

What if I turn my project into a shared task participation?

Not a bad idea

Shared task

- Organised to tackle problems that are challenging to be addressed by a single (group of) researcher(s)

What if I turn my project into a shared task participation?

Not a bad idea

Shared task

- Organised to tackle problems that are challenging to be addressed by a single (group of) researcher(s)
- They are challenging by definition: they propose non-trivial problems, still they are specific and doable

What if I turn my project into a shared task participation?

Not a bad idea

Shared task

- Organised to tackle problems that are challenging to be addressed by a single (group of) researcher(s)
- They are challenging by definition: they propose non-trivial problems, still they are specific and doable
- They provide an evaluation framework: dataset + evaluation metrics (+ submission infrastructure)

What if I turn my project into a shared task participation?

Not a bad idea

Shared task

- Organised to tackle problems that are challenging to be addressed by a single (group of) researcher(s)
- They are challenging by definition: they propose non-trivial problems, still they are specific and doable
- They provide an evaluation framework: dataset + evaluation metrics (+ submission infrastructure)
- They give a deadline: so one does not end up tweaking the model forever

What if I turn my project into a shared task participation?

Not a bad idea

Shared task

- Organised to tackle problems that are challenging to be addressed by a single (group of) researcher(s)
- They are challenging by definition: they propose non-trivial problems, still they are specific and doable
- They provide an evaluation framework: dataset + evaluation metrics (+ submission infrastructure)
- They give a deadline: so one does not end up tweaking the model forever
- Visit the lesson website to see some interesting alternatives

What if I turn my project into a shared task participation?

+ You could turn a COLI star!

L'algoritmo di Arianna: "Così su Twitter do la caccia ai post contro le donne"

di Emanuela Giampaoli



La Repubblica (Bologna); 31 March, 2021



RAI News 24; 8 April, 2021

92586 Computational Linguistics

Lesson 12. “More than One” Neuron

Alberto Barrón-Cedeño

Alma Mater Studiorum-Università di Bologna

a.barron@unibo.it

@albarron_

30/03/2022



Previously

- The perceptron
- Intro to neural networks

Table of Contents

1 Backpropagation (brief)

2 Keras

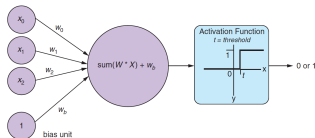
3 Some Guidelines

Chapter 5 of ?

Backpropagation (brief)

Weight Updating

Learning in a “simple” perceptron¹ vs a fully-connected network

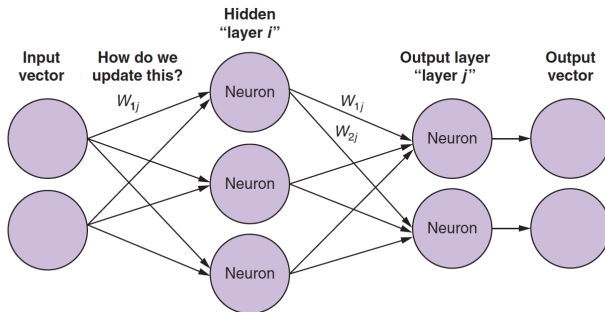
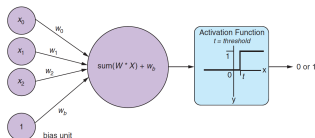


(?, p. 158, 168)

¹Remember: aka linear regression

Weight Updating

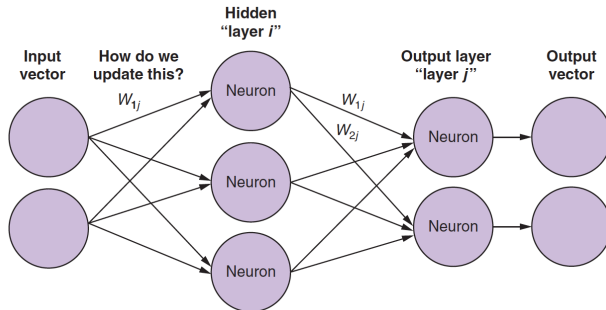
Learning in a “simple” perceptron¹ vs a fully-connected network



(?, p. 158, 168)

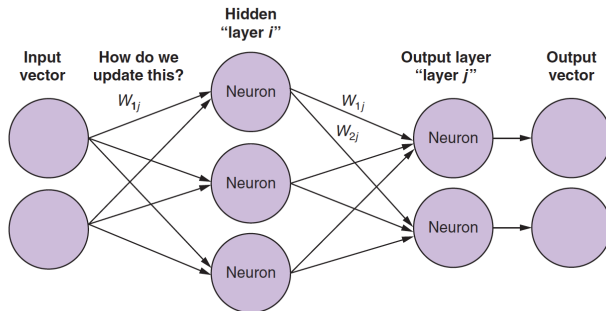
¹Remember: aka linear regression

Backpropagation (of the errors)



²Notice that the first W_{1j} should be W_{1i}

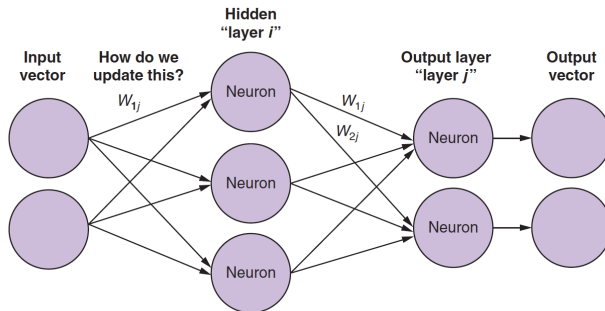
Backpropagation (of the errors)



- The error is computed on the output vector
- How much error did W_{1i} "contribute"?²

²Notice that the first W_{1j} should be W_{1i}

Backpropagation (of the errors)



- The error is computed on the output vector
- How much error did W_{1i} "contribute"?²
- "Path": $W_{1i} \rightarrow [W_{1j}, W_{2j}] \rightarrow output$

²Notice that the first W_{1j} should be W_{1i}

Backpropagation (of the errors)

A better activation function

Step function: $f(\vec{x}) = \begin{cases} 1 & \text{if } \sum_{i=0}^n x_i w_i > \text{threshold} \\ 0 & \text{otherwise} \end{cases}$

³The change of the output is not proportional to the change of the input.

Backpropagation (of the errors)

A better activation function

Step function: $f(\vec{x}) = \begin{cases} 1 & \text{if } \sum_{i=0}^n x_i w_i > \text{threshold} \\ 0 & \text{otherwise} \end{cases}$

Sigmoid function: non-linear³ and continuously differentiable

$$S(x) = \frac{1}{1 + e^{-x}} \quad (1)$$

³The change of the output is not proportional to the change of the input.


Backpropagation (of the errors)

A better activation function

Step function: $f(\vec{x}) = \begin{cases} 1 & \text{if } \sum_{i=0}^n x_i w_i > \text{threshold} \\ 0 & \text{otherwise} \end{cases}$

Sigmoid function: non-linear³ and continuously differentiable

$$S(x) = \frac{1}{1 + e^{-x}} \quad (1)$$

 Let us see

³The change of the output is not proportional to the change of the input. 


Backpropagation (of the errors)

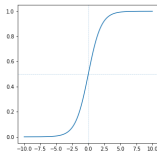
A better activation function

Step function: $f(\vec{x}) = \begin{cases} 1 & \text{if } \sum_{i=0}^n x_i w_i > \text{threshold} \\ 0 & \text{otherwise} \end{cases}$

Sigmoid function: non-linear³ and continuously differentiable

$$S(x) = \frac{1}{1 + e^{-x}} \quad (1)$$

 Let us see



Non-linear → model non-linear relationships
Continuously differentiable → partial derivatives wrt various variables to
_____ update the weights

³The change of the output is not proportional to the change of the input. 

Backpropagation

Differentiating to adjust

Squared error (in (? , p. 171) they say this is MSE; wrong)

$$SE = (y - f(x))^2 \quad (2)$$

Backpropagation

Differentiating to adjust

Squared error (in (? , p. 171) they say this is MSE; wrong)

$$SE = (y - f(x))^2 \quad (2)$$

Mean squared error

$$MSE = \frac{1}{n} \sum_{i=1}^n (y - f(x))^2 \quad (3)$$

Backpropagation

Differentiating to adjust

Squared error (in (? , p. 171) they say this is MSE; wrong)

$$SE = (y - f(x))^2 \quad (2)$$

Mean squared error

$$MSE = \frac{1}{n} \sum_{i=1}^n (y - f(x))^2 \quad (3)$$

Calculus chain rule

$$f(g(x))' = F'(x) = f'(g(x)) g'(x) \quad (4)$$

Backpropagation

Differentiating to adjust

Squared error (in (? , p. 171) they say this is MSE; wrong)

$$SE = (y - f(x))^2 \quad (2)$$

Mean squared error

$$MSE = \frac{1}{n} \sum_{i=1}^n (y - f(x))^2 \quad (3)$$

Calculus chain rule

$$f(g(x))' = F'(x) = f'(g(x)) g'(x) \quad (4)$$

With (4) we can find the derivative of the actfunct \forall neuron wrt its input.

Backpropagation

Differentiating to adjust

Squared error (in (? , p. 171) they say this is MSE; wrong)

$$SE = (y - f(x))^2 \quad (2)$$

Mean squared error

$$MSE = \frac{1}{n} \sum_{i=1}^n (y - f(x))^2 \quad (3)$$

Calculus chain rule

$$f(g(x))' = F'(x) = f'(g(x)) g'(x) \quad (4)$$

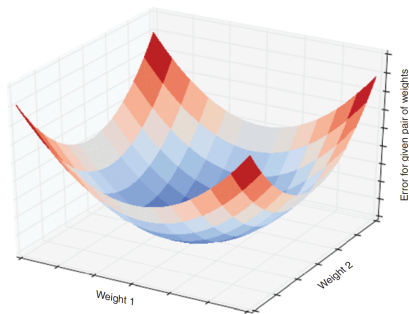
With (4) we can find the derivative of the actfunct \forall neuron wrt its input.

Plain words: find the contribution of a weight to the error and adjust it!

(no further math)

Backpropagation (of the errors)

~Gradient descent: minimising the error

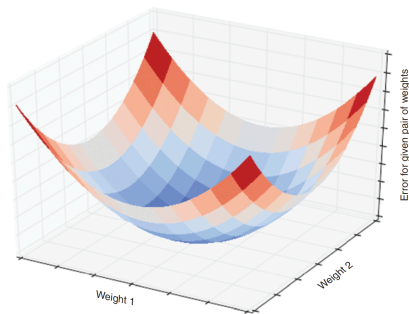


Convex error curve

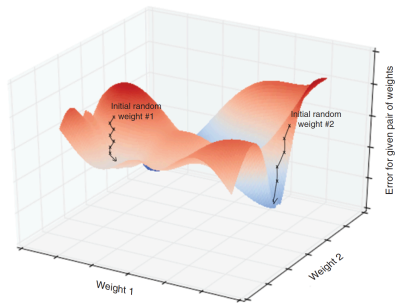
(?, p. 173–174)

Backpropagation (of the errors)

~Gradient descent: minimising the error



Convex error curve



Non-convex error curve

(?, p. 173–174)

Addressing Local minima

Batch learning

- Aggregate the error for the batch
- Update the weight at the end

Addressing Local minima

Batch learning

- Aggregate the error for the batch
- Update the weight at the end
- → hard to find global minimum

Addressing Local minima

Batch learning

- Aggregate the error for the batch
- Update the weight at the end
- → hard to find global minimum

Stochastic gradient descent

- Look at the error for each single instance
- Update the weights right away

Addressing Local minima

Batch learning

- Aggregate the error for the batch
- Update the weight at the end
- → hard to find global minimum

Stochastic gradient descent

- Look at the error for each single instance
- Update the weights right away
- → more likely to make it to the global minimum

Addressing Local minima

Batch learning

- Aggregate the error for the batch
- Update the weight at the end
- → hard to find global minimum

Stochastic gradient descent

- Look at the error for each single instance
- Update the weights right away
- → more likely to make it to the global minimum

Mini-batch

- Much smaller batch, combining the best of the two worlds
- → Fast as batch, resilient as stochastic gradient descent

Addressing Local minima

Batch learning

- Aggregate the error for the batch
- Update the weight at the end
- → hard to find global minimum

Stochastic gradient descent

- Look at the error for each single instance
- Update the weights right away
- → more likely to make it to the global minimum

Mini-batch

- Much smaller batch, combining the best of the two worlds
- → Fast as batch, resilient as stochastic gradient descent

Important parameter: learning rate α

A parameter to define at what extent should we “correct” the error

Keras

Some Available Libraries

There are many high- and low-level libraries in many languages

- **PyTorch**

Community-driven; <https://pytorch.org/>

- **Theano**

MILA (UdeM); www.deeplearning.net/software/theano/⁴

- **TensorFlow**

Google Brain; <https://www.tensorflow.org/>

- **Others**

⁴Non active

Some Available Libraries

There are many high- and low-level libraries in many languages

- **PyTorch**

Community-driven; <https://pytorch.org/>

- **Theano**

MILA (UdeM); www.deeplearning.net/software/theano/⁴

- **TensorFlow**

Google Brain; <https://www.tensorflow.org/>

- **Others**

We will use **Keras**; <https://keras.io/>

⁴Non active

What is Keras

- High-level wrapper with an accessible API for Python

What is Keras

- High-level wrapper with an accessible API for Python
- Gives access to three alternative backends
 - ▶ Theano
 - ▶ TensorFlow
 - ▶ CNTK (MS)

Keras


Logical exclusive OR in Keras

input		output
0	0	0
0	1	1
1	0	1
1	1	0

Keras

Logical exclusive OR in Keras


input		output
0	0	0
0	1	1
1	0	1
1	1	0

 Let us see

Keras

Logical exclusive OR in Keras

input		output
0	0	0
0	1	1
1	0	1
1	1	0


 Let us see

- First dense layer
 - ▶ 2 inputs, 10 neurons
 - ▶ 30 parameters

Keras

Logical exclusive OR in Keras

input		output
0	0	0
0	1	1
1	0	1
1	1	0


 Let us see

- First dense layer
 - ▶ 2 inputs, 10 neurons
 - ▶ 30 parameters
 - ▶ $2 \times 10 \rightarrow 20$

Keras

Logical exclusive OR in Keras

input		output
0	0	0
0	1	1
1	0	1
1	1	0


 Let us see

- First dense layer
 - ▶ 2 inputs, 10 neurons
 - ▶ 30 parameters
 - ▶ $2 \times 10 \rightarrow 20$
 - ▶ But we also have the bias! That's 10 more weights
- Second dense layer

Keras

Logical exclusive OR in Keras

input		output
0	0	0
0	1	1
1	0	1
1	1	0


 Let us see

- First dense layer
 - ▶ 2 inputs, 10 neurons
 - ▶ 30 parameters
 - ▶ $2 \times 10 \rightarrow 20$
 - ▶ But we also have the bias! That's 10 more weights
- Second dense layer
 - ▶ 10 inputs, 1 neuron
 - ▶ 11 parameters

Keras


Logical exclusive OR in Keras

input		output
0	0	0
0	1	1
1	0	1
1	1	0

 Let us see

- First dense layer
 - ▶ 2 inputs, 10 neurons
 - ▶ 30 parameters
 - ▶ $2 \times 10 \rightarrow 20$
 - ▶ But we also have the bias! That's 10 more weights
- Second dense layer
 - ▶ 10 inputs, 1 neuron
 - ▶ 11 parameters

Now we can compile the model

 Let us see

Some Guidelines

Design Decisions

Activation functions

Sigmoid

RelU Rectified linear unit (and variations)

tanh Hyperbolic tangent

Figure source: (?)

Design Decisions

Activation functions

Sigmoid

ReLU Rectified linear unit (and variations)

tanh Hyperbolic tangent

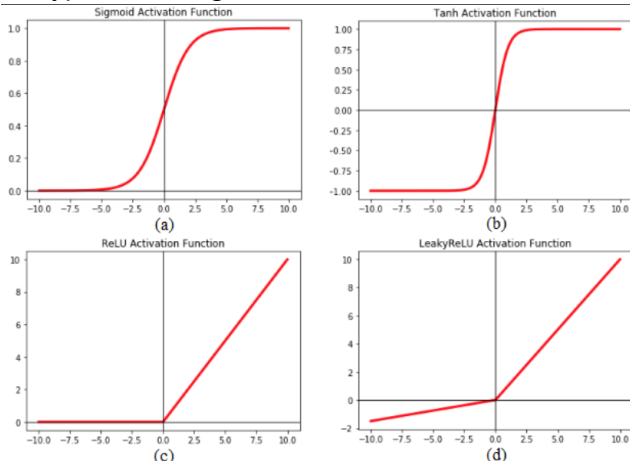


Figure source: (?)

Design Decisions

Activation functions

- Sigmoid
- ReLU (rectified linear unit)
- tanh (hyperbolic tangent)

Learning rate

- Choosing one in advance
- Use **momentum** to perform dynamic adjustments

Design Decisions

Activation functions

- Sigmoid
- ReLU (rectified linear unit)
- tanh (hyperbolic tangent)

Learning rate

- Choosing one in advance
- Use **momentum** to perform dynamic adjustments

Dropout

- Ignore randomly-chosen weights in a training pass to prevent overfitting

Design Decisions

Activation functions

- Sigmoid
- ReLU (rectified linear unit)
- tanh (hyperbolic tangent)

Learning rate

- Choosing one in advance
- Use **momentum** to perform dynamic adjustments

Dropout

- Ignore randomly-chosen weights in a training pass to prevent overfitting

Regularisation

- Dampen a weight from growing/shrinking too far from the rest to prevent overfitting

Normalisation

Example House classification.

Input number of bedrooms, last selling price

Output Likelihood of selling

Vector `input_vec` = `[2, 90000]`

Normalisation

Example House classification.

Input number of bedrooms, last selling price

Output Likelihood of selling

Vector `input_vec` = `[2, 90000]`

All input dimensions should have comparable values

Ideally, all features should be in the range $[-1, 1]$ or $[0, 1]$

Normalisation

Example House classification.

Input number of bedrooms, last selling price

Output Likelihood of selling

Vector `input_vec` = `[2, 90000]`

All input dimensions should have comparable values

Ideally, all features should be in the range $[-1, 1]$ or $[0, 1]$

Typical normalisation: mean normalization, feature scaling, coefficient of variation

Normalisation

Example House classification.

Input number of bedrooms, last selling price

Output Likelihood of selling

Vector `input_vec` = `[2, 90000]`

All input dimensions should have comparable values

Ideally, all features should be in the range $[-1, 1]$ or $[0, 1]$

Typical normalisation: mean normalization, feature scaling, coefficient of variation

NLP uses TF-IDF, one-hot encoding, word2vec (already normalised!)

References

Kandel, I. and M. Castelli

2020. Transfer learning with convolutional neural networks for diabetic retinopathy image classification. a review. *Applied Sciences*, 10(6).

Lane, H., C. Howard, and H. Hapkem

2019. *Natural Language Processing in Action*. Shelter Island, NY: Manning Publication Co.