

91258 / B0385 Natural Language Processing

Lesson 11. "More than One" Neuron

Alberto Barrón-Cedeño a.barron@unibo.it

03/11/2025

DIT, LM SpecTra

Table of Contents

- 1. Backpropagation (brief)
- 2. Keras
- 3. Some Guidelines

Chapter 5 of Lane et al. (2019)

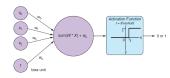
Previously

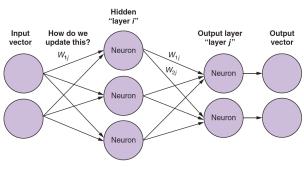
• The perceptron
• Intro to neural networks



Weight Updating

Learning in a "simple" perceptron vs a fully-connected network





(Lane et al. 2019 p. 158-168)

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

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

Let us see



Non-linear → model non-linear relationships

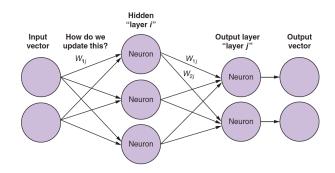
Bounded → constrained to lower and upper boundaries

Continuously differentiable → partial derivatives wrt various variables to

update the weights are possible

A. Barrón-Cedeño DIT, LM SpecTra

Backpropagation (of the errors)



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

²Errata: W_{1i} (between the input and the hidden layer) should be W_{1i}

DIT, LM SpecTra

2025 6 / 19

Backpropagation

Differentiating to adjust

Squared error⁴

$$SE = (y - f(x))^2 \tag{2}$$

Mean squared error

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

Calculus chain rule

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

With (4) we can find the derivative of the actfunct \forall unit wrt its input. Plain words: find the contribution of a weight to the error and adjust it!

(no more math)

⁴In (Lane et al., 2019, p. 171) they say this is MSE; but there is no mean

DIT, LM SpecTra

Backpropagation (of the errors) Caradient descent: minimising the error Supplies to part and part of the error of the er

Keras A. Barrón-Cedeño DIT, LM SpecTra 2025 11 / 19

Addressing Local minima

Batch learning

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

Stochastic gradient descent

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

Mini-batch

- Much smaller batch, combining the best of the two worlds
- ullet 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

A. Barrón-Cedeño

DIT, LM SpecTra

025 10 / 1

Some Popular Libraries

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

PyTorch

Community-driven; https://pytorch.org

TensorFlow

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

Others

We will use Keras; https://keras.io

A. Barrón-Cedeño DIT, LM SpecTra 2025 12 /

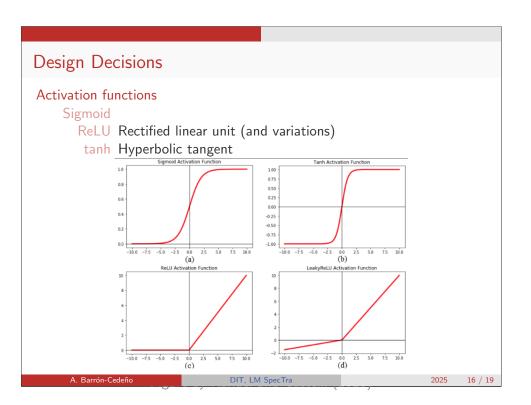
What is Keras A high-level wrapper with an accessible API for Python It gives access to three alternative backends TensorFlow CNTK (MS)

A. Barrón-Cedeño



Keras Logical exclusive OR (XOR) in Keras input output 0 0 1 1 1 1 0 Let us see First dense layer Second dense layer • 2 inputs, 10 units • 10 inputs, 1 unit • 30 parameters • 11 parameters • $2 \times 10 \rightarrow 20$ • But we also have the bias! (10 more weights) Now we can compile the model Let us see

A. Barrón-Cedeño



Design Decisions

Activation functions

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

Learning rate

- Choose 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

A. Barrón-Cedeño

DIT, LM SpecTra

025 17 /

References

Kandel, I. and M. Castelli

A. Barrón-Cedeño

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

DIT, LM SpecTra

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

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

Normalisation

Example House classification.

Input number of bedrooms, last selling price

Output Likelihood of selling

Vector input_vec = [4, 12000]

All input dimensions should have comparable values

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

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

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

Barrón-Cedeño DIT, LM SpecTra 2025 18 / 19