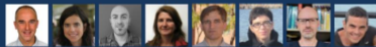


INTRODUCTION TO DEEP LEARNING

Winter School at UPC TelecomBCN Barcelona. 22-30 January 2018.



Instructors



Xavier
Giró-i-Nieto

Maria R.
Costa Jussà

Neel
Causa

Elise
Sayrol

Antonio
Bonafonte

Verónica
Vilaplana

Ramon
Marro

Javier
Ruiz

Organizers



Supporters

aws educate

GitHub Education

+ info: <https://telecombcn-dl.github.io/2018-idl/>

[\[course site\]](#)



#DLUPC

Day 1 Lecture 3 Perceptron



Antonio Bonafonte
antonio.bonafonte@upc.edu

Associate Professor
Universitat Politècnica de Catalunya



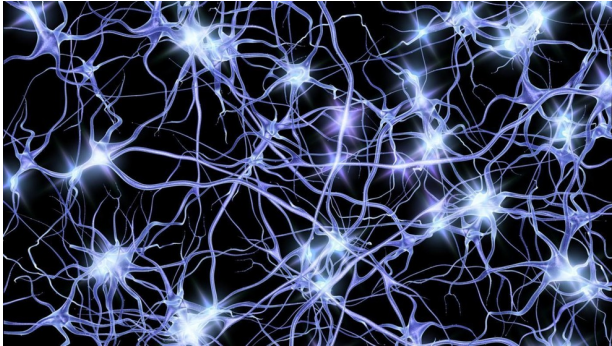
1. Artificial Neural Networks (ANN) and Perceptron
2. Linear Regression and Linear Classification
3. Activation Functions
4. Estimation of the Perceptron parameters (θ)

ANN: Artificial Neural Networks

Proposed in 1965.

Renaissance \approx 2010 (Deep Learning)

Inspired on biological Neural Networks



Biological Neural Networks

Human Brain

86 billion neurons

In average: 7 000 synapses
(connexions)

Purkinje cells:

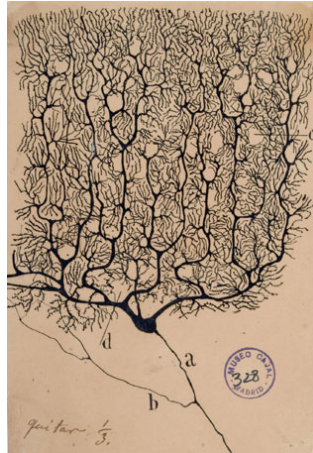
Input: 1000 dendritic branches,
tens of thousands of connexions

Output: 100 connexions

Total synapses: $\approx 10^{12}$

300 million synapses per mm^3

Plasticity (learning)



Ramon y Cajal Drawings

(\mathbf{X}, \mathbf{Y}) training examples

$\mathbf{X} = \mathbf{x}_1, \dots, \mathbf{x}_N$ input features

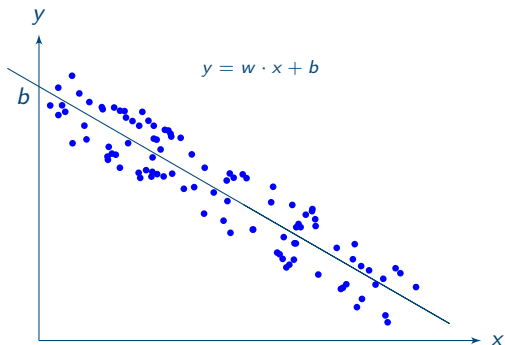
\mathbf{x}_n vector with components $\mathbf{x}_n(j) \in \mathbb{R}$ or $\mathbf{x}_n(j) \in \{0, 1\}$

$\mathbf{Y} = y_1, \dots, y_N$ labels (regression: \mathbb{R} ; class. $\{0, 1\}$)

$\hat{y} = f_{\theta}(\mathbf{x})$ prediction by model $f_{\theta}(\cdot)$

Input features and output: binary (1,0), *one-hot* encoding or *normalized* real-values.

Basic component: linear regression



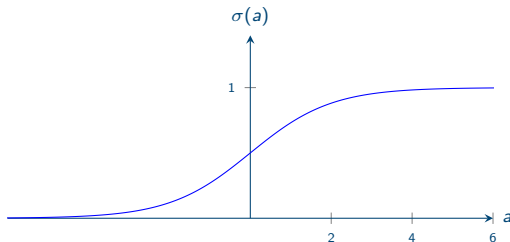
$$y = \mathbf{w}^T \cdot \mathbf{x} + b \equiv \bar{\mathbf{w}}^T \cdot \bar{\mathbf{x}}$$

(In the plot, x is a scalar value)

Example: x : mean speed; y : fuel consumption km/l

Basic component: logistic regression

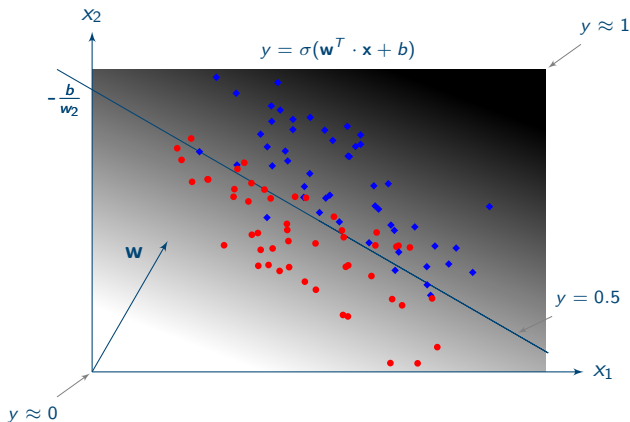
$$y = \sigma(\mathbf{w}^T \cdot \mathbf{x} + b) = \frac{1}{1 + e^{-(\mathbf{w}^T \cdot \mathbf{x} + b)}}$$



$$a = \mathbf{w}^T \cdot \mathbf{x} + b = \|\mathbf{w}\| \cdot \|\mathbf{x}\| \cdot \cos(\alpha) + b$$

Small $\|\mathbf{w}\| \rightarrow$ smoother transition boundaries.

Basic component: logistic regression



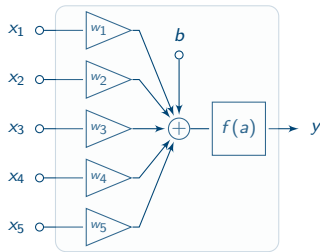
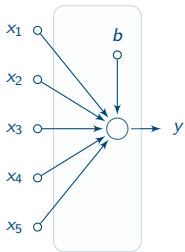
$$\theta = \{\mathbf{w}, b\} \quad \mathbf{w} = (w_1, w_2)^T$$

Example: x : (mean_spectrum, mean_f0); y : IS_MUSIC

Perceptron (Neuron)

Activation: $a = \mathbf{w}^T \cdot \mathbf{x} + b$

Output: $y = f(a)$



$f(\cdot)$: activation function

Linear regression: $f(a) = a$

Logistic regression: $f(a) = \sigma(a)$

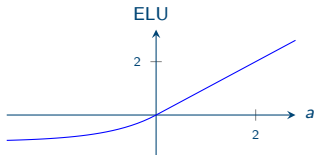
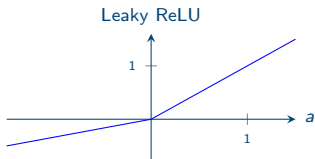
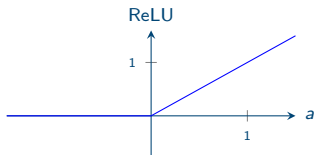
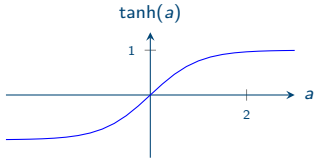
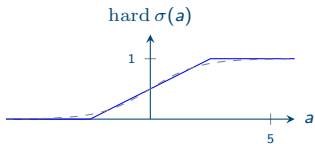
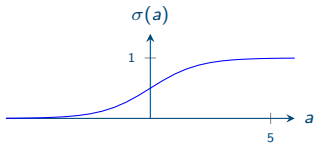
Activation functions

Activation function (non-linearities) are fundamental to combine neurons to model complex functions.

Some activation functions:

- Sigmoid $\sigma(\cdot)$, $\tanh(\cdot) = 2 \cdot \sigma(\cdot) - 1$
- ReLU (rectified linear unit), Leaky ReLU, Parametric Leaky ReLU, ELU (exponential linear unit)
- softmax, maxout

Some Activation functions



Training: estimate the parameters θ of model $f_\theta()$ from training data.

(\mathbf{X}, \mathbf{Y}) training examples

$\mathbf{X} = \mathbf{x}_1, \dots, \mathbf{x}_N$ input features

\mathbf{x}_n vector with components $\mathbf{x}_n(j) \in \mathbb{R}$ or $\mathbf{x}_n(j) \in \{0, 1\}$

$\mathbf{Y} = y_1, \dots, y_N$ labels (regression: \mathbb{R} ; class. $\{0, 1\}$)

$\hat{y} = f_\theta(\mathbf{x})$ prediction by model $f_\theta(\cdot)$

$\mathcal{L}(\mathbf{Y}, \hat{\mathbf{Y}})$ loss function

$$\theta^* = \underset{\theta}{\operatorname{argmin}} \mathcal{L}(\mathbf{Y}, \hat{\mathbf{Y}}) \equiv \underset{\theta}{\operatorname{argmin}} \frac{1}{N} \sum_{n=1}^N \ell(y_n, f_\theta(\mathbf{x}_n))$$

The loss function:

- Related with final metric (e.g., classification rate, number of clicks, revenue, ...)
- Nice mathematical properties (to find $\min_{\theta}(\cdot)$).

Example (Regression loss function)

Squared error:

$$\ell(y, \hat{y}) = (y - \hat{y})^2$$

Example (Classification loss function)

Binary cross entropy:

$$\ell(y, \hat{y}) = -(y \cdot \ln \hat{y} + (1 - y) \cdot \ln(1 - \hat{y}))$$

Solving $\operatorname{argmin}_{\theta}(\cdot)$

First idea:

1. Take derivatives of the loss with respect each parameter and equalize with zero.
2. Solve set of $|\theta|$ equations with $|\theta|$ variables.

This works for *losses* (e.g. squared error), easy *activation function* (identity) and architectures (e.g. one perceptron).

In general, millions of non-linear equations.

Finding the minimum in θ is a key aspect on deep learning.¹

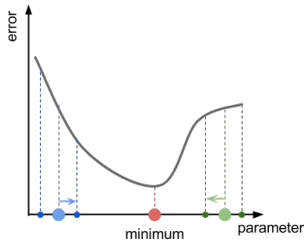
Most methods are based on the iterative-method *gradient descent*

1. Initialize: $\theta^{(0)}$, $i = 0$
2. Repeat till convergence from $i = 0, 1, \dots$:
 - Compute gradient of the loss function with respect θ at $\theta^{(i)}$
 - $\theta^{(i+1)} = \theta^{(i)} - \alpha \nabla \mathcal{L}(Y, f_{\theta^{(i)}}(X))$

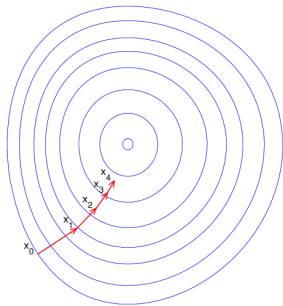
If the *learning rate* α is small enough,

$$\mathcal{L}(Y, f_{\theta^{(i+1)}}(X)) \leq \mathcal{L}(Y, f_{\theta^{(i)}}(X))$$

¹An introduction in this course, *D2L3 - Optimization*



$$\theta^{(i+1)} = \theta^{(i)} - \alpha \nabla \mathcal{L}(Y, f_{\theta^{(i)}}(X))$$
$$\mathcal{L}(Y, f_{\theta^{(i+1)}}(X)) \leq \mathcal{L}(Y, f_{\theta^{(i)}}(X))$$



From *Wikipedia*

$$\theta^{(i+1)} = \theta^{(i)} - \alpha \nabla \mathcal{L}(Y, f_{\theta^{(i)}}(X))$$

$$\mathcal{L}(Y, f_{\theta^{(i+1)}}(X)) \leq \mathcal{L}(Y, f_{\theta^{(i)}}(X))$$

Gradient computation

Let's assume the model $\hat{y}_n = f_{\theta}(\mathbf{x}_n)$ is just a perceptron:

$$\hat{y}_n = f(z_n) = f(\mathbf{w} \cdot \mathbf{x}_n) = f\left(\sum_{j=1}^J w_j \cdot x_n(j)\right)$$

Derivative with respect to a particular parameter: w_k :

$$\begin{aligned}\frac{\partial \mathcal{L}(\mathbf{Y}, \hat{\mathbf{Y}})}{\partial w_k} &= \frac{1}{N} \sum_{n=1}^N \frac{\partial \ell(y_n, \hat{y}_n)}{\partial w_k} \\ &= \frac{1}{N} \sum_{n=1}^N \frac{\partial \ell(y_n, \hat{y}_n)}{\partial \hat{y}_n} \cdot \frac{\partial \hat{y}_n}{\partial w_k} \\ &= \frac{1}{N} \sum_{n=1}^N \frac{\partial \ell(y_n, \hat{y}_n)}{\partial \hat{y}_n} \cdot f'(z_n) \cdot \frac{\partial z_n}{\partial w_k} \\ &= \frac{1}{N} \sum_{n=1}^N \frac{\partial \ell(y_n, \hat{y}_n)}{\partial \hat{y}_n} \cdot f'(z_n) \cdot x_k\end{aligned}$$

(See exercises with particular cases at the end of these slides)

- The perceptron is the basic unit in ANN
- It can be used both for linear regression and binary classification
- The weights can be estimated with iterative algorithms for any differentiable activation function and loss.
- Can be extended to neural networks (D2L1) and better optimization algorithms (D2L2)

Exercises

Exercise 1: regression

Lets assume:

Loss function, squared error $\ell(y, \hat{y}) = (y - \hat{y})^2$

Activation function, identity $f(z) = z$

Exercise

1. Show that the derivative of the \mathcal{L} with respect the weight vector is:

$$\frac{\partial \mathcal{L}(\mathbf{Y}, \hat{\mathbf{Y}})}{\partial \mathbf{w}_k} = \frac{2}{N} \sum_{n=1}^N (\hat{y} - y_n) \cdot \mathbf{x}_n(k)$$

2. Using previous equation for $k = 1, \dots, J$, and equalizing with 0, you can get the set of J linear equations and J variables²:

$$\mathbf{A} \cdot \mathbf{w} = \mathbf{b}$$

Find the values a_{ij} and b_j for $1 \leq i, j \leq J$.

²This part of the exercise only shows that in this particular case you can find a closed solution. It is not relevant to ANN

Exercise 2: classification – BCE

Lets assume:

Loss function, binary cross entropy $\ell(y, \hat{y}) = -(y \cdot \ln \hat{y} + (1 - y) \cdot \ln(1 - \hat{y}))$

Activation function, sigmoid $f(z) = \sigma(z)$

Exercise

1. Show that the derivative of the sigmoid

$$\sigma'(z) = \sigma(z) \cdot (1 - \sigma(z))$$

2. Show that the derivative of the \mathcal{L} with respect the weight vector is:

$$\frac{\partial \mathcal{L}(\mathbf{Y}, \hat{\mathbf{Y}})}{\partial w_k} = \frac{1}{N} \sum_{n=1}^N (\hat{y}_n - y_n) \cdot x_n(k)$$

Exercise 3: classification – MSE

Lets assume:

Loss function, squared error $\ell(y, \hat{y}) = (y - \hat{y})^2$

Activation function, sigmoid $f(z) = \sigma(z)$

Exercise

1. Show that the derivative of the \mathcal{L} with respect the weight vector is:

$$\frac{\partial \mathcal{L}(\mathbf{Y}, \hat{\mathbf{Y}})}{\partial w_k} = \frac{2}{N} \sum_{n=1}^N (\hat{y}_n - y_n) \cdot \sigma'(z_n) \cdot x_n(k)$$

2. Comparing this gradient and the one using binary cross entropy, the magnitude of both increases with difference between labels and predictions, $\hat{y}_n - y_n$. Do you think this is a good or bad characteristic of the algorithm?
3. However, the MSE loss introduces the term $\sigma'(z_n)$. Looking to the terms with big errors (e.g.: $y_n = 1$ and $\hat{y}_n \approx 0$), do you think this term helps in the convergence?