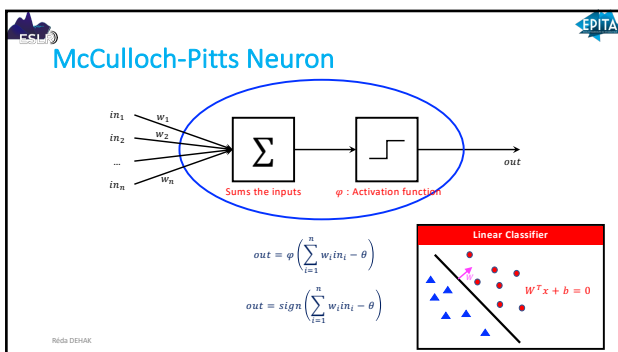
- A set of synapses (i.e. connections) brings in activations from other neurons.
- A processing unit sums the inputs, and then applies a non-linear activation function (i.e. squashing/transfer/threshold function).
- An output line transmits the result to other neurons.
- Any number of McCulloch-Pitts neurons can be **connected together in any way we like**

Réda DEHAË
12

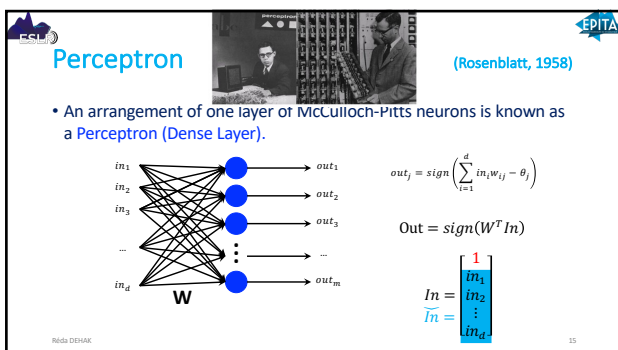
12




13



14



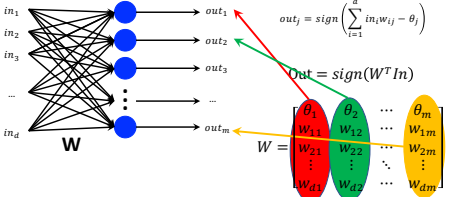
15




Perceptron

(Rosenblatt, 1958)

• An arrangement of one layer of McCulloch-Pitts neurons is known as a **Perceptron (Dense Layer)**.



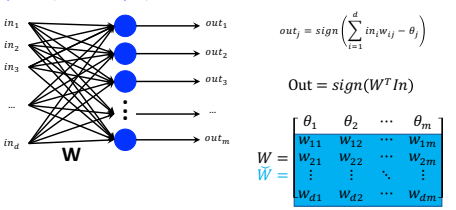
16




Perceptron

(Rosenblatt, 1958)

• An arrangement of one layer of McCulloch-Pitts neurons is known as a **Perceptron (Dense Layer)**.



17

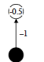


Implementation of Logical operator

NOT

in	out
0	1
1	0

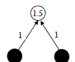
Thresholds \Rightarrow (-0.5)
Weights \Rightarrow -1



AND

in ₁	in ₂	out
0	0	0
0	1	0
1	0	0
1	1	1

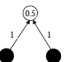
Thresholds \Rightarrow (1.5)
Weights \Rightarrow $1, 1$



OR

in ₁	in ₂	out
0	0	0
0	1	1
1	0	1
1	1	1

Thresholds \Rightarrow (0.5)
Weights \Rightarrow $1, 1$



18

XOR

in_1	in_2	out
0	0	0
0	1	1
1	0	1
1	1	0

19

Architectures, Structures or Topology

- **Single-Layer Feed-forward NNs:** One input layer and one output layer of processing units. No feed-back connections. (For example, a simple Perceptron.)
- **Multi-Layer Feed-forward NNs:** One input layer, one output layer, and one or more hidden layers of processing units. No feed-back connections. The hidden layers sit in between the input and output layers, and are thus hidden from the outside world. (For example, a Multi-Layer Perceptron.)
- **Recurrent NNs:** Any network with at least one feed-back connection. It may, or may not, have hidden units. (For example, a Simple Recurrent Network.)

20

Examples

Single Layer Feed-forward

Single-Layer Perceptron

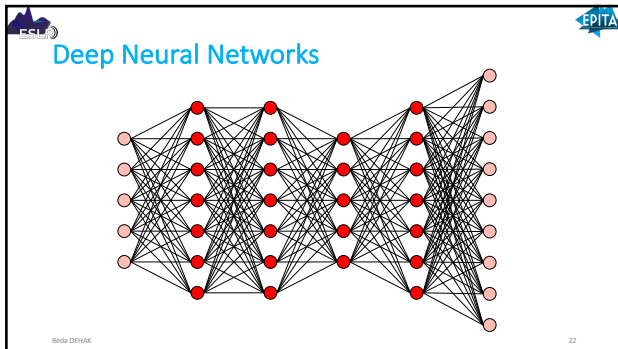
Multi-Layer Feed-forward

Multi-Layer Perceptron

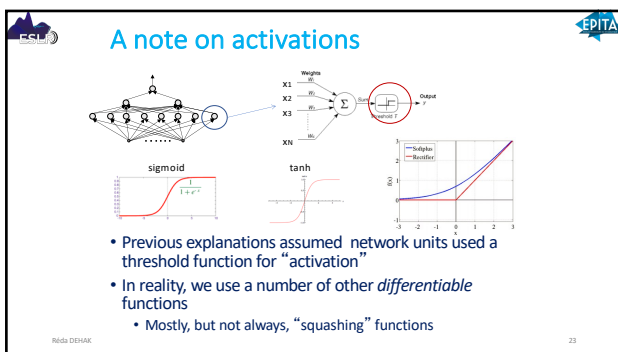
Recurrent Network

Simple Recurrent Network

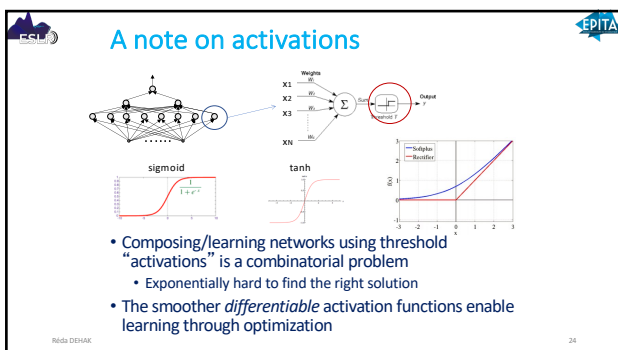
21



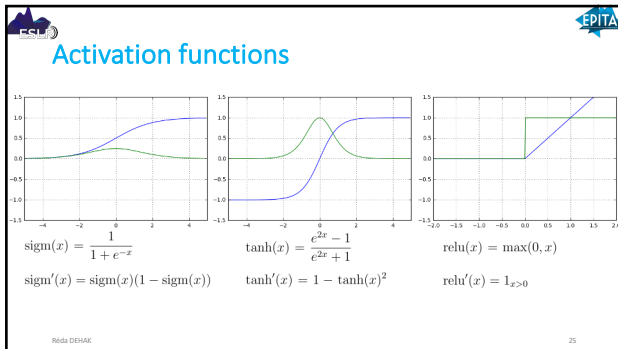
22



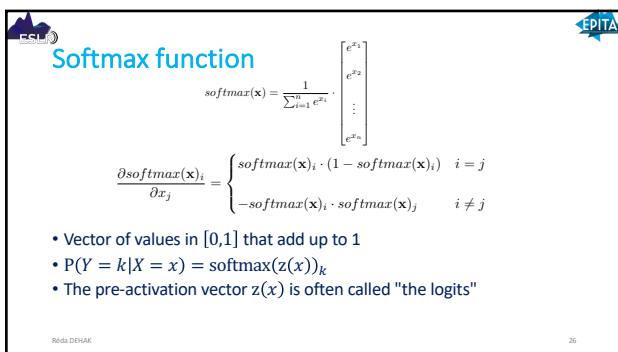
23



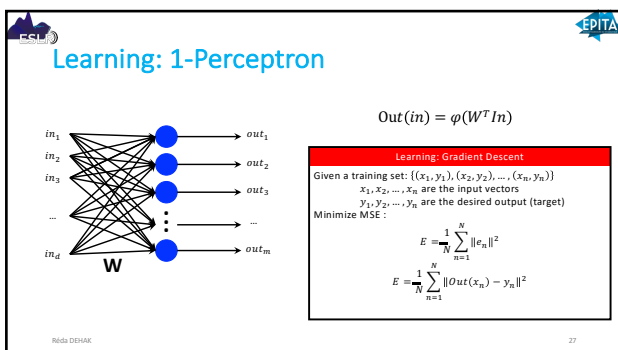
24



25



26



27

Derivatives: Chain Rule

- $z = u(v(x))$ with $v(x)$ is a vector

$$\frac{\partial z}{\partial x_{ij}} = \sum_k \frac{\partial z}{\partial v_k} \frac{\partial v_k}{\partial x_{ij}} = \left(\frac{\partial v}{\partial x_{ij}} \right)^T @ \frac{\partial z}{\partial v}$$
- $z = u(v(x))$ with $v(x)$ is a matrix

$$\frac{\partial z}{\partial x_{ij}} = \sum_k \sum_l \frac{\partial z}{\partial v_{kl}} \frac{\partial v_{kl}}{\partial x_{ij}}$$

28

Learning: 1-Perceptron

- $\frac{\partial E(x)}{\partial e(x)_i} = \frac{\partial (e^T e)}{\partial e(x)_i} = 2 e(x)_i$
- $\frac{\partial e(x)_j}{\partial \varphi(x)_l} = \frac{\partial (\varphi(x) - y_j)}{\partial \varphi(x)_l} = \mathbb{I}_{l=j}$
- $\frac{\partial E(x)}{\partial \varphi(x)_i} = 2(\vec{\mathbb{I}}_i)^T @ e(x) = 2e(x)_i$

$\frac{\partial E(x)}{\partial e(x)} = 2e(x)$
 $\frac{\partial e(x)}{\partial \varphi(x)} = \vec{\mathbb{I}}_i$
 $\frac{\partial E(x)}{\partial \varphi(x)} = 2e(x)$

$\vec{\mathbb{I}}_i = \begin{bmatrix} 0 \\ \vdots \\ 0 \\ 1 \\ 0 \\ \vdots \end{bmatrix}$ ith row

29

Learning: 1-Perceptron

- $\frac{\partial \varphi(z(x))_j}{\partial z(x)_l} = \frac{\partial (\varphi(z(x)_l))}{\partial z(x)_l} = \varphi'(z(x)_l) \mathbb{I}_{l=j} = \varphi'(z(x))_j \cdot \mathbb{I}_{l=j}$
- $\frac{\partial E(x)}{\partial z(x)_i} = 2(\varphi'(z(x)) \cdot \vec{\mathbb{I}}_i)^T @ e(x) = 2e(x)_i \varphi'(z(x))_i$

$\frac{\partial E(x)}{\partial z(x)} = 2e(x) \cdot \varphi'(z(x))$

30

Learning: 1-Perceptron

$$\frac{\partial z(x)_k}{\partial W_{ij}} = \frac{\partial (W^T In)_k}{\partial W_{ij}} = \frac{\partial \sum_{d=1}^D W_{dk} In_d}{\partial W_{ij}} = In_i \mathbb{1}_{k=j}$$

$$\frac{\partial z(x)}{\partial W_{ij}} = In_i \vec{\mathbb{1}}_j$$

$$\frac{\partial E(x)}{\partial W_{ij}} = 2(In_i \vec{\mathbb{1}}_j)^T @ (e(x) \cdot \varphi'(z(x)))$$

$$\frac{\partial E(x)}{\partial W_{ij}} = 2In_i (e(x) \cdot \varphi'(z(x)))_j$$

Réda DEHAENE 31

31

Learning: 1-Perceptron

$$\frac{\partial z(x)_k}{\partial W_{ij}} = \frac{\partial (W^T In)_k}{\partial W_{ij}} = \frac{\partial \sum_{d=1}^D W_{dk} In_d}{\partial W_{ij}} = In_i \mathbb{1}_{k=j}$$

$$\frac{\partial z(x)}{\partial W_{ij}} = In_i \vec{\mathbb{1}}_j$$

$$\frac{\partial E(x)}{\partial W_{ij}} = 2(In_i \vec{\mathbb{1}}_j)^T @ (e(x) \cdot \varphi'(z(x)))$$

$$\frac{\partial E(x)}{\partial W} = 2In @ (e(x) \cdot \varphi'(z(x)))^T$$

Réda DEHAENE 32

32

Learning: 1-Perceptron

$$\frac{\partial z(x)_k}{\partial W_{ij}} = \frac{\partial (W^T In)_k}{\partial W_{ij}} = \frac{\partial \sum_{d=1}^D W_{dk} In_d}{\partial W_{ij}} = In_i \mathbb{1}_{k=j}$$

$$\frac{\partial z(x)}{\partial W_{ij}} = In_i \vec{\mathbb{1}}_j$$

$$\frac{\partial E(x)}{\partial W_{ij}} = 2(In_i \vec{\mathbb{1}}_j)^T @ (e(x) \cdot \varphi'(z(x)))$$

$$\frac{\partial E(x)}{\partial W} = In @ \left(\frac{\partial E(x)}{\partial z(x)} \right)^T$$

Réda DEHAENE 33

33

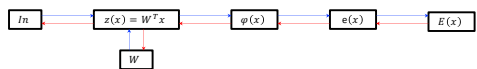
Learning Multi Layers Neural Networks

Backpropagation Algorithm

34

Learning: 2-Multilayers Neural Networks

Backpropagation Algorithm



$$\bullet \frac{\partial z(x)_k}{\partial \tilde{I}n_i} = \frac{\partial (W^T \tilde{I}n)_k}{\partial \tilde{I}n_i} = \frac{\partial \sum_{d=1}^D W_{dk} \tilde{I}n_d}{\partial \tilde{I}n_i} = \tilde{W}_{ik} \quad \frac{\partial z(x)}{\partial \tilde{I}n_i} = (\tilde{W}_{i\bullet})^T$$

$$\frac{\partial E(x)}{\partial \tilde{I}n_i} = 2\tilde{W}_{i\bullet} @ (e(x) \cdot \varphi'(z(x)))$$

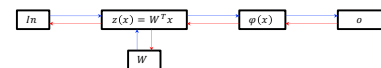
$$\frac{\partial E(x)}{\partial \tilde{I}n} = \tilde{W} @ \frac{\partial E(x)}{\partial z(x)}$$

Réda DEHAENE 35

35

Learning: 2-MultiLayers Neural Networks

Backpropagation Algorithm



$$\bullet \frac{\partial o_j}{\partial z(x)_i} = \frac{\partial (\varphi(z(x)))_j}{\partial z(x)_i} = \varphi'(z(x))_{\parallel_{i=j}} = \varphi'(z(x))_j \cdot \mathbb{I}_{i=j}$$

$$\frac{\partial o}{\partial z(x)_i} = \varphi'(z(x)) \cdot \mathbb{I}_i$$

$$\frac{\partial E(x)}{\partial z(x)_i} = (\varphi'(z(x)) \cdot \mathbb{I}_i)^T @ \frac{\partial E(x)}{\partial o} = \left(\frac{\partial E(x)}{\partial o} \right)_i \varphi'(z(x))_i$$

Réda DEHAENE 36

36

Learning: 2-MultiLayers Neural Networks
Backpropagation Algorithm

$$\frac{\partial o_j}{\partial z(x)_i} = \frac{\partial (\varphi(z(x)))_j}{\partial z(x)_i} = \varphi'(z(x))_{ij} = \varphi'(z(x))_i \cdot \mathbb{I}_{i=j}$$

$$\frac{\partial E}{\partial z(x)} = \left(\frac{\partial E(x)}{\partial o} \cdot \varphi'(z(x)) \right)$$

$$\frac{\partial E}{\partial z(x)} = \left(\frac{\partial E(x)}{\partial I_{n+1}} \cdot \varphi'(z(x)) \right)$$

Réda DEHAË

37

Learning: 2-MultiLayers Neural Networks
Backpropagation Algorithm

$$\frac{\partial z(x)_k}{\partial w_{ij}} = \frac{\partial (w^T I_n)_k}{\partial w_{ij}} = \frac{\partial \sum_{d=1}^D w_{dk} I_{nd}}{\partial w_{ij}} = I_{ni} \mathbb{I}_{k=j}$$

$$\frac{\partial E(x)}{\partial w_{ij}} = 2(I_{ni} \vec{\mathbb{I}}_j)^T @ \left(\frac{\partial E(x)}{\partial z(x)} \right)$$

$$\frac{\partial z(x)}{\partial w_{ij}} = I_{ni} \vec{\mathbb{I}}_j$$

$$\frac{\partial E(x)}{\partial W} = I_n @ \left(\frac{\partial E(x)}{\partial z(x)} \right)^T$$

Réda DEHAË

38

Learning: 2-MultiLayers Neural Networks
Backpropagation Algorithm

$$\frac{\partial z(x)_k}{\partial I_{ni}} = \frac{\partial (w^T I_n)_k}{\partial I_{ni}} = \frac{\partial \sum_{d=1}^D w_{dk} I_{nd}}{\partial I_{ni}} = \tilde{W}_{ik}$$

$$\frac{\partial E(x)}{\partial I_{ni}} = 2 \tilde{W}_{i \cdot} @ \left(\varphi'(z(x)) \right)$$

$$\frac{\partial z(x)}{\partial I_{ni}} = (\tilde{W}_{i \cdot})^T$$

$$\frac{\partial E(x)}{\partial I_n} = \tilde{W} @ \frac{\partial E(x)}{\partial z(x)}$$

Réda DEHAË

39

Learning: 2-MultiLayers Neural Networks
Backpropagation Algorithm

For a Multilayer NN with L Layers

- First Step: Forward Propagation of the input**
 - Compute and store the activation and the output of each layer in order from the input layer to the output layer

$$In^{(l)} = \begin{bmatrix} 1 \\ o^{(l-1)} \end{bmatrix}$$

$$z^{(l)} = W^{(l)T} in^{(l)}$$

$$o^{(l)} = \varphi(z^{(l)})$$

Réda DENAE 40

40

Learning: 2-MultiLayers Neural Networks
Backpropagation Algorithm

For a Multilayers NN with L Layers

- Second Step: Backward Propagation of the gradient**
 - Compute and store the gradient of each layer in order from the output layer to the input layer

Compute : $\frac{\partial E(x)}{\partial o^{(l)}} = 2e$ with $e = o^{(l)} - y$ in the case of **MSE**

and $\delta^{(L)} = \frac{\partial E(x)}{\partial z(x)^{(L)}} = \frac{\partial E(x)}{\partial o^{(L)}} \cdot \varphi'(z(x)^{(L)})$

Iterate : $\frac{\partial E(x)}{\partial W^{(l)}} = In^{(l)} @ (\delta^{(l)})^T$

$$\delta^{(l-1)} = \frac{\partial E(x)}{\partial z(x)^{(l-1)}} = \varphi'(z(x)^{(l-1)}) \cdot (\tilde{W}^{(l)} @ \delta^{(l)})$$

Réda DENAE 41


41

Loss Function

- Regression:**
 - Use MSE cost function with a Identity activation function for output layer.
$$E = \frac{1}{N} \sum_{i=1}^N \|o^{(L)}(x_i) - y_i\|^2 \quad \frac{\partial E(x)}{\partial o^{(L)}} = 2(o^{(L)} - y)$$
- Classification:**
 - Use Cross Entropy cost function with sigmoid (in the case of binary classifier) or softmax (in the case of multiclass classifier) activation function for output layer
$$E = -\frac{1}{N} \sum_{i=1}^N \sum_{k=1}^K \mathbb{1}_{y_i=k} \log(o^{(L)}(x_i))_k \quad \frac{\partial E(x)}{\partial o^{(L)}} = -\frac{\tilde{y}}{o^{(L)}}$$

Réda DENAE 42

42




Gradient descent variants

- Three variants of gradient descent:
 - **Stochastic gradient descent**: uses **one random sample** of the training dataset to compute the gradient at each iteration
 - **Batch gradient descent**: uses **the full training dataset** to compute the gradient at each iteration
 - **Mini-batch gradient descent**: use **a small random portion** of the training dataset to compute the gradient at each iteration

Réda DENAE 43

43




Heuristic for making the Back Propagation algorithm perform better

It is more of an art than a science

- Input data should be normalized to have approx. same range: standardization or quantile normalization
- Initializing W^h and W^o :
 - Zero is a saddle point: no gradient, no learning
 - Constant init: hidden units collapse by symmetry
 - Solution: random init, ex: $W \sim \mathcal{N}(0, 0.01)$
 - Better inits:
 - Xavier Glorot and Kaming He
 - orthogonal
- Biases can (should) be initialized to zero

Réda DENAE 44

44

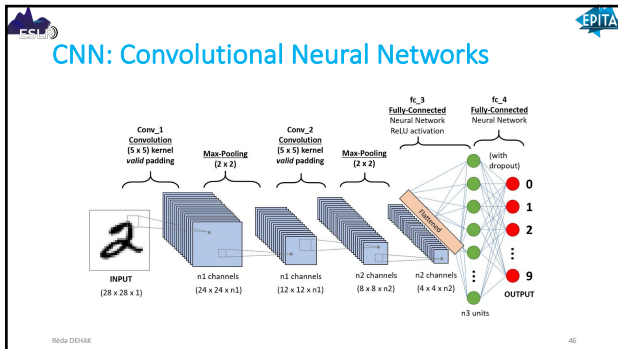


Conclusions:

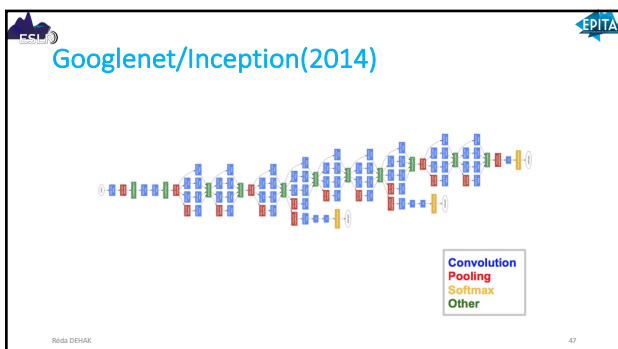
- Neural Networks is a complex function.
- NN Learning is based on Gradient Descent and Backpropagation algorithm.
- New Regulation algorithms were proposed recently: dropout, ...
- New topologies emerge recently.

Réda DENAE 45

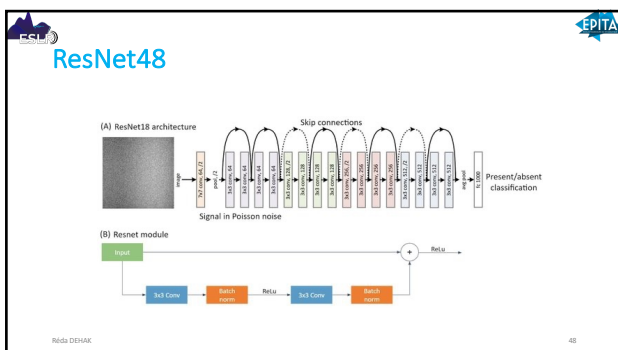
45



46



47



48

DL Applications: Speech to Text

Input: English

Output: CAT

Input: Mandarin

Output: 猫

Legend:
 - Convolution Layer
 - Recurrent Layer
 - Fully Connected Layer

[Baidu 2014]

49

DL Applications: Vision

[Krizhevsky 2012]

[Ciresan et al. 2013]

[Faster R-CNN - Ren 2015]

[NVIDIA dev blog]

50

DL Applications: Vision


[Stanford 2017]

[Nvidia Dev Blog 2017]


[FaceNet - Google 2015]

[Facial landmark detection CUHK 2014]

51



DL Applications: NLP



Translate


From: English -> French Target language:

deep data
deep learning


[Google Translate System - 2016]

Structure

The structure extends from about 10kms to about 100kms in altitude.



Network



Réda DENAEK

[Encher 2016.1]

52

52

[illegible]






[Google Inbox Smart Reply]
ChatGPT, MISTRAL, ...



[Amazon Echo / Alexa]

53



[illegible]

The diagram illustrates a sequence-to-sequence model using a Recurrent Neural Network (RNN) for image captioning. The input is an image of a man playing guitar, and the output is the caption "There is black shirt is playing guitar". The model consists of an encoder (RNN) that processes the image and a decoder (RNN) that generates the caption. The encoder processes the image and outputs a sequence of hidden states h_1, h_2, \dots, h_T . The decoder takes the final hidden state h_T and the previous hidden state h_{t-1} as input and generates the next word in the caption. The process is repeated until the end of the caption is reached.

54



DL Applications: Image translation



original

Isabelle
(21.94801/6423)

SRGAN
(73.44180/7777)



SRGAN
(70.94801/6423)

[Ledig 2016]


Réda DENAE

55

55



DL Applications: Generative Models





Sampled celebrities [Nvidia 2017]

Réda DENAE

56

56



DL Applications: Generative Models

Text description

This bird is blue with white wings that are brown and has a very short beak

This bird has wings that are brown and has a yellow belly

A white bird with a black crown and yellow beak








This bird is white, black, and brown in color, with a brown beak

The bird has small beak, with reddish brown crown and gray belly








This is a small, black bird with a white breast and white on the wingbars.

This bird is white black and yellow in color, with a short black beak

Stage-I images



Stage-II images



StackGAN v2 [Zhang 2017]

Réda DENAE

57

57

19

DL Applications: Generative Models

Output

Hidden Layer

Hidden Layer

Hidden Layer

Input

Sound generation with WaveNet
[DeepMind 2017]

Réda DENAE 58

58

DL Applications: Generative Models

- Guess which one is generated?

- Tacotron 2 Natural TTS Synthesis by Conditioning WaveNet on Mel Spectrogram Predictions, 2017

Réda DENAE 59

59

DL Applications: Games

[Deepmind AlphaGo / Zero 2017]

[Atari Games - DeepMind 2016]

[Starcraft 2 for AI research]

Réda DENAE 60

60
