

UNIVERSITÀ DEGLI STUDI DI BERGAMO

Dipartimento di Ingegneria Gestionale, dell'Informazione e della Produzione

Lesson 9.

Neural networks

DATA SCIENCE AND AUTOMATION COURSE

MASTER DEGREE SMART TECHNOLOGY ENGINEERING

TEACHER Mirko Mazzoleni

PLACE University of Bergamo

Outline

1. Neural networks: basic concepts

2. Neural networks: learning

3. Tips & tricks

4. Deep learning



Outline

1. Neural networks: basic concepts

2. Neural networks: learning

3. Tips & tricks

4. Deep learning



Nonlinear hypothesis

Linear methods can be used to learn nonlinear function

• the important thing is that the model is **linear in the parameters**



Logistic regression

 $h_{\theta}(x) = s(\theta_0 + \theta_1 x_1 + \theta_2 x_2 + \theta_3 x_1 x_2 + \theta_4 x_1^2 x_2 + \theta_5 x_1^3 x_2 + \cdots)$

- Logistic regression in **OK if we have a 2-dim feature space** (x_1, x_2) , so q = d 1 = 2. We can compute all the quadratic terms $O(q^2) \approx \frac{q^2}{2}$
- If we have 100 features, we end up with about ≈ 5000 features
- This is **not a good way** to learn nonlinear functions



A Neural Network (NN) is a nonlinear model

- NN originate from trying to mimic the brain
- Very widely used in 80s and 90s, then their popularity diminished in late 90s
- Today they are state-of-the-art technique for many applications (vision, audio, text)



- If a neuron wants to send a message, it sends an electric impulse to other neurons
- Which, in turn, do some computation with the received input to send its messages



Neuron model: logistic unit



UNIVERSITÀ DEGLI STUDI **DI BERGAMO**

di Ingegneria Gestionale, dell'Informazione e della Produzione

Neural networks can be seen as combination of logistic units



- Values in the **hidden layer** are not observed in the training set
- The network computes its own features

$$h_{\Theta}(\mathbf{x}) = s \left(\theta_0 a_0^{(2)} + \theta_1 a_1^{(2)} + \theta_2 a_2^{(2)} + \theta_3 a_3^{(2)} \right)$$

• It is like logistic regression, but it uses $a_0^{(2)}, \dots, a_3^{(2)}$ instead of x_0, \dots, x_3





 $a_i^{(l)}$ = "activation" of unit *i* in layer *l*

 $\Theta^{(l)}$ = matrix of weights controlling function

mapping from layer l to layer l + 1

$$a_{1}^{(2)} = s \left(\Theta_{10}^{(1)} x_{0} + \Theta_{11}^{(1)} x_{1} + \Theta_{12}^{(1)} x_{2} + \Theta_{13}^{(1)} x_{3} \right)$$

$$a_{2}^{(2)} = s \left(\Theta_{20}^{(1)} x_{0} + \Theta_{21}^{(1)} x_{1} + \Theta_{22}^{(1)} x_{2} + \Theta_{23}^{(1)} x_{3} \right)$$

(a) $a_{2}^{(1)} = a \left(A_{20}^{(1)} x_{0} + A_{21}^{(1)} x_{1} + A_{22}^{(1)} x_{2} + A_{23}^{(1)} x_{3} \right)$

$$a_{3}^{(2)} = s \left(\Theta_{30}^{(1)} x_{0} + \Theta_{31}^{(1)} x_{1} + \Theta_{32}^{(1)} x_{2} + \Theta_{33}^{(1)} x_{3} \right)$$

$$h_{\Theta}(\mathbf{x}) = a_1^{(3)} = s \left(\Theta_{10}^{(2)} a_0^{(2)} + \Theta_{11}^{(2)} a_1^{(2)} + \Theta_{12}^{(2)} a_2^{(2)} + \Theta_{13}^{(2)} a_3^{(2)} \right)$$

If network has w_l units in layer l, w_{l+1} units in layer l + 1, then $\Theta^{(j)}$ will be of dimension $w_{l+1} \times (w_l + 1)$ In our example: $\Theta^{(1)} \in \mathbb{R}^{3 \times 4}$ and $\Theta^{(2)} \in \mathbb{R}^{1 \times 4}$





3 units (do not consider bias units) **3 units** (do not consider bias units)

$$a_i^{(l)} =$$
 "activation" of unit *i* in layer *l*

 $\Theta^{(l)}$ = matrix of weights controlling function mapping from layer *l* to layer *l* + 1

$$a_{1}^{(2)} = s \left(\Theta_{10}^{(1)} x_{0} + \Theta_{11}^{(1)} x_{1} + \Theta_{12}^{(1)} x_{2} + \Theta_{13}^{(1)} x_{3} \right)$$

$$a_{2}^{(2)} = s \left(\Theta_{20}^{(1)} x_{0} + \Theta_{21}^{(1)} x_{1} + \Theta_{22}^{(1)} x_{2} + \Theta_{23}^{(1)} x_{3} \right)$$

$$a_{3}^{(2)} = s \left(\Theta_{30}^{(1)} x_{0} + \Theta_{31}^{(1)} x_{1} + \Theta_{32}^{(1)} x_{2} + \Theta_{33}^{(1)} x_{3} \right)$$

$$(\mathbf{x}) = a_{1}^{(3)} = s \left(\Theta_{10}^{(2)} a_{0}^{(2)} + \Theta_{11}^{(2)} a_{1}^{(2)} + \Theta_{12}^{(2)} a_{2}^{(2)} + \Theta_{13}^{(2)} a_{3}^{(2)} \right)$$

If network has w_l units in layer l, w_{l+1} units in layer l + 1, then $\Theta^{(j)}$ will be of dimension $w_{l+1} \times (w_l + 1)$ In our example: $\Theta^{(1)} \in \mathbb{R}^{3 \times 4}$ and $\Theta^{(2)} \in \mathbb{R}^{1 \times 4}$

 h_{Θ}







$$a_i^{(l)} =$$
 "activation" of unit *i* in layer *l*

 $\Theta^{(l)}$ = matrix of weights controlling function mapping from layer *l* to layer *l* + 1

$$a_{1}^{(2)} = s \left(\Theta_{10}^{(1)} x_{0} + \Theta_{11}^{(1)} x_{1} + \Theta_{12}^{(1)} x_{2} + \Theta_{13}^{(1)} x_{3} \right)$$

$$a_{2}^{(2)} = s \left(\Theta_{20}^{(1)} x_{0} + \Theta_{21}^{(1)} x_{1} + \Theta_{22}^{(1)} x_{2} + \Theta_{23}^{(1)} x_{3} \right)$$

$$a_{3}^{(2)} = s \left(\Theta_{30}^{(1)} x_{0} + \Theta_{31}^{(1)} x_{1} + \Theta_{32}^{(1)} x_{2} + \Theta_{33}^{(1)} x_{3} \right)$$

$$(\mathbf{x}) = a_{1}^{(3)} = s \left(\Theta_{10}^{(2)} a_{0}^{(2)} + \Theta_{11}^{(2)} a_{1}^{(2)} + \Theta_{12}^{(2)} a_{2}^{(2)} + \Theta_{13}^{(2)} a_{3}^{(2)} \right)$$

If network has w_l units in layer l, w_{l+1} units in layer l + 1, then $\Theta^{(j)}$ will be of dimension $w_{l+1} \times (w_l + 1)$ In our example: $\Theta^{(1)} \in \mathbb{R}^{3 \times 4}$ and $\Theta^{(2)} \in \mathbb{R}^{1 \times 4}$

 h_{Θ}





bias units)

$$a_i^{(l)}$$
 = "activation" of unit *i* in layer *l*

 $\Theta^{(l)}$ = matrix of weights controlling function mapping from layer *l* to layer *l* + 1

$$a_{1}^{(2)} = s \left(\Theta_{10}^{(1)} x_{0} + \Theta_{11}^{(1)} x_{1} + \Theta_{12}^{(1)} x_{2} + \Theta_{13}^{(1)} x_{3} \right)$$

$$a_{2}^{(2)} = s \left(\Theta_{20}^{(1)} x_{0} + \Theta_{21}^{(1)} x_{1} + \Theta_{22}^{(1)} x_{2} + \Theta_{23}^{(1)} x_{3} \right)$$

$$a_{3}^{(2)} = s \left(\Theta_{30}^{(1)} x_{0} + \Theta_{31}^{(1)} x_{1} + \Theta_{32}^{(1)} x_{2} + \Theta_{33}^{(1)} x_{3} \right)$$

$$(\mathbf{x}) = a_{1}^{(3)} = s \left(\Theta_{10}^{(2)} a_{0}^{(2)} + \Theta_{11}^{(2)} a_{1}^{(2)} + \Theta_{12}^{(2)} a_{2}^{(2)} + \Theta_{13}^{(2)} a_{3}^{(2)} \right)$$

If network has w_l units in layer l, w_{l+1} units in layer l + 1, then $\Theta^{(j)}$ will be of dimension $w_{l+1} \times (w_l + 1)$ In our example: $\Theta^{(1)} \in \mathbb{R}^{3 \times 4}$ and $\Theta^{(2)} \in \mathbb{R}^{1 \times 4}$

 h_{Θ}



bias units)



3 units (do not consider (do not consider bias units) bias units)

$$a_i^{(l)} =$$
 "activation" of unit *i* in layer *l*

 $\Theta^{(l)}$ = matrix of weights controlling function mapping from layer *l* to layer *l* + 1

$$a_{1}^{(2)} = s \left(\Theta_{10}^{(1)} x_{0} + \Theta_{11}^{(1)} x_{1} + \Theta_{12}^{(1)} x_{2} + \Theta_{13}^{(1)} x_{3} \right)$$

$$a_{2}^{(2)} = s \left(\Theta_{20}^{(1)} x_{0} + \Theta_{21}^{(1)} x_{1} + \Theta_{22}^{(1)} x_{2} + \Theta_{23}^{(1)} x_{3} \right)$$

$$a_{3}^{(2)} = s \left(\Theta_{30}^{(1)} x_{0} + \Theta_{31}^{(1)} x_{1} + \Theta_{32}^{(1)} x_{2} + \Theta_{33}^{(1)} x_{3} \right)$$

$$h_{\Theta}(\mathbf{x}) = a_{1}^{(3)} = s \left(\Theta_{10}^{(2)} a_{0}^{(2)} + \Theta_{11}^{(2)} a_{1}^{(2)} + \Theta_{12}^{(2)} a_{2}^{(2)} + \Theta_{13}^{(2)} a_{3}^{(2)} \right)$$

If network has w_l units in layer l, w_{l+1} units in layer l + 1, then $\Theta^{(l)}$ will be of dimension $w_{l+1} \times (w_l + 1)$ In our example: $\Theta^{(1)} \in \mathbb{R}^{3 \times 4}$ and $\Theta^{(2)} \in \mathbb{R}^{1 \times 4}$



Forward computation



Forward
computation
$$x_{1}^{(2)} = s\left(\begin{array}{c} a_{10}^{(1)} x_{0} + \theta_{11}^{(1)} x_{1} + \theta_{12}^{(1)} x_{2} + \theta_{13}^{(1)} x_{3} \right) \Rightarrow a_{1}^{(2)} = s\left(z_{1}^{(2)}\right)$$

$$a_{1}^{(2)} = s\left(\begin{array}{c} a_{10}^{(1)} x_{0} + \theta_{11}^{(1)} x_{1} + \theta_{12}^{(1)} x_{2} + \theta_{13}^{(1)} x_{3} \right) \Rightarrow a_{1}^{(2)} = s\left(z_{1}^{(2)}\right)$$

$$a_{1}^{(2)} = s\left(\begin{array}{c} a_{10}^{(2)} x_{0} + \theta_{11}^{(1)} x_{1} + \theta_{12}^{(2)} x_{2} + \theta_{13}^{(2)} x_{3} \right) \Rightarrow a_{2}^{(2)} = s\left(z_{2}^{(2)}\right)$$

$$a_{2}^{(2)} = s\left(\begin{array}{c} \theta_{20}^{(1)} x_{0} + \theta_{21}^{(1)} x_{1} + \theta_{22}^{(2)} x_{2} + \theta_{23}^{(2)} x_{3} \right) \Rightarrow a_{2}^{(2)} = s\left(z_{2}^{(2)}\right)$$

$$a_{2}^{(2)} = s\left(\begin{array}{c} \theta_{20}^{(1)} x_{0} + \theta_{21}^{(1)} x_{1} + \theta_{22}^{(2)} x_{2} + \theta_{23}^{(2)} x_{3} \right) \Rightarrow a_{2}^{(2)} = s\left(z_{2}^{(2)}\right)$$

$$a_{3}^{(2)} = s\left(\begin{array}{c} \theta_{30}^{(1)} x_{0} + \theta_{31}^{(1)} x_{1} + \theta_{32}^{(1)} x_{2} + \theta_{33}^{(2)} x_{3} \right) \Rightarrow a_{3}^{(2)} = s\left(z_{3}^{(2)}\right)$$

$$a_{3}^{(2)} = s\left(\begin{array}{c} \theta_{30}^{(1)} x_{0} + \theta_{31}^{(1)} x_{1} + \theta_{32}^{(2)} x_{2} + \theta_{33}^{(2)} x_{3} \right) \Rightarrow a_{3}^{(2)} = s\left(z_{3}^{(2)}\right)$$

$$a_{3}^{(2)} = s\left(\begin{array}{c} \theta_{30}^{(1)} x_{0} + \theta_{31}^{(1)} x_{1} + \theta_{32}^{(2)} x_{2} + \theta_{33}^{(2)} x_{3} \right) \Rightarrow a_{3}^{(2)} = s\left(z_{3}^{(2)}\right)$$

$$a_{3}^{(2)} = s\left(\begin{array}{c} \theta_{30}^{(1)} x_{0} + \theta_{31}^{(1)} x_{1} + \theta_{32}^{(2)} x_{2} + \theta_{33}^{(2)} x_{3} \right) \Rightarrow a_{3}^{(2)} = s\left(z_{3}^{(2)}\right)$$

$$a_{3}^{(2)} = s\left(\begin{array}{c} \theta_{30}^{(1)} x_{0} + \theta_{31}^{(1)} x_{1} + \theta_{32}^{(2)} x_{2} + \theta_{33}^{(2)} x_{3} \right) \Rightarrow a_{3}^{(2)} = s\left(z_{3}^{(2)}\right)$$

$$a_{3}^{(2)} = s\left(\begin{array}{c} \theta_{30}^{(1)} x_{0} + \theta_{31}^{(1)} x_{1} + \theta_{32}^{(2)} x_{2} + \theta_{33}^{(2)} x_{3} \right)$$

$$a_{3}^{(2)} = s\left(\begin{array}{c} \theta_{30}^{(1)} x_{0} + \theta_{31}^{(1)} x_{1} + \theta_{32}^{(2)} x_{2} + \theta_{33}^{(2)} x_{3} \right) = a_{3}^{(2)} = s\left(z_{3}^{(2)}\right)$$

$$a_{3}^{(2)} = s\left(\begin{array}{c} \theta_{30}^{(1)} x_{0} + \theta_{31}^{(1)} x_{1} + \theta_{32}^{(2)} x_{2} + \theta_{33}^{(2)} x_{3} \right) = a_{3}^{(2)} = s\left(z_{3}^{(2)}\right)$$

$$a_{3}^{(2)} = s\left(\begin{array}(a_{3}^{(2)} x_{0} + \theta_{31}^{(2)} x_{1} + \theta_{32}^{(2)} x_{2} + \theta_{33}^{(2)} x_{3} \right) = a_{3}^{(2)} = s\left(z_{3}^{(2)} + \theta_{3}^{(2)} x_{3}^{(2)} + \theta_{3}^{(2)} x_{3}^{(2)} + \theta_{3}^{($$



Outline

1. Neural networks: basic concepts

2. Neural networks: learning

3. Tips & tricks

4. Deep learning



Multi-class classification

- Neural networks can deal straightforwardly with multiple outputs
- Suppose we want to classify in **4 classes**



• Then, we want that the following to hold:





Neural networks for classification

 $\mathcal{D} = \left\{ \left(\boldsymbol{x}(1), \boldsymbol{y}(1) \right), \dots, \left(\boldsymbol{x}(N), \boldsymbol{y}(N) \right) \right\}$

L = total number of layers in network

 w_l = number of units (not counting bias unit) in layer l



Binary classification

$$y = 0 \text{ or } y = 1$$

1 output unit

Multiclass classification (C classes)

$$y \in \mathbb{R}^{C} \qquad \begin{bmatrix} 1\\0\\0\\0\\0 \end{bmatrix} \qquad \begin{bmatrix} 0\\1\\0\\0\\0 \end{bmatrix} \qquad \begin{bmatrix} 0\\0\\1\\0\\0 \end{bmatrix} \qquad \begin{bmatrix} 0\\0\\0\\1\\0 \end{bmatrix} \qquad \begin{bmatrix} 0\\0\\0\\1\\0 \end{bmatrix}$$

C output units



Cost function for classification

Logistic regression

$$J(\theta) = -\sum_{i=1}^{N} (y(i) \ln h_{\theta}(x(i)) + (1 - y(i)) \ln[1 - h_{\theta}(x(i))]) + \frac{\lambda}{2} \sum_{j=1}^{d-1} \theta_{j}^{2}$$

Neural network

$$h_{\Theta}(\mathbf{x}) \in \mathbb{R}^{C} \qquad \{h_{\Theta}(\mathbf{x})\}_{c} = c \text{-th output}$$

$$J(\Theta) = -\sum_{i=1}^{N} \sum_{c=1}^{C} \left(\{y\}_{c}(i) \ln\{h_{\Theta}(\mathbf{x}(i))\}_{c} + (1 - \{y\}_{c}(i)) \ln\left[1 - \{h_{\Theta}(\mathbf{x}(i))\}_{c}\right]\right) + \frac{\lambda}{2} \sum_{l=1}^{L-1} \sum_{i=1}^{W_{l}} \sum_{j=1}^{W_{l}+1} \left(\Theta_{ji}^{(l)}\right)^{2}$$



Additional

regularization term

Gradient computation

Given the cost function $J(\Theta)$, we now have to minimize it. Need to compute:

- $J(\Theta)$: cost function
- $\frac{\partial}{\partial \Theta_{ij}^{(l)}} J(\Theta)$: gradient with respect to the network parameters (weights), $\Theta_{ij}^{(l)} \in \mathbb{R}$

Neural networks models use an algorithm called **Backpropagation**

Backpropagation is an algorithm to compute gradients. It relies on the **chain rule** of derivatives



Outline

1. Neural networks: basic concepts

2. Neural networks: learning

3. Tips & tricks

4. Deep learning



Activation functions

There are different (nonlinear) activation functions that can be used. The following can be used in the **hidden layers**. In the **output layer**, we used **sigmoid** for classification and a **linear** function for regression. In this case, the cost function can be the squared error.



Weights random initialization

It is not correct to inizialize all the parameters to zero. It turns out that in this way the network learns the same feature for each neuron, since $\frac{\partial J(\Theta)}{\partial \Theta_{ij}^{(l)}}$ is the same for each weight

 A proper initialization is mandatory to alleviate the vanishing or exploding gradient probems [13]

Common heuristics: generate weights from a Gaussian distribution of proper dimension, and then multiply for a number that depends on the network dimension and activation

<u>ReLU</u>

$$\Theta^{(l)} = \operatorname{randn}(w_l, w_{l-1}) \cdot \sqrt{\frac{2}{w_{l-1}}}$$

ReLU units exibit less vanishing or exploding gradients problems





Input normalization

It is often useful to **normalize** the inputs to the neural net (and in general to all optimization algorithms) in order to lead to faster convergence of the gradient descent





Input normalization



Normalized







Dipartimento di Ingegneria Gestionale, dell'Informazione e della Produzione

Mini batch gradient descent

Nowdays, neural networks are not trained by computing the full gradient, that requires all examples

- The stochastic gradient descent does a gradient descent step each time a new data is available
 - ✓ This is both faster to compute and also leads to improved regularization properties. However, the gradient based on only one sample tends to be very noisy
- For this reason, modern implementations allow to specify a **batch** of data that will be used to compute the gradient. This is called **mini-batch gradient descent**



Dropout

Dropout is a **regularization technique** specifically designed for neural networks [12]

- It consists into (probabilistically) turn off (drop) random neurons of the network during the training phase
- These units are not considered during forward and backward passes
- It is like the are performing an averaging of less complex models







(b) After applying dropout.

Dropout forces a neural network to learn more robust features that are useful in conjunction with many different random subsets of the other neurons.



Other tricks and heuristics

There exists different tricks that can be used to improve the learning of the neural net:

• Batch normalization: normalization step performed on every batch of data

• Learning rate decay: decrease the learning rate progressively as we reach towards the minimum

Advanced optimization schemes: Momentum, RMSprop, ADAM





https://towardsdatascience.com/types-of-optimization-algorithms-used-in-neural-networks-and-ways-to-optimize-gradient-95ae5d39529f

Stochastic Gradient Descent (SGD) is the slowest

SGD gets stuck in saddle point



A Dipartimento
 Di di Ingegneria Gestionale,
 O dell'Informazione e della Produzione

Outline

1. Neural networks: basic concepts

2. Neural networks: learning

3. Tips & tricks

4. Deep learning



Deep learning

The rise of available **data**, computational **power** and research **innovations** (such as the improvements in optimization just reviewed) allowed to train **more complex models**

- Neural networks with many layers and particular architectures were named deep learning models
- They became state-of-art in many tasks, most often for the more unstructed ones (computer vision, Natural Language Processing,...)
- In the next lessons we will talk about Convolutional Neural Networks, a particular type of deep architecture used for image analysis



Deep learning and bias-variance tradeoff

In principle, deep learning models require a **lot of data**, since they are very complex (remember the VC analysis)

- Therefore, most deep learning projects have the train and test data that come from different distributions (and different applicative domain)
- On the first domain we have a lot of data and we train the model. But, the trained model needs to be **employed on the second** domain, for which we could have few examples
- The train/validation/test set paradigm needs to slightly change to reflect this data mismatch problem [14]



Traditional train/validation/test and bias-variance

Say you want to build a human-level speech recognition system. You split your data into train/validation/test: Human error is often an "upper bound" for





Basic recipe for machine learning





Say you want to build a **speech recognition** system for a new **in-car** rearview mirror product. You have 50.000 hours of general speech data, and 10 hours of in-car data. How do you split your data?

Bad way



Having **mismatched val and test** distributions is not a good idea. Your team may spend months optimizing for dev set performance only to find it doesn't work well on the test set



Say you want to build a **speech recognition** system for a new **in-car** rearview mirror product. You have 50.000 hours of general speech data, and 10 hours of in-car data. How do you split your data?

Better way



Hence, a smarter way of splitting the above dataset would be just like the second line of

the diagram. In this way, the validation and test set are from the **same domain**



Say you want to build a **speech recognition** system for a new **in-car** rearview mirror product. You have 50.000 hours of general speech data, and 10 hours of in-car data. How do you split your data?

<u>Best way</u>



Create validation sets from **both data distributions**: a train-val and test-val set. In this

way, any gap between the different errors can help you tackle the problem more clearly







New recipe for machine learning





More general formulation

| | General speech data (50,000 hours) | In-car speech data (10 hours) | In this way we can compare performance in different settings, being able to better assess the model |
|--|--|---|---|
| Performance of humans | Human-level error | (Carry out human evaluation to measure.) | "Avoidable bias" |
| Performance on examples you've trained on | Training error | (Insert some in-car data into training set to measure.) | |
| Performance on examples you haven't trained on | Training-Val error | Val/Test error | "Variance"/degree of overfitting |
| Data mismatch | | | |



References

- 1. Provost, Foster, and Tom Fawcett. "*Data Science for Business: What you need to know about data mining and dataanalytic thinking*". O'Reilly Media, Inc., 2013.
- 2. Brynjolfsson, E., Hitt, L. M., and Kim, H. H. *"Strength in numbers: How does data driven decision making affect firm performance?"* Tech. rep., available at SSRN: <u>http://ssrn.com/abstract=1819486</u>, 2011.
- 3. Pyle, D. "*Data Preparation for Data Mining*". Morgan Kaufmann, 1999.
- 4. Kohavi, R., and Longbotham, R. "Online experiments: Lessons learned". Computer, 40 (9), 103–105, 2007.
- 5. Abu-Mostafa, Yaser S., Malik Magdon-Ismail, and Hsuan-Tien Lin. "Learning from data". AMLBook, 2012.
- 6. Andrew Ng. "Machine learning". Coursera MOOC. (https://www.coursera.org/learn/machine-learning)
- 7. Domingos, Pedro. "*The Master Algorithm*". Penguin Books, 2016.
- 8. Christopher M. Bishop, "Pattern recognition and machine learning", Springer-Verlag New York, 2006.
- 9. Hastie, T., Tibshirani, R.,, Friedman, J. *"The Elements of Statistical Learning"*. New York, NY, USA: Springer New York Inc, 2001.
- 10. Tom Fawcett, *"An introduction to ROC analysis",* Pattern Recognition Letters, Volume 27, Issue 8, 2006, Pages 861–874, ISSN 0167-8655,
- 11. Gareth James, Daniela Witten, Trevor Hastie, and Robert Tibshirani. *"An Introduction to Statistical Learning: With Applications in R".* Springer Publishing Company, Incorporated, 2014.
- 12. Nitish Srivastava, Geoffrey Hinton, Alex Krizhevsky, Ilya Sutskever, and Ruslan Salakhutdinov. "*Dropout: a simple way to prevent neural networks from overfitting*". J. Mach. Learn. Res. 15, 1 (January 2014), 1929–1958.
- 13. Andrew Ng. "Deep learning specialization". Coursera MOOC (<u>https://www.coursera.org/specializations/deep-learning</u>)
- 14. Andrew Ng, "The nuts and bolts of deep learning", <u>https://media.nips.cc/Conferences/2016/Slides/6203-Slides.pdf</u>

