

# Introduction to Machine Learning

## Supervised Learning

Luís Nunes   Ricardo Ribeiro   Sancho Oliveira

Iscte – Instituto Universitário de Lisboa

2022/2023



# Learning Types

Reinforcement Learning	Learn the best (sequence of) action(s) given state(s)
Search	Given a set of possible solutions, find the best (a reasonable) one
Unsupervised Learning	Group similar things
Supervised Learning	Predict outcome. Learn a rule/model given examples

# Outline

## 1 Supervised Learning Problems

Regression

Classification

## 2 Linear Approximations

## 3 Perceptron and Delta Rule

## 4 Learning

Simple Networks

Backpropagation

## 5 Classical Algorithms

K-Nearest Neighbours

Naïve Bayes Classifier

Decision Trees

Emsembles

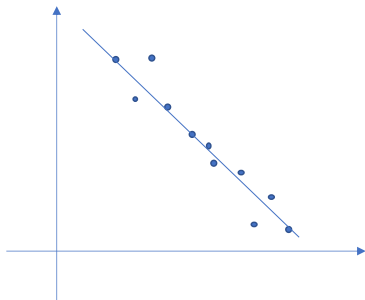
# Regression problem

- Given a set of examples
  - Each example has
    - A set of values, one for each attribute
    - A desired output: a **continuous** value

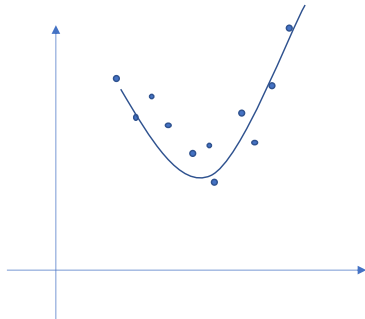
## Example

Attributes	Square meters	140
	Number of rooms	4
	...	...
Output	Current house price	€150,000

# Regression



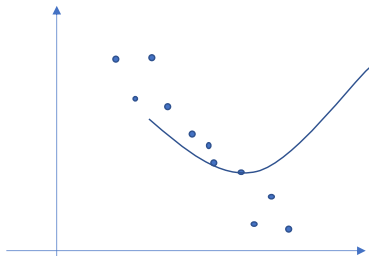
Linear



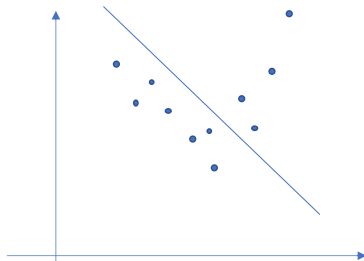
Non-linear

# Regression

when the model does not fit data



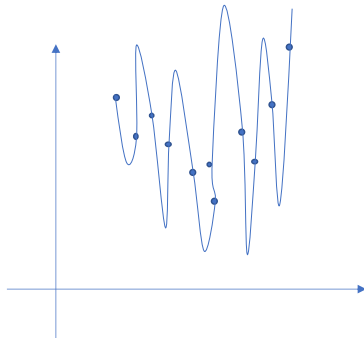
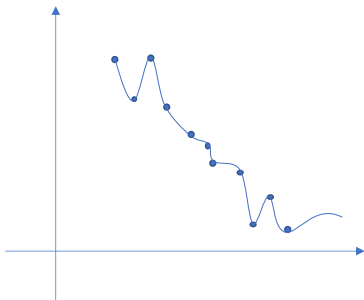
Model too complex



Model too simple

# Regression

when the model does not generalize well (overfit)



# Classification

problem

- Given a set of examples
  - Each example has
    - A set of values, one for each attribute
    - A desired output: a **category**

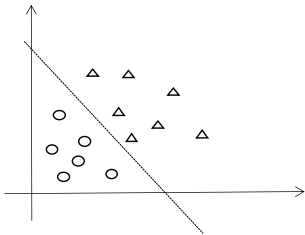
## Example

Attributes	Age	40
	Current balance	5,000
	Had previous loans	No
	Loan value	10,000
Output	Will pay current loan?	No



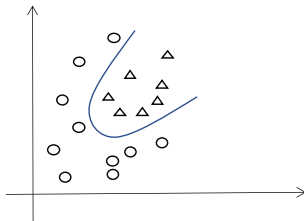
# Classification

separable classes



Linear separation

Non-linear separation



# Bias / Variance Dilemma

**Bias error** Model does/can not correctly represent the concept (**underfit**)

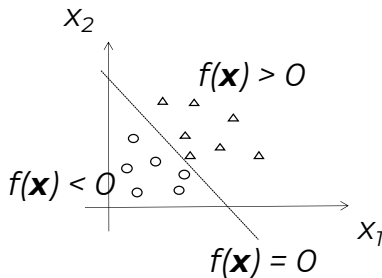
**Variance error** Model specializes in training set (**overfit**)

**Mitigating variance error** Regularization (favor smoother functions – output varies slowly with input)

# Classification

linearly separable classes

$$f(\mathbf{x}) = w_0 + w_1 \cdot x_1 + w_2 \cdot x_2$$



# Multilinear Regression

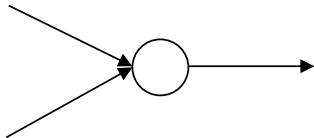
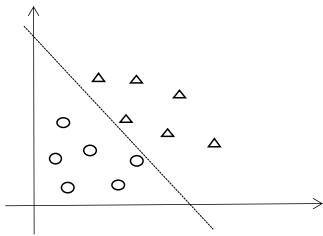
$$y = w_0 + w_1x_1 + w_2x_2 + \cdots + w_nx_n + \varepsilon$$

$$Y = WX + \varepsilon$$

- Assumptions
  - Relation between  $x_i$  and  $y$  is linear
  - All variables ( $x$ ) have Normal distributions
  - Variables are independent and residual / error ( $y(x_i) - \hat{y}(x_i)$ ) is constant
- Least Mean Squares

$$\hat{W} = (X^T X)^{-1} (X^T Y)$$

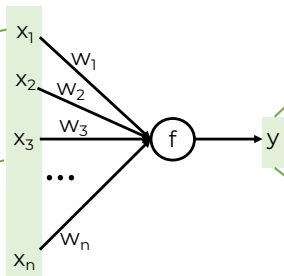
# Classification with Perceptron



# Artificial Neuron

An artificial neuron, in a feedforward neural network, receives a set of inputs and generates a single output.

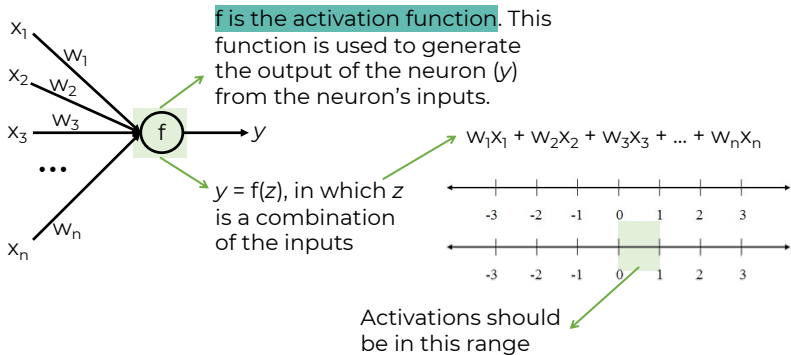
The input of an artificial neuron comes from all neurons of the previous layer or it is an external input



The output of an artificial neuron is a simple combination of the received inputs

The output of an artificial neuron is sent to all neurons of the next layer or is (part of) the network output.

# Artificial Neuron



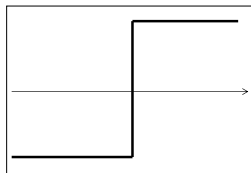
# Perceptron and Delta Rule

simple function and update rule for classification

$$f(x) = \begin{cases} 1, & w_0 + w_1x_1 + w_2x_2 + \dots > 0 \\ -1, & w_0 + w_1x_1 + w_2x_2 + \dots \leq 0 \end{cases}$$

$$\Delta w_i = \alpha(f(x) - d)x_i$$

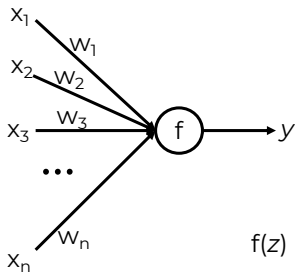
where  $d$  is the desired value



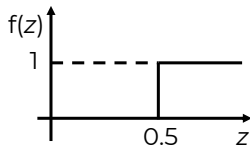


# Boolean Operations Network

example



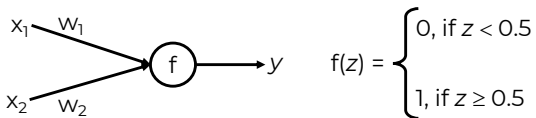
$$f(z) = \begin{cases} 0, & \text{if } z < 0.5 \\ 1, & \text{if } z \geq 0.5 \end{cases}$$



In which,  $z = \sum_{i=1}^n (w_i x_i)$

# Boolean Operations Network

AND

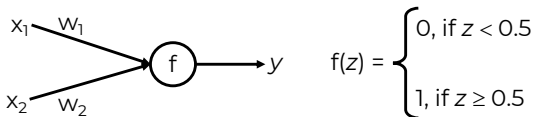


$x_1$	$x_2$	AND
0	0	0
0	1	0
1	0	0
1	1	1

$x_1$	$x_2$	$W_1 = W_2 = 0.1$	$W_1 = W_2 = 0.3$
		$Y = f(0.1 \times x_1 + 0.1 \times x_2)$	$Y = f(0.3 \times x_1 + 0.3 \times x_2)$
0	0	$f(0.1 \times 0 + 0.1 \times 0) = f(0) = 0$	$f(0.3 \times 0 + 0.3 \times 0) = f(0) = 0$
0	1	$f(0.1 \times 0 + 0.1 \times 1) = f(0.1) = 0$	$f(0.3 \times 0 + 0.3 \times 1) = f(0.3) = 0$
1	0	$f(0.1 \times 1 + 0.1 \times 0) = f(0.1) = 0$	$f(0.3 \times 1 + 0.3 \times 0) = f(0.3) = 0$
1	1	$f(0.1 \times 1 + 0.1 \times 1) = f(0.2) = 0$	$f(0.3 \times 1 + 0.3 \times 1) = f(0.6) = 1$

# Boolean Operations Network

OR



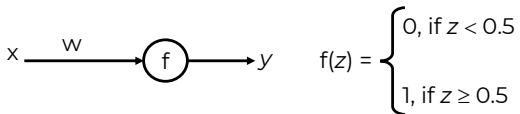
$$f(z) = \begin{cases} 0, & \text{if } z < 0.5 \\ 1, & \text{if } z \geq 0.5 \end{cases}$$

$x_1$	$x_2$	OR
0	0	0
0	1	1
1	0	1
1	1	1

		$w_1 = w_2 = 0.6$
$x_1$	$x_2$	$Y = f(0.6 \times x_1 + 0.6 \times x_2)$
0	0	$f(0.6 \times 0 + 0.6 \times 0) = f(0) = 0$
0	1	$f(0.6 \times 0 + 0.6 \times 1) = f(0.6) = 1$
1	1	$f(0.6 \times 1 + 0.6 \times 0) = f(0.6) = 1$
1	1	$f(0.6 \times 1 + 0.6 \times 1) = f(1.2) = 1$

# Boolean Operations Network

NOT



X	NOT
0	1
1	0

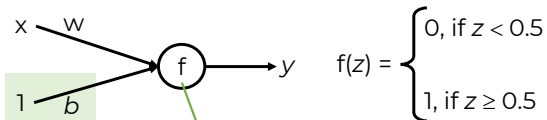
If the network correctly implemented the Boolean NOT, whenever  $x = 1, y = 0$ ; and whenever  $x = 0, y = 1$

But  $f(w \times 0) = f(0) = 0$ . Independently of the value of the weight  $w$ ,  $f(w \times 0) = 0$

Thus, this cannot be used to implement the Boolean NOT; something else is required: **bias value**

# Boolean Operations Network

NOT



X	NOT
0	1
1	0

$b$  is called the bias of the neuron

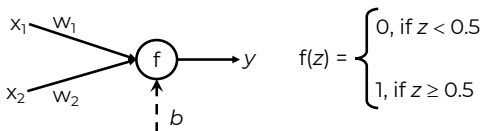
General expression of the output:

$$y = f(z), \text{ with } z = b + \sum_{i=1}^n (w_i x_i)$$

		$W = -0.5; b = 0.5$
X	1	$Y = f(-0.5 \times X + 0.5)$
0	1	$f(-0.5 \times 0 + 0.5) = f(0.5) = 1$
1	1	$f(-0.5 \times 1 + 0.5) = f(0) = 0$

# Boolean Operations Network

## NAND



$x_1$	$x_2$	NAND
0	0	1
0	1	1
1	0	1
1	1	0

With different network parameters (weights and biases), it is possible to implement different relations between input and output

$x_1$	$x_2$	$B=0.75; W_1 = W_2 = -0.25$
0	0	$Y = f(0.75 - 0.25 \times (x_1 + x_2))$
0	1	$f(0.75 - 0.25 \times 0) = f(0.75) = 1$
1	0	$f(0.75 - 0.25) = f(0.5) = 1$
1	1	$f(0.75 - 0.25 \times (1+1)) = f(0.25) = 0$

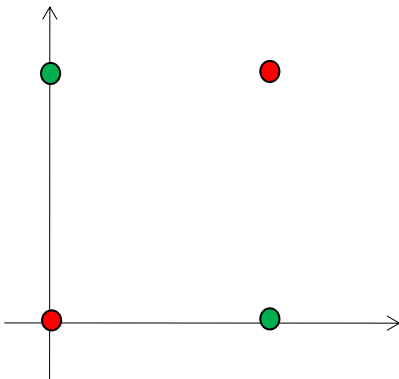
# Learning

changing the network parameters

- Changing the parameters of the network (weights and biases), it is possible to implement different relations between input and output
- Thus, if we want a network to learn the input/output relation implicit in a given training set, the learning algorithm just needs to change the network parameters (connection weights and neuron biases)
- But there are certain input/output relations that required more complex networks than a single neuron besides the input layer

## Non-linearly Separable Classes

Perceptron draws linear separation between classes, not suitable for all problems (XOR) [Minsky & Papert 69]



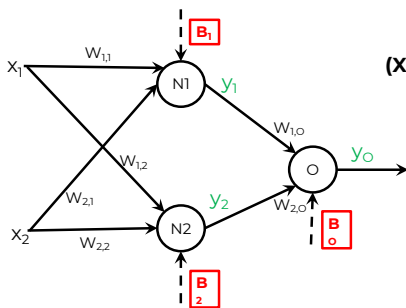


# Boolean Operations Network

XOR

$$\text{XOR}(x_1, x_2) = \text{AND}(\text{OR}(x_1, x_2), \text{NAND}(x_1, x_2))$$

$x_1$	$x_2$	XOR
0	0	0
0	1	1
1	0	1
1	1	0



$$(x_1 \vee x_2) \wedge \neg(x_1 \wedge x_2)$$

$N1$  implements the Boolean OR of the two inputs  $x_1$  and  $x_2$


$N2$  implements the Boolean NAND of the two inputs  $x_1$  and  $x_2$

$O$  implements the Boolean AND of the outputs of  $N1$  and  $N2$ ,  $y_1$  and  $y_2$

# Boolean Operations Network

XOR

$$\text{XOR}(x_1, x_2) = \text{AND}(\text{OR}(x_1, x_2), \text{NAND}(x_1, x_2))$$

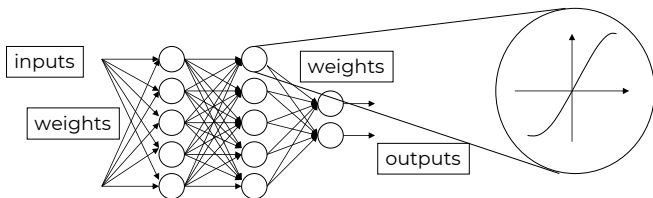

$$(X_1 \vee X_2) \wedge \neg(X_1 \wedge X_2)$$

		OR	NAND	AND	XOR
$X_1$	$X_2$	$Y_1 = f(0.6 \times X_1 + 0.6 \times X_1)$	$Y_2 = f(0.75 - 0.25 \times (X_1 + X_2))$	$Y_O = f(0.3 \times Y_1 + 0.3 \times Y_2)$	
0	0	$f(0) = 0$	$f(0.75) = 1$	$f(0.3) = 0$	0
0	1	$f(0.6) = 1$	$f(0.5) = 1$	$f(0.6) = 1$	1
1	0	$f(0.6) = 1$	$f(0.5) = 1$	$f(0.6) = 1$	1
1	1	$f(1.2) = 1$	$f(0.25) = 0$	$f(0.3) = 0$	0

# Artificial Neural Networks

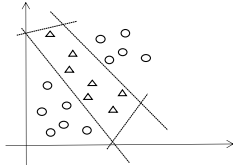
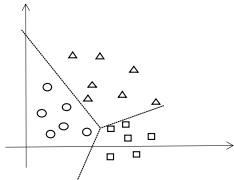
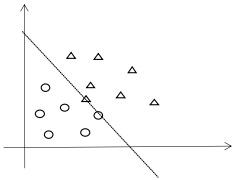
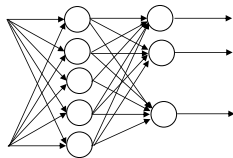
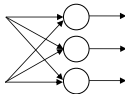
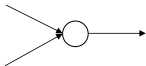
[Rumelhart, Hinton, Williams 86]

- Multilayer Perceptron
- Performance depends only on number of hidden layer units



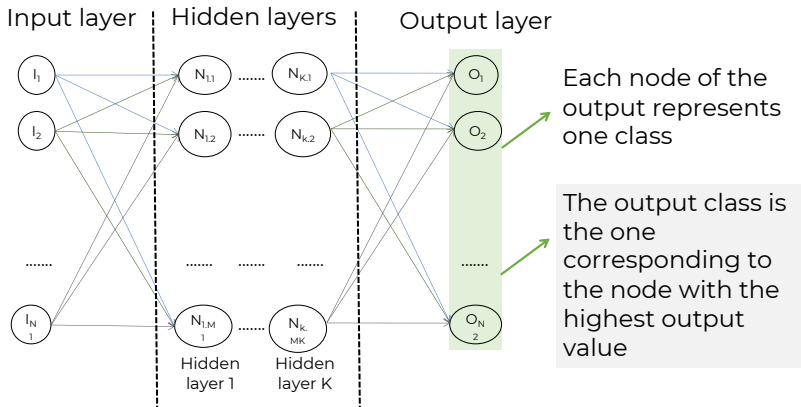
# Classification with MLP

potential for space division



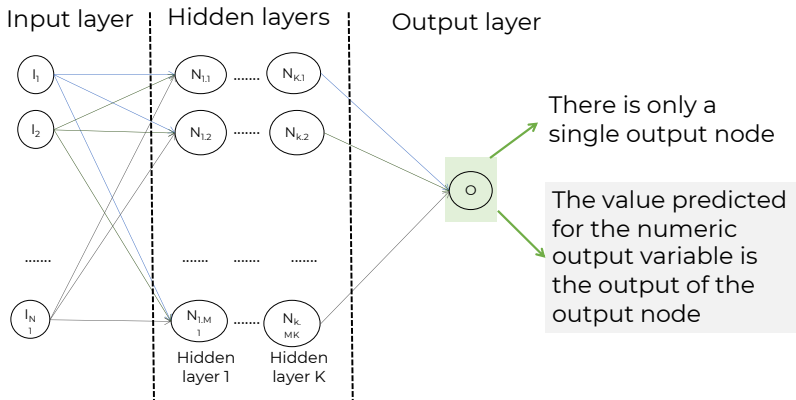
# Feedforward Neural Network with Backpropagation

Classification

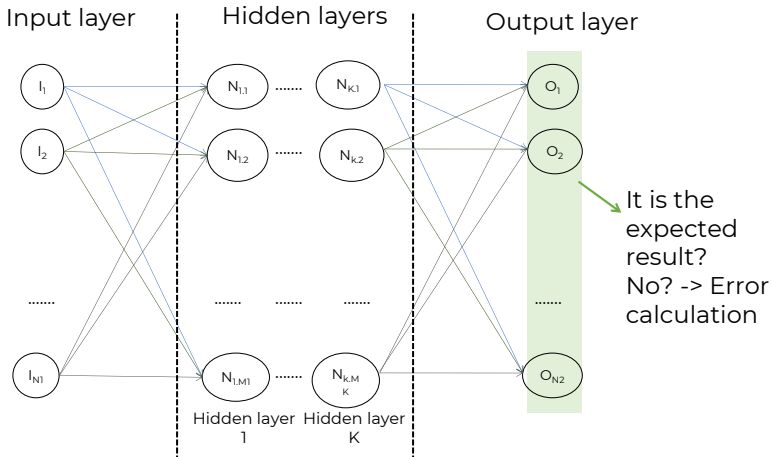


# Feedforward Neural Network with Backpropagation

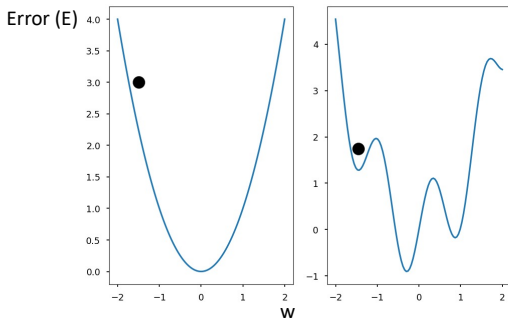
Regression



# Feedforward Neural Network with Backpropagation



# Gradient Descent Algorithm



$$\frac{\partial E}{\partial w}$$

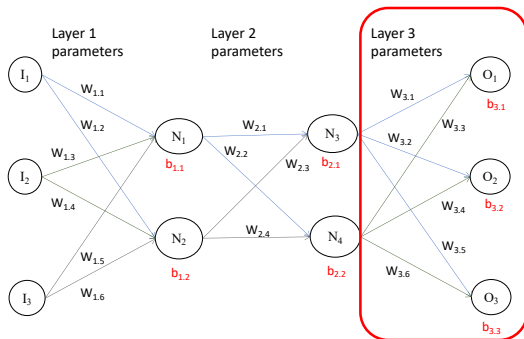


# Derivative of the error

$$\Delta_P = -(\eta \times \frac{\partial E}{\partial P}),$$

in which P is a parameter,  $\Delta_P$  is the update to be added to P, E is the error and  $\eta$  is a proportionality constant called the **learning rate**

# Backpropagation



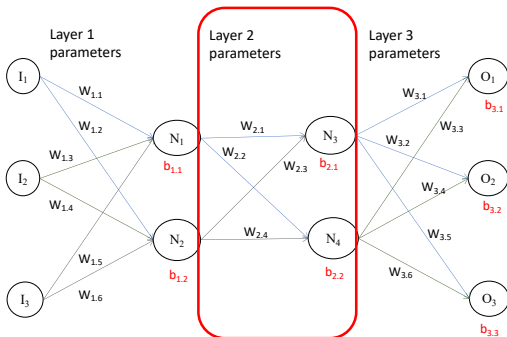
For instance,

$$w_{3.1} = w_{3.1} + \Delta_{w_{3.1}}$$

$$\Delta_{w_{3.1}} = -(\eta \times \frac{\partial E}{\partial w_{3.1}})$$

$\frac{\partial E}{\partial w_{3.1}}$  depends on the network outputs

# Backpropagation



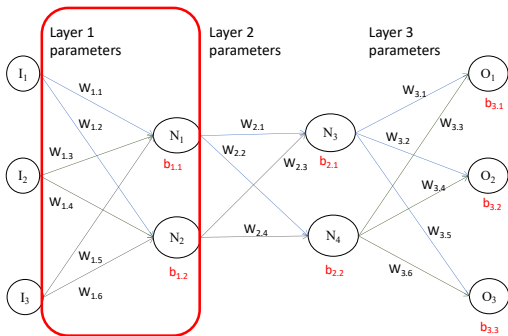
For instance,

$$W_{2.1} = W_{2.1} + \Delta_{w2.1}$$

$$\Delta_{w2.1} = -\left(\eta \times \frac{\partial E}{\partial w_{2.1}}\right)$$

$\frac{\partial E}{\partial w_{2.1}}$  depends only on the network outputs and on the derivatives of the error with respect to the parameters of layer 3 (already computed)

# Backpropagation



For instance,

$$b_{1.1} = b_{1.1} + \Delta_{b_{1.1}}$$

$$\Delta_{b_{1.1}} = -(\eta \times \frac{\partial E}{\partial b_{1.1}})$$

$\frac{\partial E}{\partial b_{1.1}}$  depends only on the network outputs and on the derivatives of the error with respect to the parameters of layer 2 (already computed)

# Backpropagation

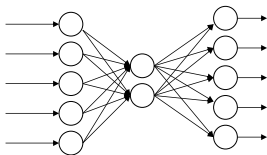
[Rumelhart, Hinton, Williams 86]

- Present each example  $(x(i), d(i))$
- Calculate network response  $x(i) : f(x(i))$
- Propagate error backwards (iteratively building error derivative at each layer)
- Save partial derivatives
- After all examples processed, update weights

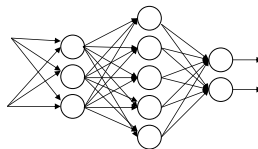
# Neural Network

## Compression vs Feature Generation

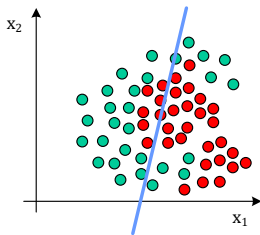
Compression



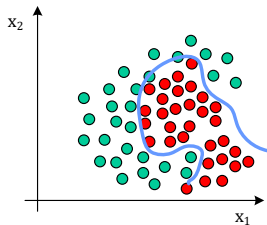
Feature generation



# Activation Functions



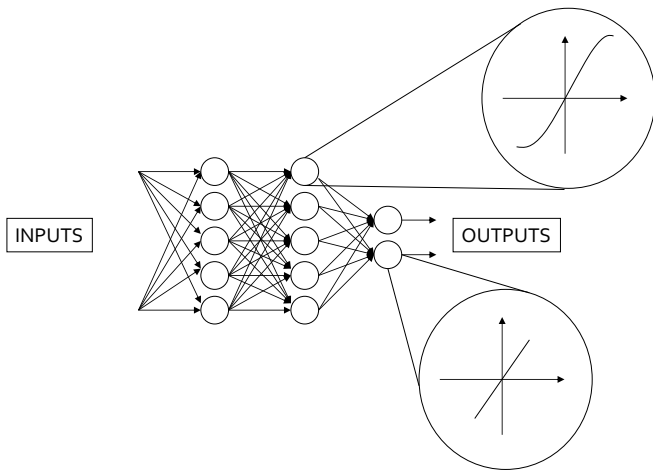
Linear activation function allows a linear separation



Non-linear activation function allow complex separations

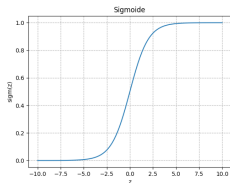
# MPL

Linear outputs for regression



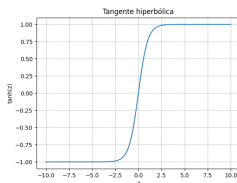


# Activation Functions



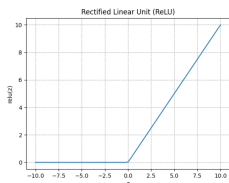
$$\text{sigm}(z) = \frac{1}{1 + e^{-z}}$$

$$\text{sigm}'(z) = \text{sigm}(z) \times \text{sigm}(1 - z)$$



$$\tanh(z) = \frac{e^z - e^{-z}}{e^z + e^{-z}}$$

$$\tanh'(z) = 1 - \tanh(z)^2$$



$$\text{ReLU}(z) = \begin{cases} x, & \text{se } x > 0 \\ 0, & \text{se } x \leq 0 \end{cases}$$

$$\text{ReLU}'(z) = \begin{cases} 1, & \text{se } x > 0 \\ 0, & \text{se } x \leq 0 \end{cases}$$

# ANN/MLP+backpropagation

- Analytic methods of classification / regression deal badly with noise and are sensitive to numerical approximations
- ANN are:
  - Robust to noise and approximations
  - Based in simplified neuron model
  - Incremental training
  - Compress information of many examples in small model

# Deep Learning

- Alternating prediction layers with feature decorrelation
- Techniques used were first thought of in the 70s and 80s
- ... now benefit from the massive computing power available

# Deep Learning

## Examples

- Object Detection with Tensorflow API:  
[https://www.youtube.com/watch?v=\\_zZe27JYi8Y](https://www.youtube.com/watch?v=_zZe27JYi8Y)
- Quadcopter Navigation in the Forest using Deep Neural Networks:  
<https://www.youtube.com/watch?v=umRdt3zGgpU>
- Real-time face recognition with Deep Learning technology:  
<https://www.youtube.com/watch?v=B4m2RVFLbME>

# References

- Almeida, L.B.D., 1990. and Multilayer Perceptrons F. Emile & B. Russell, eds. Analysis, 12(12), pp.1167-1178.
- Minsky, M. & Papert, S., 1969. Perceptrons: An Introduction to Computational Geometry, MIT Press. Available at: <http://www.computer.org/portal/web/csdl/doi/10.1109/T-C.1969.222718>.
- Minsky, M.L. & Papert, S.A., 1988. Perceptrons: Expanded edition, The MIT Press. Available at: <http://mitpress.mit.edu/book-home.tcl?isbn=0262631113>.
- Rumelhart, D E, Hinton, G.E. & Williams, R.J., 1986. Learning internal representation by error propagation. In D E Rumelhart & J. L. McClelland, eds. Parallel Distributed Processing Explorations in the Microstructure of Cognition. Cambridge, MA: MIT Press, pp. 318-362.
- Rumelhart, David E, Widrow, B. & Lehr, M.A., 1994. The basic ideas in neural networks. Communications of the ACM, 37(3), pp.87-92. Available at: <http://portal.acm.org/citation.cfm?doid=175247.175256>.
- Silva, F.M. & Almeida, L.B., 1990. Acceleration Techniques for the Backpropagation Algorithm. Neural Networks Proc EURASIP Workshop, 412, pp.110-119.

# $k$ -Nearest Neighbours

- Find  $k$  patterns in the set most similar to the one to classify
- Select a class between those of known patterns (how? Most common? Only consider majority? ties?)
- Problems:
  - Define distance,
  - Define class selection,
  - Non-linear problems

# $k$ -Nearest Neighbours

## Attribute Values

Sky	Temperature	Humidity	Wind	Sea	Prediction
Clear	Warm	Normal	Strong	Warm	Stable
Cloudy	Cold	High	Weak	Cold	Unstable
Rain					

# $k$ -Nearest Neighbours

Learn

- Can we find the pattern?

#	Sky	Temp.	Humid	Wind	Sea	Pred.	Go surf?
1	Clear	Warm	Normal	Strong	Warm	Stable	Yes
2	Clear	Warm	Normal	Strong	Warm	Unstable	No
3	Cloudy	Cold	High	Strong	Cold	Unstable	No
4	Clear	Warm	High	Strong	Cold	Stable	Yes
5	Rain	Cold	High	Strong	Warm	Stable	Yes
6	Rain	Cold	High	Weak	Warm	Unstable	No



# Naïve Bayes Classifier

- Calculate probabilities and use Bayes theorem:

Diagram illustrating the components of Bayes' theorem:

$$P(A | B) \equiv \frac{P(B | A)P(A)}{P(B)}$$

Labels and connections:

- posterior** points to  $P(A | B)$
- prior** points to  $P(A)$
- $P(B | A)/P(B)$ : support** points to the fraction  $\frac{P(B | A)P(A)}{P(B)}$

$$P(go = Yes | x = \dots) \equiv \frac{P(x | go = Yes)P(go = Yes)}{P(x)}$$

$x = \{\text{cloudy, cold, normal, strong, warm, unstable}\}$

# Naïve Bayes Classifier

$$P(\text{go} = \text{Yes} | x = \dots) \equiv \frac{P(x | \text{go} = \text{Yes})P(\text{go} = \text{Yes})}{P(x)}$$

$x = \{\text{rain, cold, high, strong, warm, stable}\}$

$P(\text{go} = \text{Yes}) = 3/6$  : #Yes / # observations

$P(x = \dots) = 1/6$  : #patterns equal to  $x$  / # observations

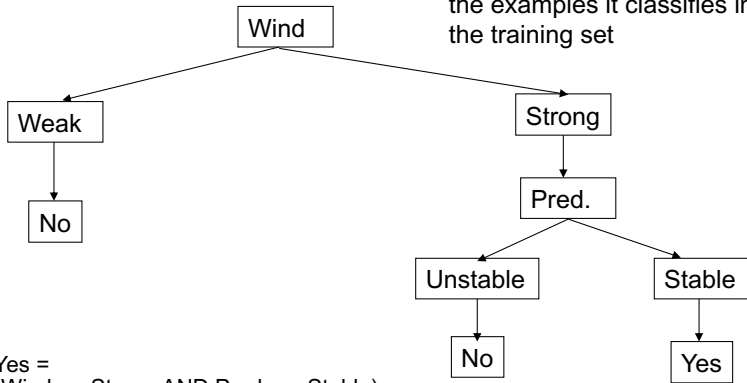
$P(x = \dots | \text{go} = \text{Yes}) = 1/3 * 1/3 * 2/3 * 3/3 * 2/3 * 3/3 = 4/81$  : probability of each attribute equal to its value in  $x$

$P(\text{go} = \text{Yes} | x = \dots) = 0,148\dots$

# Decision Tree

Equivalent

Each node is associated to the examples it classifies in the training set



Yes =  
(Wind == Strong AND Pred. == Stable)

# Decision Tree

## Homogeneity of a set

The **entropy** of a set is the measure of the diversity of its elements.

A set has the largest entropy if each of its elements belongs to a different class.

The smaller the entropy of a set, the larger its homogeneity!

$$\text{Entropy}(\text{Set}) = -\sum_{i=1, \dots, n} [P(x_i) \times \log_2(P(x_i))]$$

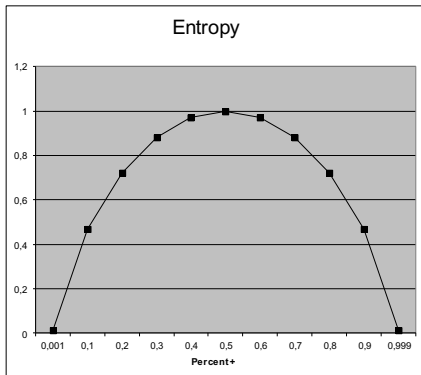
in which,  $P(x_i)$  is the probability of picking an element of class  $x_i$ .

Example Classes: go\_surf(**no** / **yes**)

$$\text{Entropy}(\text{Set}) = -[P(\text{no}) \times \log_2(P(\text{no})) + P(\text{yes}) \times \log_2(P(\text{yes}))]$$

# Decision Tree

## Homogeneity of a set



Maximum entropy  
when  $p_+ = p_-$ , i.e.  
lower probability  
to predict example  
class

# Decision Tree

## The best split

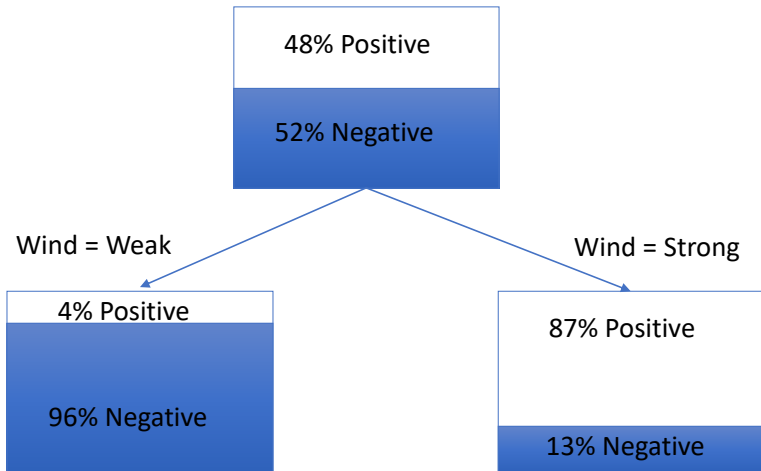
The **best split** is the split that results in the largest entropy reduction, that is, the largest **information gain (IG)**.

$$\text{IG}(\text{Set}/\{\text{Subset1}, \text{Subset2}\}) = \text{Entropy}(\text{Set}) - ((1/\#\text{Set}) \times (\#\text{Subset1} \times \text{Entropy}(\text{Subset1}) + \#\text{Subset2} \times \text{Entropy}(\text{Subset2})))$$

in which,  $S/\{S_1, S_2\}$  is the split of  $S$  into the two subsets  $S_1$  and  $S_2$ , and  $\#S$  is the cardinality of set  $S$ .

# Decision Tree

## Subsets and Gain



# Decision Tree

ID3(Examples, Target-Attribute, Attributes)

- Create root
  - If  $p_+ = 1$ : root = +
  - If  $p_- = 1$ : root = -
  - If Attributes =  $\emptyset$ , root = most common target-value in examples
- $A \leftarrow$  Attribute with best information gain
- Root = A
- For each (v) possible for A:
  - Add branch A = v
  - ExamplesV = Set of examples where A=v
    - If ExamplesV ==  $\emptyset$ : add branch with most common target-value in Examplesv
    - else branch = ID3(ExamplesV, A, Attributes - A)



# Decision Tree

C4.5 / C5.0 (Quinlan 96)

- Similar to ID3, but ...
  - Support for continuous attributes: discretizes continuous attributes
  - Allows missing values: examples not used when calculating entropy
  - Allows different costs for attributes
  - Pruning

# Learning ensembles

- Boosting (Kearns 88)
  - Can a set of weak learners create a single strong learner?
  - Classification combines the results of all the subtrees
  - Misclassified examples become more important for the error in each iteration
  - New trees are trained to fit the residual error
- Bagging - Bootstrap aggregating: (Breiman 96)
  - Selects randomly the subsets
  - Trains several learners,
  - Classification by voting, regression by averaging

# XGBoost (eXtreme Gradient Boosting)

(Chen Guestrin 2016)

- An optimized Gradient Boosting Machine
- Uses many small trees
- Classifies an example by joining the scores of each of the various trees
- Train by adding trees that improve the result or pruning
- Trees are fitted to predict the residual error

# Examples ML

- <https://www.youtube.com/watch?v=yeS8TJwBAFs>  
Not on today's subject ...
- <https://www.youtube.com/watch?v=wL7tSgUpy8w&t>

# References

- <https://xgboost.readthedocs.io/en/latest/tutorials/model.html>
- Chen, Tianqi; Guestrin, Carlos (2016). "XGBoost: A Scalable Tree Boosting System". In Krishnapuram, Balaji; Shah, Mohak; Smola, Alexander J.; Aggarwal, Charu C.; Shen, Dou; Rastogi, Rajeev (eds.). Proceedings of the 22nd ACM SIGKDD International Conference on Knowledge Discovery and Data Mining, San Francisco, CA, USA, August 13-17, 2016. ACM. pp. 785–794. arXiv:1603.02754. doi:10.1145/2939672.2939785.

# Summary

## 1 Supervised Learning Problems

Regression

Classification

## 2 Linear Approximations

## 3 Perceptron and Delta Rule

## 4 Learning

Simple Networks

Backpropagation

## 5 Classical Algorithms

K-Nearest Neighbours

Naïve Bayes Classifier

Decision Trees

Emsembles