# Classification

# Classification

**Source: Towards Data Science**

# Classification

CAT

(LABELED PHOTOS)
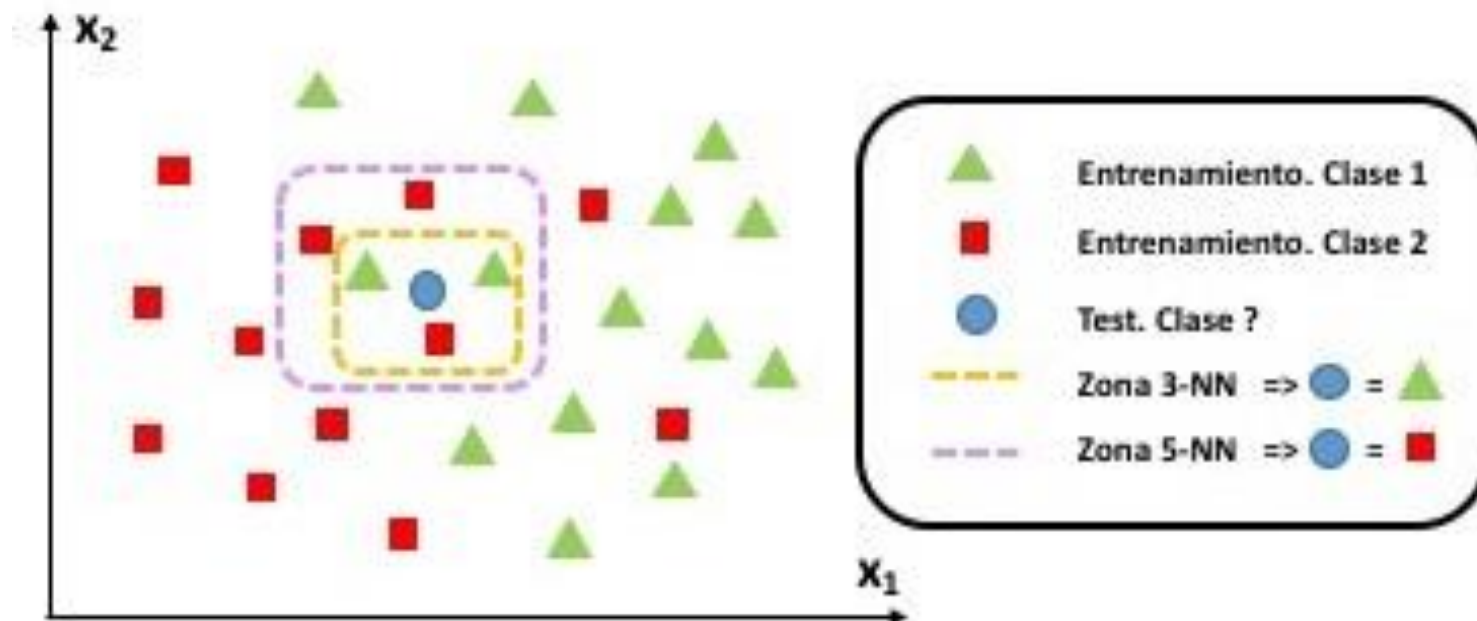
DOG

OUTPUT

**Input
Data**

**Mathematical
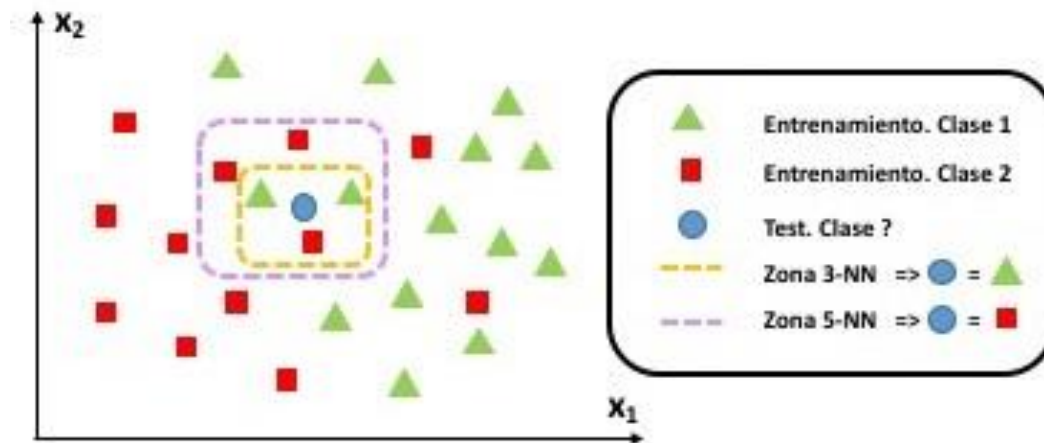Model**

**Class
Prediction**
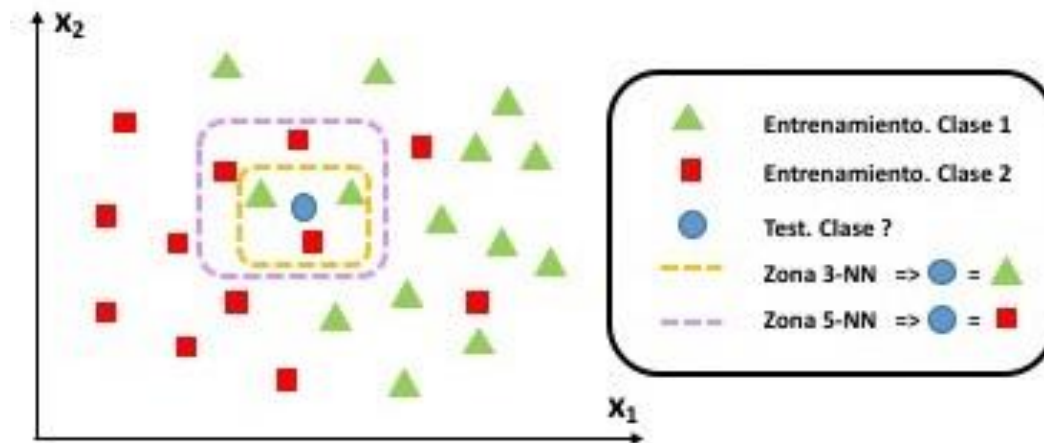
# K-Nearest Neighbors

# K-Nearest Neighbors

- Given a training set and a sample, for which we want to obtain a prediction

# K-Nearest Neighbors

- Given a training set and a sample, for which we want to obtain a prediction:
  - Calculate all pairwise "distances" between the test samples and the train samples



Legend:
- ▲ Entrenamiento. Clase 1
- ▇ Entrenamiento. Clase 2
- ● Test. Clase ?
- Zona 3-NN => ● = ▲
- Zona 5-NN => ● = ▇

# K-Nearest Neighbors

- Given a training set and a sample, for which we want to obtain a prediction:
    - Calculate all pairwise "distances" between the test samples and the train samples
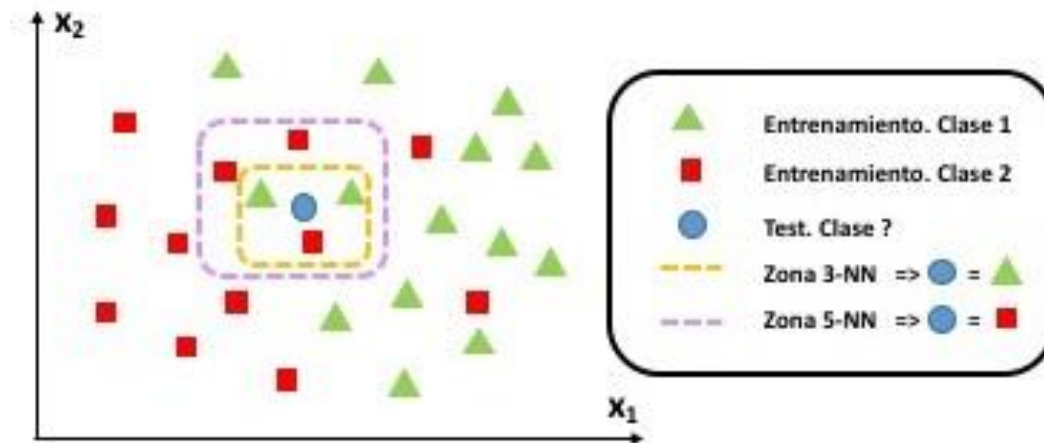    - Select the K smaller distances

# K-Nearest Neighbors

- Given a training set and a sample, for which we want to obtain a prediction:
    - Calculate all pairwise "distances" between the test samples and the train samples
    - Select the K smaller distances
    - Predict the class for the test sample by means of majority voting among its K nearest neighbors
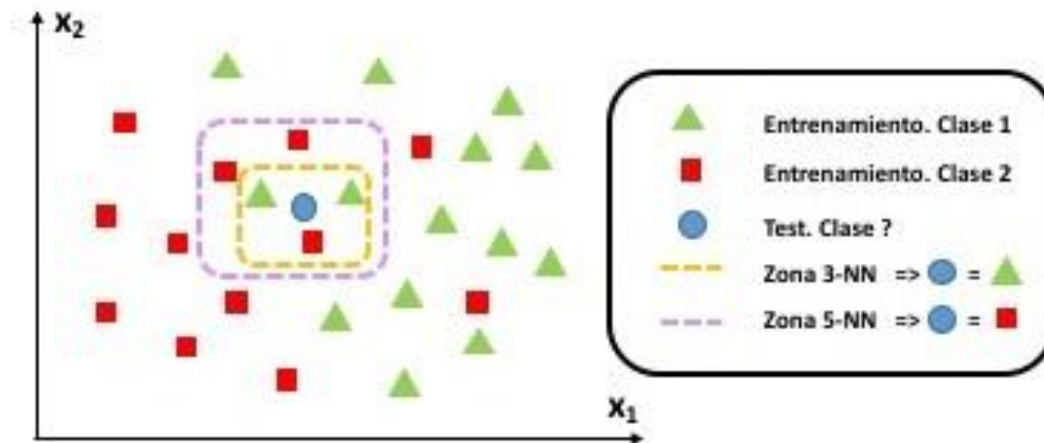
# K-Nearest Neighbors

- Given a training set and a sample, for which we want to obtain a prediction:
    - Calculate all pairwise "distances" between the test samples and the train samples
    - Select the K smaller distances
    - Predict the class for the test sample by means of majority voting among its K nearest neighbors
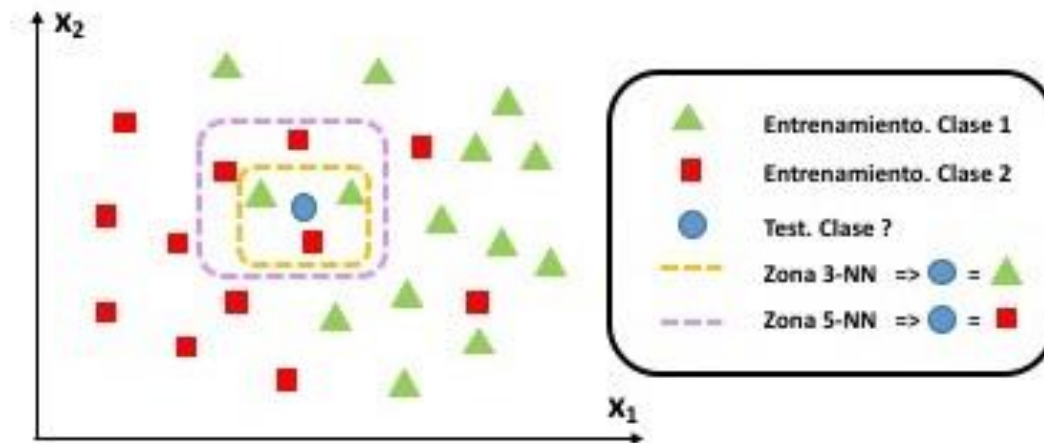


- Variants

# K-Nearest Neighbors

- Given a training set and a sample, for which we want to obtain a prediction:
    - Calculate all pairwise "distances" between the test samples and the train samples
    - Select the K smaller distances
    - Predict the class for the test sample by means of majority voting among its K nearest neighbors



- Variants:
    - Use a weighted mean, which weights are inversely proportional to the distance
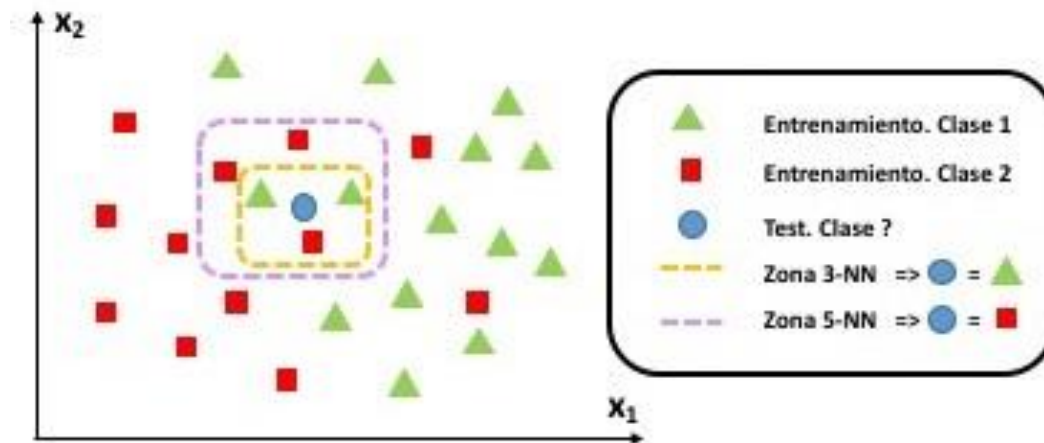
# K-Nearest Neighbors

- Given a training set and a sample, for which we want to obtain a prediction:
  - Calculate all pairwise "distances" between the test samples and the train samples
  - Select the K smaller distances
  - Predict the class for the test sample by means of majority voting among its K nearest neighbors



- Variants:
  - Use a weighted mean, which weights are inversely proportional to the distance
  - Without prefixing K, choose a threshold for defining the neighborhood (radius) and consider as neighbors all training samples that are not further than that radius
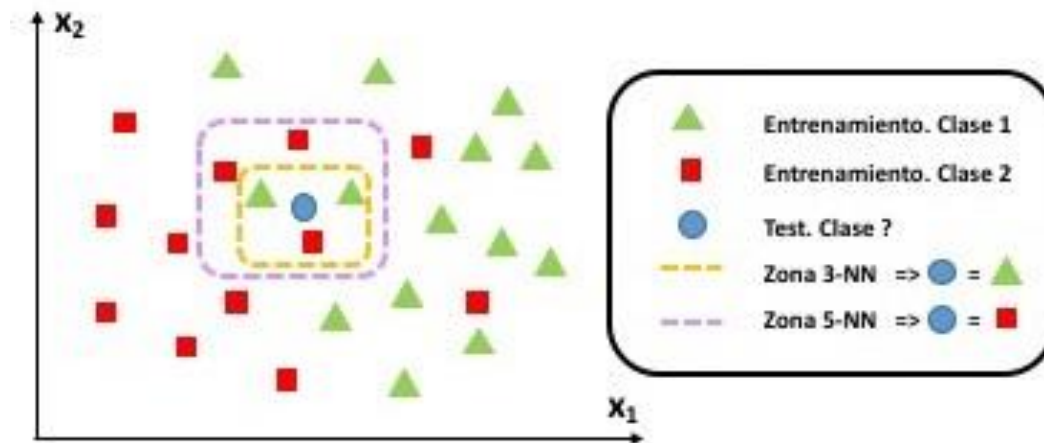
# K-Nearest Neighbors

- Given a training set and a sample, for which we want to obtain a prediction:
    - Calculate all pairwise "distances" between the test samples and the train samples
    - Select the K smaller distances
    - Predict the class for the test sample by means of majority voting among its K nearest neighbors



- Variants:
    - Use a weighted mean, which weights are inversely proportional to the distance
    - Without prefixing K, choose a threshold for defining the neighborhood (radius) and consider as neighbors all training samples that are not further than that radius
- **Prefixing K** is a tricky task to solve
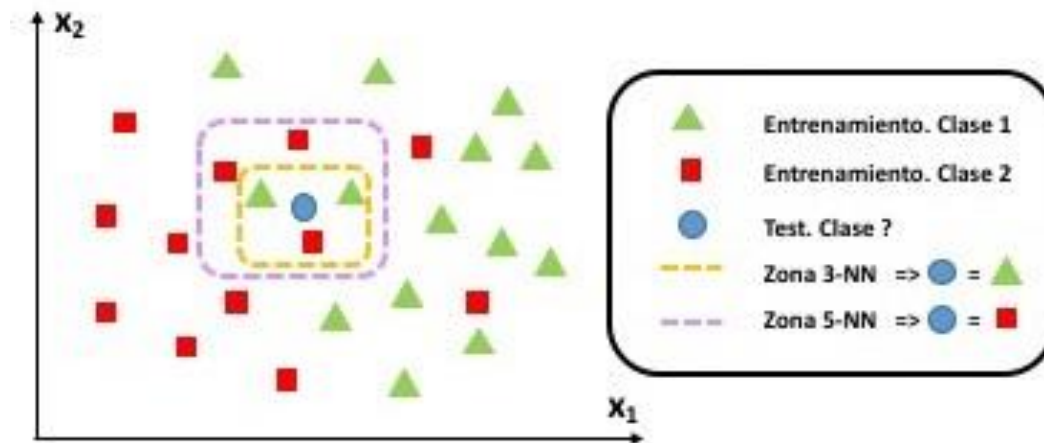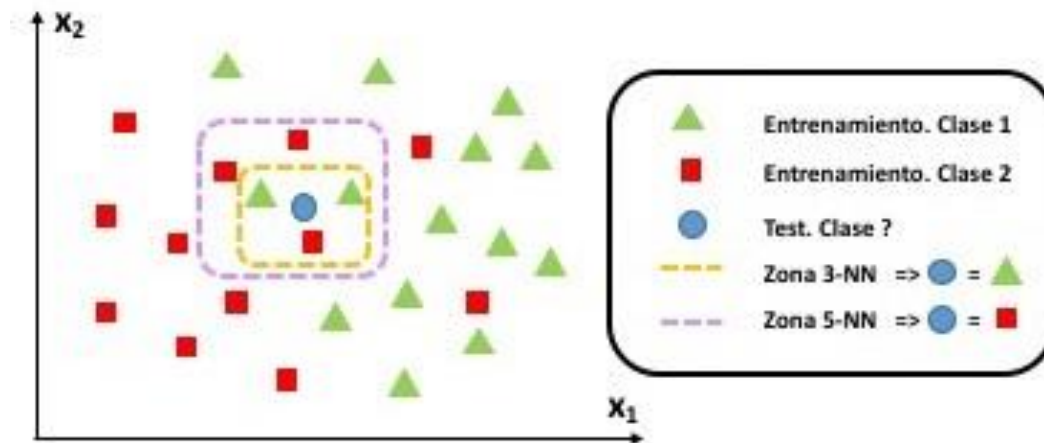
# K-Nearest Neighbors

- Given a training set and a sample, for which we want to obtain a prediction:
    - Calculate all pairwise "distances" between the test samples and the train samples
    - Select the K smaller distances
    - Predict the class for the test sample by means of majority voting among its K nearest neighbors
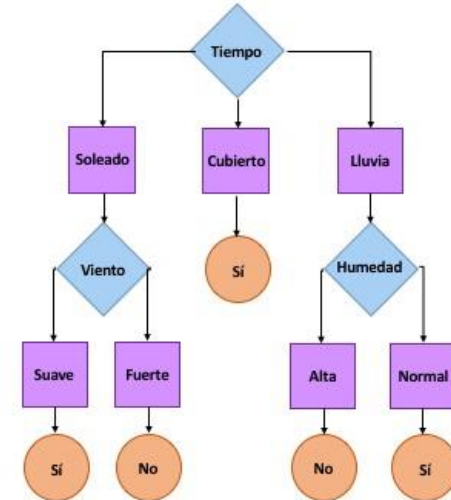


- Variants:
    - Use a weighted mean, which weights are inversely proportional to the distance
    - Without prefixing K, choose a threshold for defining the neighborhood (radius) and consider as neighbors all training samples that are not further than that radius
- **Prefixing K** is a tricky task to solve
- We must design an strategy for **breaking ties**

# Decision trees

- Decision Trees (C4.5)

  - They can handle both continuous and descrete variables

  - They are highly interpretable if the number of variables is not much high

  - They can suffer from *overfitting* (adapt too much to the training data loosing extrapolation power to unseen data)

| Tiempo | Temp. | Humedad | Viento | Jugar golf |
|---|---|---|---|---|
| Lluvia | Calor | Alta | Suave | No |
| Lluvia | Calor | Alta | Fuerte | No |
| Cubierto | Calor | Alta | Suave | Sí |
| Soleado | Templado | Alta | Suave | Sí |
| Soleado | Frío | Normal | Suave | Sí |
| Soleado | Frío | Normal | Fuerte | No |
| Cubierto | Frío | Normal | Fuerte | Sí |
| Lluvia | Templado | Alta | Suave | No |
| Lluvia | Frío | Normal | Suave | Sí |
| Soleado | Templado | Normal | Suave | Sí |
| Lluvia | Templado | Normal | Fuerte | Sí |
| Cubierto | Templado | Alta | Fuerte | Sí |
| Cubierto | Calor | Normal | Suave | Sí |
| Soleado | Templado | Alta | Fuerte | No |

# Decision trees

- Decision Trees (C4.5)
  - They can handle both continuous and descrete variables
  - They are highly interpretable if the number of variables is not much high
  - They can suffer from *overfitting* (adapt too much to the training data loosing extrapolation power to unseen data)

# Decision trees

- Decision Trees (C4.5)

  – They can handle both continuous and descrete variables

  – They are highly interpretable if the number of variables is not much high

  – They can suffer from *overfitting* (adapt too much to the training data loosing extrapolation power to unseen data)



depth = 1

# Decision trees

- Decision Trees (C4.5)
  - They can handle both continuous and descrete variables
  - They are highly interpretable if the number of variables is not much high
  - They can suffer from *overfitting* (adapt too much to the training data loosing extrapolation power to unseen data)
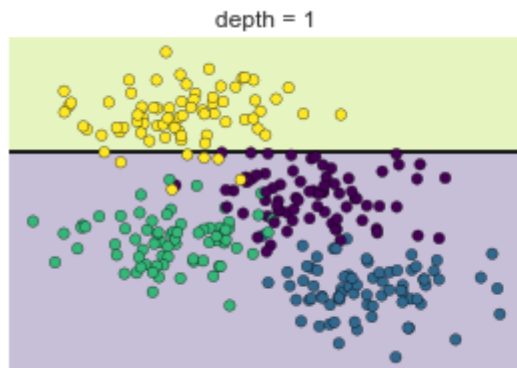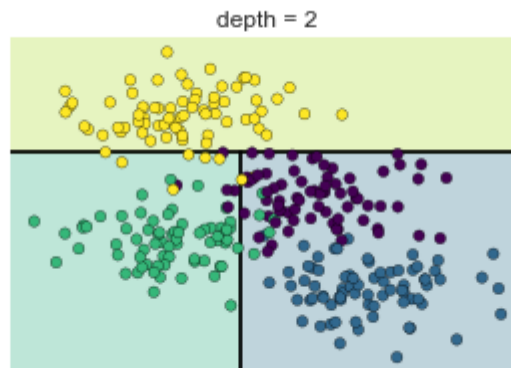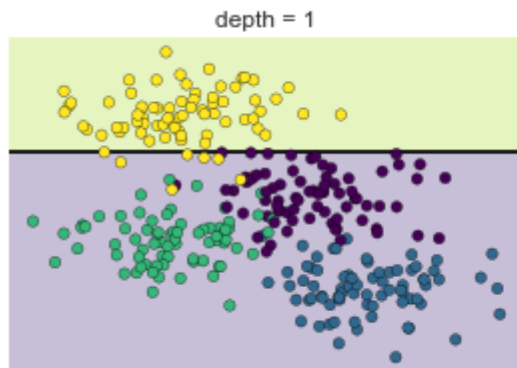
# Decision trees

- Decision Trees (C4.5)
  - They can handle both continuous and descrete variables
  - They are highly interpretable if the number of variables is not much high
  - They can suffer from *overfitting* (adapt too much to the training data loosing extrapolation power to unseen data)



depth = 1     depth = 2     depth = 3

# Logistic regression

- Logit function

$$p(X) = \frac{e^{(\beta_0 + \beta_1 X)}}{1 + e^{(\beta_0 + \beta_1 X)}}$$



Logistic Regression: 1 Feature

- samples (y=1)
- samples (y=0)
- logistic: g(z)
- linear: z

# Logistic regression

- Logit function

$$p(X) = \frac{e^{(\beta_0 + \beta_1 X)}}{1 + e^{(\beta_0 + \beta_1 X)}}$$

- Horizontal asymptotes
  - y = 0, towards the left
  - y = 1, towards the right

Logistic Regression: 1 Feature



- samples (y=1)
- samples (y=0)
- logistic: g(z)
- linear: z

# Logistic regression

- Logit function

$$p(X) = \frac{e^{(\beta_0 + \beta_1 X)}}{1 + e^{(\beta_0 + \beta_1 X)}}$$

- Horizontal asymptotes
  - y = 0, towards the left
  - y = 1, towards the right

- It estimates classes probabilities

Logistic Regression: 1 Feature

samples (y=1)
samples (y=0)
logistic: g(z)
linear: z

y

Feature

# Logistic regression

- Logit function

$$p(X) = \frac{e^{(\beta_0 + \beta_1 X)}}{1 + e^{(\beta_0 + \beta_1 X)}}$$

- Horizontal asymptotes
  - y = 0, towards the left
  - y = 1, towards the right

- It estimates classes probabilities

- It is not recommendable for multiclass classification, i.e. not binary



Logistic Regression: 1 Feature

Legend:
- samples (y=1)
- samples (y=0)
- logistic: g(z)
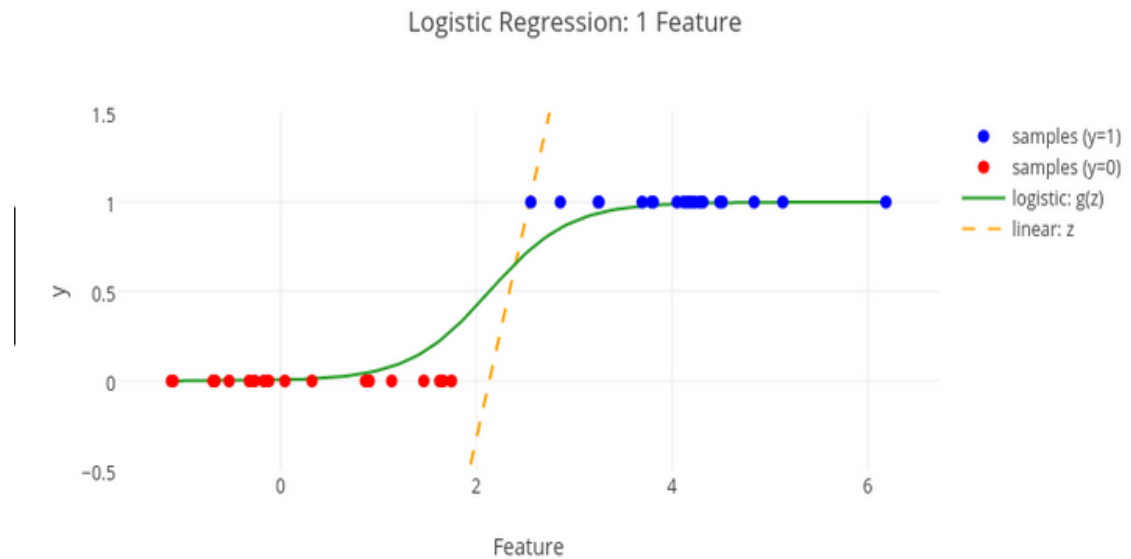- linear: z

# Naïve Bayes

**Prior**

$P(C+) = \alpha$
$P(C-) = 1-\alpha$

**Train inputs**

| $X_1$ | $X_2$ | | $X_{n-1}$ | $X_n$ |
|---|---|---|---|---|
| 3 | 6 | ... | 5 | 9 |
| 5 | 1 | ... | 5 | 6 |
| 4 | 6 | ... | 5 | 5 |
| 7 | 89 | ... | 23 | 85 |
| 3 | 435 | ... | 3 | 1 |
| ... | ... | ... | ... | ... |
| ... | ... | ... | ... | ... |
| 8 | 1 | ... | 77 | 321 |
| 9 | 8 | ... | 6 | 8 |
| 4 | 77 | ... | 3 | 132 |
| 8 | 9 | ... | 1 | 8 |

**Conditional independence**

**Bayes Theorem**

**Prediction**

# Naïve Bayes

We need to calculate

$$p(C|X_1, \ldots, X_n)$$



**Prior**

P(C+) = α
P(C−) = 1−α

**Train inputs**

| X₁ | X₂ | | X_{n-1} | X_n |
|---|---|---|---|---|
| 3 | 6 | ... | 5 | 9 |
| 5 | 1 | ... | 5 | 6 |
| 4 | 6 | ... | 5 | 5 |
| 7 | 89 | ... | 23 | 85 |
| 3 | 435 | ... | 3 | 1 |
| ... | ... | ... | ... | ... |
| ... | ... | ... | ... | ... |
| 8 | 1 | ... | 77 | 321 |
| 9 | 8 | ... | 6 | 8 |
| 4 | 77 | ... | 3 | 132 |
| 8 | 9 | ... | 1 | 8 |

**Conditional independence**

Bayes Theorem

**Prediction**

# Naïve Bayes

We need to calculate

$$p(C|X_1, \ldots, X_n)$$

Bayes Theorem,

$$p(C|X_1, \ldots, X_n) = \frac{p(C) \cdot p(X_1, \ldots, X_n|C)}{p(X_1, \ldots, X_n)}$$

**Prior**

| P(C+) = α |
| P(C-) = 1-α |

**Train inputs**

| X₁ | X₂ | | Xₙ₋₁ | Xₙ |
|----|----|----|----|----|
| 3 | 6 | ... | 5 | 9 |
| 5 | 1 | ... | 5 | 6 |
| 4 | 6 | ... | 5 | 5 |
| 7 | 89 | ... | 23 | 85 |
| 3 | 435 | ... | 3 | 1 |
| ... | ... | ... | ... | ... |
| ... | ... | ... | ... | ... |
| 8 | 1 | ... | 77 | 321 |
| 9 | 8 | ... | 6 | 8 |
| 4 | 77 | ... | 3 | 132 |
| 8 | 9 | ... | 1 | 8 |

**Conditional independence**

Bayes Theorem → **Prediction**

# Naïve Bayes

We need to calculate

$$p(C|X_1, \ldots, X_n)$$

$p(C, X_1, \ldots, X_n)$

Joint probability

Bayes Theorem,

$$p(C|X_1, \ldots, X_n) = \frac{p(C) \cdot p(X_1, \ldots, X_n|C)}{p(X_1, \ldots, X_n)}$$

**Prior**

P(C+) = α
P(C-) = 1-α

**Train inputs**

| X₁ | X₂ | | Xₙ₋₁ | Xₙ |
|---|---|---|---|---|
| 3 | 6 | ... | 5 | 9 |
| 5 | 1 | ... | 5 | 6 |
| 4 | 6 | ... | 5 | 5 |
| 7 | 89 | ... | 23 | 85 |
| 3 | 435 | ... | 3 | 1 |
| ... | ... | ... | ... | ... |
| ... | ... | ... | ... | ... |
| 8 | 1 | ... | 77 | 321 |
| 9 | 8 | ... | 6 | 8 |
| 4 | 77 | ... | 3 | 132 |
| 8 | 9 | ... | 1 | 8 |

**Conditional independence**

**Bayes Theorem** → **Prediction**

# Naïve Bayes

We need to calculate

$$p(C|X_1,\ldots,X_n)$$

**Joint probability** → $p(C,X_1,\ldots,X_n)$

Bayes Theorem,

$$p(C|X_1,\ldots,X_n) = \frac{p(C) \cdot p(X_1,\ldots,X_n|C)}{p(X_1,\ldots,X_n)}$$

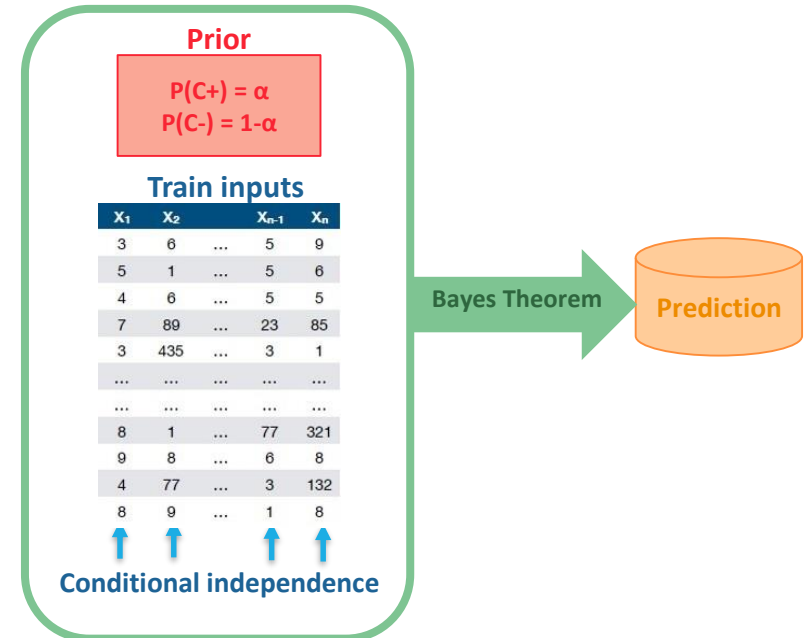Applying the chain rule in conditional probability

$$p(C,X_1,\ldots,X_n) = p(C) \cdot p(X_1|C) \cdot p(X_2|C,X_1) \cdot \ldots \cdot p(X_n|C,X_1,\ldots,X_{n-1})$$

**Prior**

P(C+) = α
P(C-) = 1-α

**Train inputs**

| X₁ | X₂ | | Xₙ₋₁ | Xₙ |
|---|---|---|---|---|
| 3 | 6 | ... | 5 | 9 |
| 5 | 1 | ... | 5 | 6 |
| 4 | 6 | ... | 5 | 5 |
| 7 | 89 | ... | 23 | 85 |
| 3 | 435 | ... | 3 | 1 |
| ... | ... | ... | ... | ... |
| ... | ... | ... | ... | ... |
| 8 | 1 | ... | 77 | 321 |
| 9 | 8 | ... | 6 | 8 |
| 4 | 77 | ... | 3 | 132 |
| 8 | 9 | ... | 1 | 8 |

**Conditional independence**

**Bayes Theorem** → **Prediction**

# Naïve Bayes

We need to calculate

$$p(C|X_1, \ldots, X_n)$$

Joint probability → $p(C, X_1, \ldots, X_n)$

Bayes Theorem,

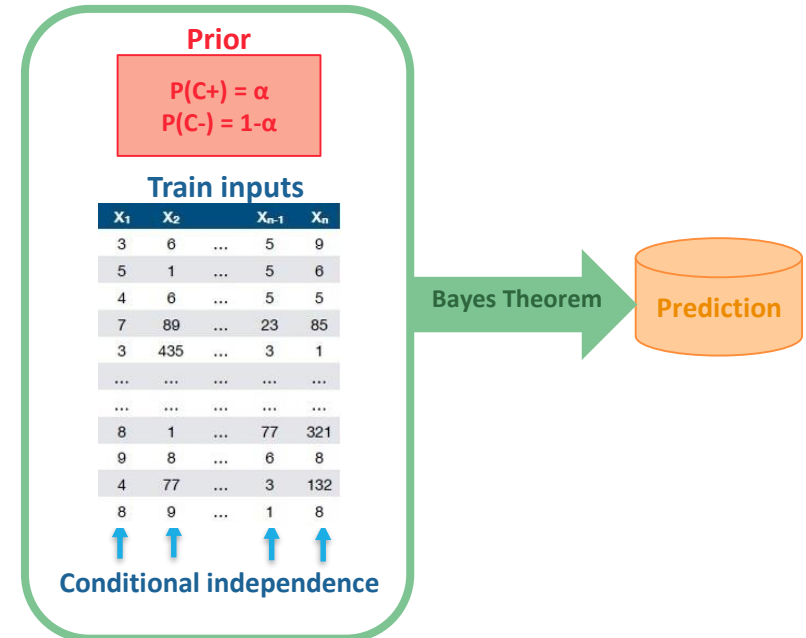$$p(C|X_1, \ldots, X_n) = \frac{p(C) \cdot p(X_1, \ldots, X_n|C)}{p(X_1, \ldots, X_n)}$$

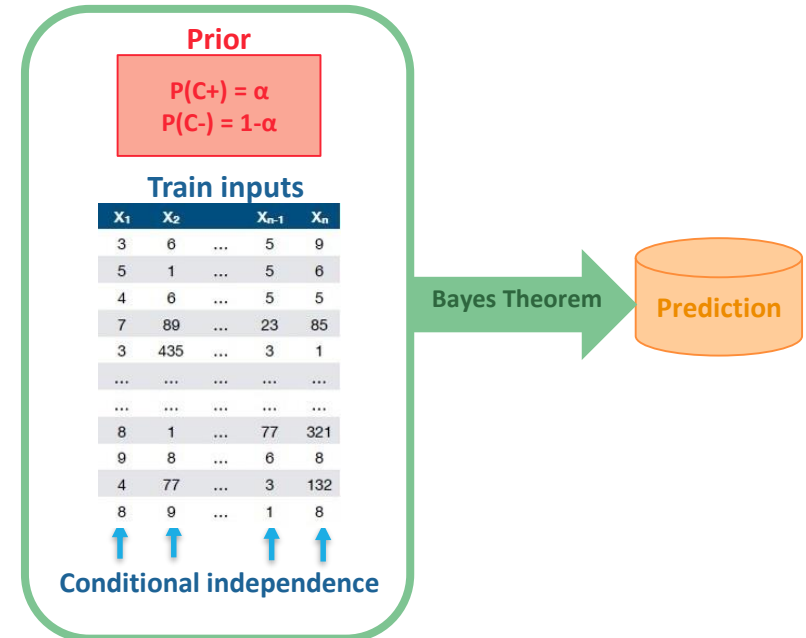Applying the chain rule in conditional probability

$$p(C, X_1, \ldots, X_n) = p(C) \cdot p(X_1|C) \cdot p(X_2|C, X_1) \cdot \ldots \cdot p(X_n|C, X_1, \ldots, X_{n-1})$$

Assuming conditional independence between $X_i$, then if $i \neq j$

$$p(X_i|C, X_j) = p(X_i|C)$$

**Prior**

| | |
|---|---|
| P(C+) = α | |
| P(C-) = 1-α | |

**Train inputs**

| X₁ | X₂ | | Xₙ₋₁ | Xₙ |
|---|---|---|---|---|
| 3 | 6 | ... | 5 | 9 |
| 5 | 1 | ... | 5 | 6 |
| 4 | 6 | ... | 5 | 5 |
| 7 | 89 | ... | 23 | 85 |
| 3 | 435 | ... | 3 | 1 |
| ... | ... | ... | ... | ... |
| ... | ... | ... | ... | ... |
| 8 | 1 | ... | 77 | 321 |
| 9 | 8 | ... | 6 | 8 |
| 4 | 77 | ... | 3 | 132 |
| 8 | 9 | ... | 1 | 8 |

**Conditional independence**

Bayes Theorem → **Prediction**

# Naïve Bayes

We need to calculate

$$p(C|X_1, \ldots, X_n)$$

Bayes Theorem,

Joint probability $\;p(C, X_1, \ldots, X_n)$

$$p(C|X_1, \ldots, X_n) = \frac{p(C) \cdot p(X_1, \ldots, X_n|C)}{p(X_1, \ldots, X_n)}$$

Applying the chain rule in conditional probability

$$p(C, X_1, \ldots, X_n) = p(C) \cdot p(X_1|C) \cdot p(X_2|C, X_1) \cdot \ldots \cdot p(X_n|C, X_1, \ldots, X_{n-1})$$

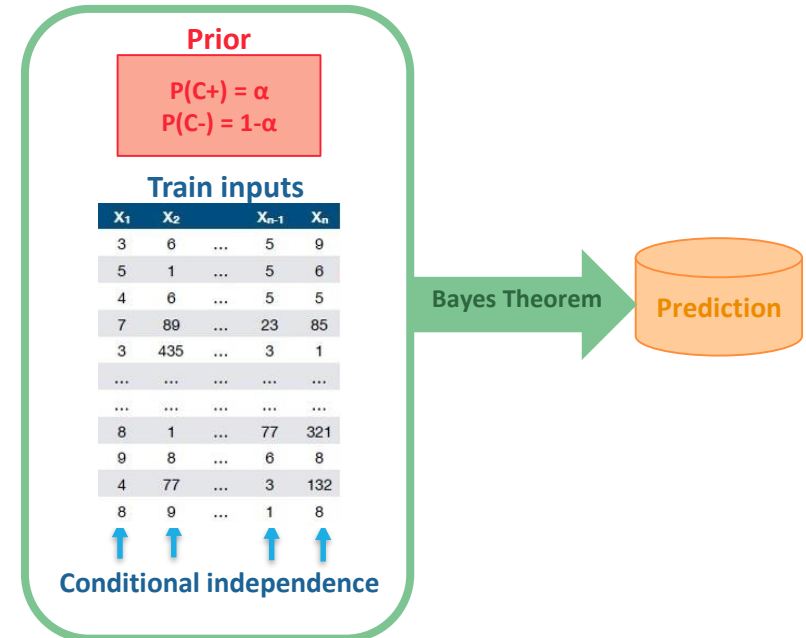Assuming conditional independence between $X_i$, then if $i \neq j$

$$p(X_i|C, X_j) = p(X_i|C)$$

Therefore,

$$p(C, X_1, \ldots, X_n) = p(C) \cdot \prod_{i=1}^{n} p(X_i|C)$$

**Prior**

P(C+) = α
P(C-) = 1-α

**Train inputs**

| X₁ | X₂ | | Xₙ₋₁ | Xₙ |
|----|----|----|----|----|
| 3 | 6 | ... | 5 | 9 |
| 5 | 1 | ... | 5 | 6 |
| 4 | 6 | ... | 5 | 5 |
| 7 | 89 | ... | 23 | 85 |
| 3 | 435 | ... | 3 | 1 |
| ... | ... | ... | ... | ... |
| ... | ... | ... | ... | ... |
| 8 | 1 | ... | 77 | 321 |
| 9 | 8 | ... | 6 | 8 |
| 4 | 77 | ... | 3 | 132 |
| 8 | 9 | ... | 1 | 8 |

**Conditional independence**

**Bayes Theorem** → **Prediction**

# Naïve Bayes

We need to calculate

$$p(C|X_1, \ldots, X_n) \qquad \xrightarrow{\text{Joint probability}} \qquad p(C, X_1, \ldots, X_n)$$

Bayes Theorem,

$$p(C|X_1, \ldots, X_n) = \frac{p(C) \cdot p(X_1, \ldots, X_n | C)}{p(X_1, \ldots, X_n)}$$

Applying the chain rule in conditional probability

$$p(C, X_1, \ldots, X_n) = p(C) \cdot p(X_1|C) \cdot p(X_2|C, X_1) \cdot \ldots \cdot p(X_n|C, X_1, \ldots, X_{n-1})$$

Assuming conditional independence between $X_i$, then if $i \neq j$
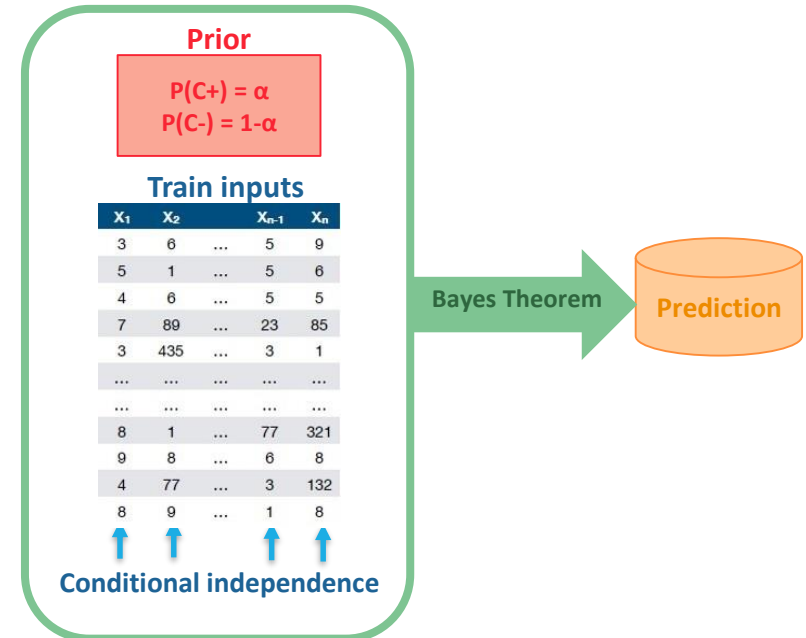
$$p(X_i | C, X_j) = p(X_i | C)$$

Therefore,

$$p(C, X_1, \ldots, X_n) = p(C) \cdot \prod_{i=1}^{n} p(X_i | C)$$

Then,

$$p(C|X_1, \ldots, X_n) \propto p(C) \cdot \prod_{i=1}^{n} p(X_i | C)$$

**Prior**

P(C+) = α
P(C-) = 1-α

**Train inputs**

| X₁ | X₂ | | Xₙ₋₁ | Xₙ |
|----|----|----|----|----|
| 3 | 6 | ... | 5 | 9 |
| 5 | 1 | ... | 5 | 6 |
| 4 | 6 | ... | 5 | 5 |
| 7 | 89 | ... | 23 | 85 |
| 3 | 435 | ... | 3 | 1 |
| ... | ... | ... | ... | ... |
| ... | ... | ... | ... | ... |
| 8 | 1 | ... | 77 | 321 |
| 9 | 8 | ... | 6 | 8 |
| 4 | 77 | ... | 3 | 132 |
| 8 | 9 | ... | 1 | 8 |

**Conditional independence**

**Bayes Theorem** → **Prediction**

# Naïve Bayes

We need to calculate

$$p(C|X_1,\ldots,X_n)$$

Bayes Theorem,

$$p(C|X_1,\ldots,X_n) = \frac{p(C) \cdot p(X_1,\ldots,X_n|C)}{p(X_1,\ldots,X_n)}$$

Joint probability $\;\to\; p(C, X_1,\ldots,X_n)$

Applying the chain rule in conditional probability

$$p(C,X_1,\ldots,X_n) = p(C) \cdot p(X_1|C) \cdot p(X_2|C,X_1) \cdot \ldots \cdot p(X_n|C,X_1,\ldots,X_{n-1})$$

Assuming conditional independence between $X_i$, then if $i \neq j$

$$p(X_i|C,X_j) = p(X_i|C)$$

Therefore,

$$p(C,X_1,\ldots,X_n) = p(C) \cdot \prod_{i=1}^{n} p(X_i|C)$$

Then,

$$p(C|X_1,\ldots,X_n) \propto p(C) \cdot \prod_{i=1}^{n} p(X_i|C)$$

where the proportionality constant is

$$\frac{1}{p(X_1,\ldots,X_n)}$$

**Prior**

P(C+) = α
P(C-) = 1-α

**Train inputs**

| X₁ | X₂ |  | Xₙ₋₁ | Xₙ |
|----|----|----|----|----|
| 3 | 6 | ... | 5 | 9 |
| 5 | 1 | ... | 5 | 6 |
| 4 | 6 | ... | 5 | 5 |
| 7 | 89 | ... | 23 | 85 |
| 3 | 435 | ... | 3 | 1 |
| ... | ... | ... | ... | ... |
| ... | ... | ... | ... | ... |
| 8 | 1 | ... | 77 | 321 |
| 9 | 8 | ... | 6 | 8 |
| 4 | 77 | ... | 3 | 132 |
| 8 | 9 | ... | 1 | 8 |

**Conditional independence**

**Bayes Theorem** → **Prediction**

# Naïve Bayes

We need to calculate

$$p(C|X_1, \ldots, X_n)$$

Bayes Theorem,

$$p(C|X_1, \ldots, X_n) = \frac{p(C) \cdot p(X_1, \ldots, X_n|C)}{p(X_1, \ldots, X_n)}$$

Joint probability $\rightarrow p(C, X_1, \ldots, X_n)$

Applying the chain rule in conditional probability

$$p(C, X_1, \ldots, X_n) = p(C) \cdot p(X_1|C) \cdot p(X_2|C, X_1) \cdot \ldots \cdot p(X_n|C, X_1, \ldots, X_{n-1})$$

Assuming conditional independence between $X_i$, then if $i \neq j$

$$p(X_i|C, X_j) = p(X_i|C)$$

Therefore,

$$p(C, X_1, \ldots, X_n) = p(C) \cdot \prod_{i=1}^{n} p(X_i|C)$$

Then,

$$p(C|X_1, \ldots, X_n) \propto p(C) \cdot \prod_{i=1}^{n} p(X_i|C)$$

where the proportionality constant is

$$\frac{1}{p(X_1, \ldots, X_n)}$$

**Prior**

P(C+) = α
P(C-) = 1-α

**Train inputs**

| X₁ | X₂ | | Xₙ₋₁ | Xₙ |
|---|---|---|---|---|
| 3 | 6 | ... | 5 | 9 |
| 5 | 1 | ... | 5 | 6 |
| 4 | 6 | ... | 5 | 5 |
| 7 | 89 | ... | 23 | 85 |
| 3 | 435 | ... | 3 | 1 |
| ... | ... | ... | ... | ... |
| ... | ... | ... | ... | ... |
| 8 | 1 | ... | 77 | 321 |
| 9 | 8 | ... | 6 | 8 |
| 4 | 77 | ... | 3 | 132 |
| 8 | 9 | ... | 1 | 8 |

**Conditional independence**

Sample data

Bayes Theorem → **Prediction**

# Support Vector Machines



Separable

# Support Vector Machines



**Separable**

**Not Separable**

# Support Vector Machines



**Separable**



**Not Separable**

**Kernel trick** $\phi$

**Input Space**

**Feature Space**

# Support Vector Machines



Kernel trick $\phi$

Input Space

Feature Space

Separable

Not Separable

Support vectors

ε-loss

# Support Vector Machines

Kernel trick $\phi$

Input Space

Feature Space

Separable

Not Separable
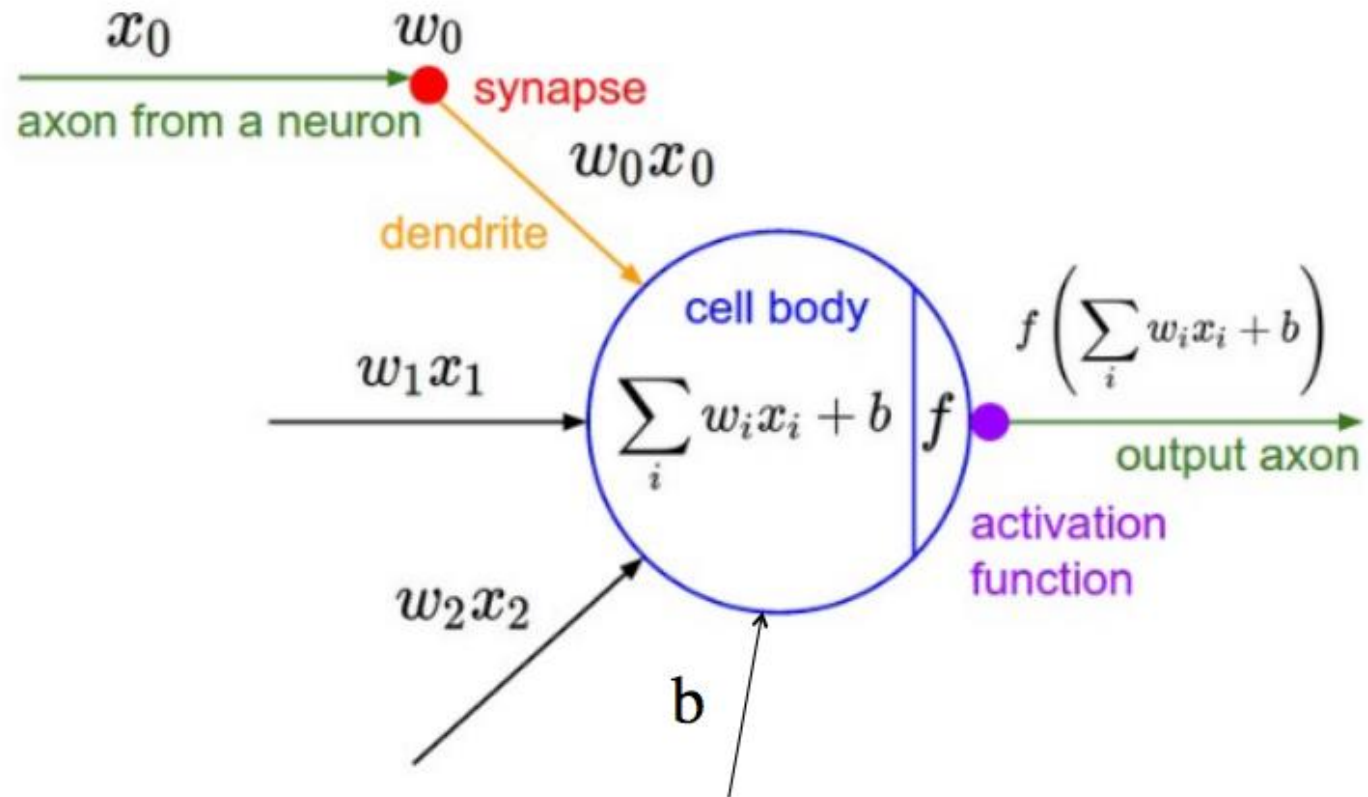
$$\arg \max_{\beta_0, \beta_1, \ldots, \beta_n} M$$

$$\text{subject to} \sum_{j=1}^{n} \beta_j^2 = 1,$$

$$y_i * (\beta_0 + \sum_j X_{ij}\beta_{ij}) \geq M(1 - \epsilon_i),$$
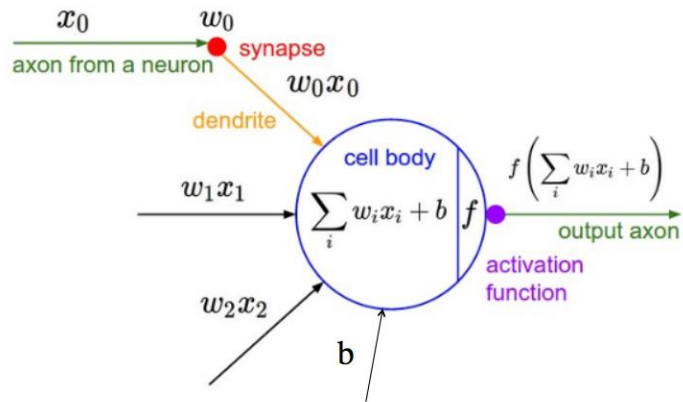
$$\epsilon_i \geq 0, \sum_i \epsilon_i \leq C$$

Support vectors

ε-loss

# Neural networks

- Neuron



Diagram showing inputs $x_0$ with weight $w_0$ along an axon from a neuron, through a synapse, along a dendrite as $w_0 x_0$. Inputs $w_1 x_1$ and $w_2 x_2$ enter the cell body which computes $\sum_i w_i x_i + b$, passes through activation function $f$, producing output $f\left(\sum_i w_i x_i + b\right)$ along the output axon. The bias is labeled $b$.
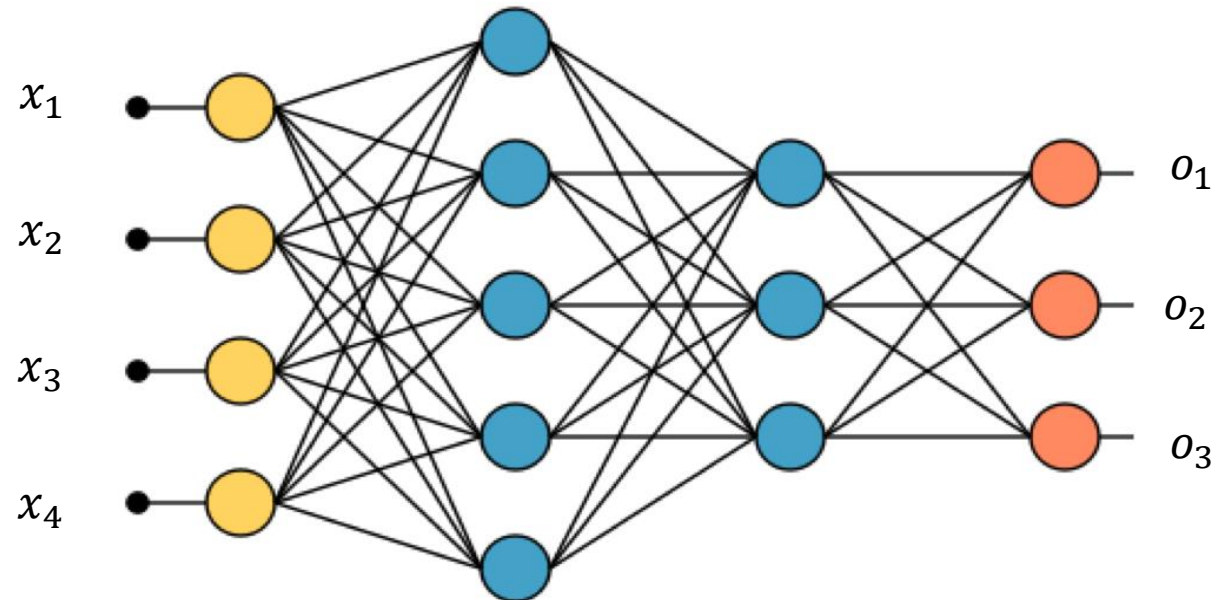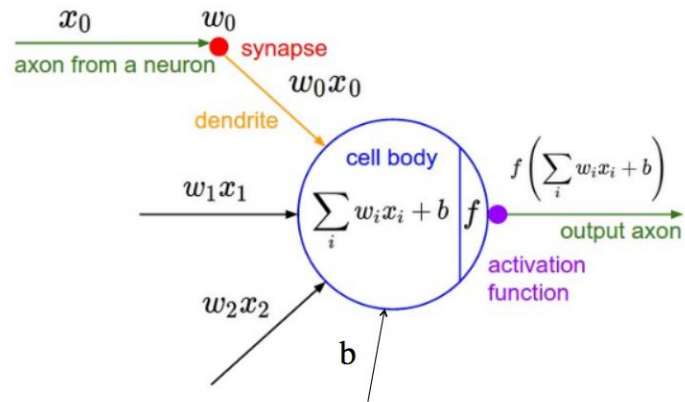
# Neural networks

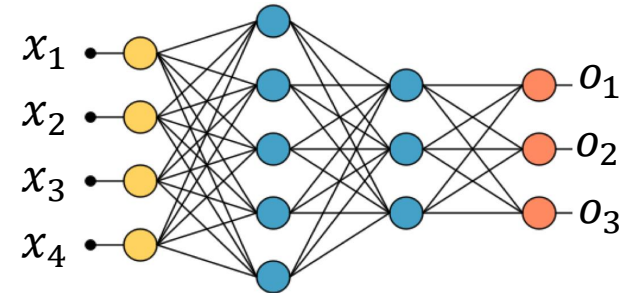- Neuron



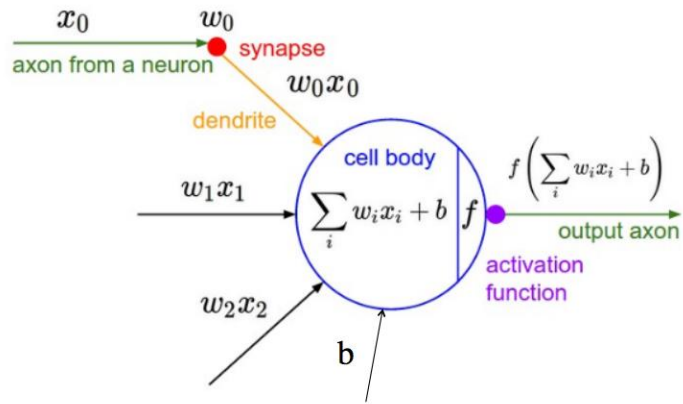- Neural Network

# Neural networks
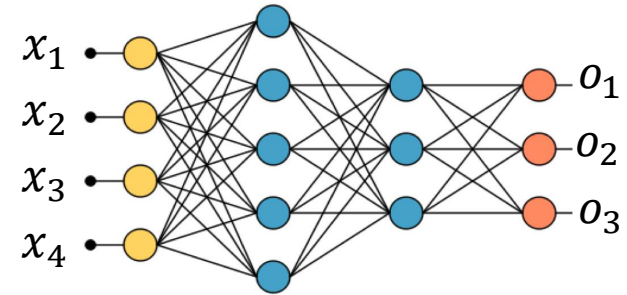
- Neuron



- Neural Network



The union, in reticular form, of neurons is **neural network**.

# Neural networks

- Neuron



- Neural Network



The union, in reticular form, of neurons is **neural network**.

It is divided into **layers**, having an **input layer**, one or more hidden layers and an output layer.

# Neural networks

- Neuron



- Neural Network



The union, in reticular form, of neurons is **neural network**.

It is divided into **layers**, having an **input layer**, one or more hidden layers and an output layer.

They are used as many neurons as variables in the input and as many as classes in the output.
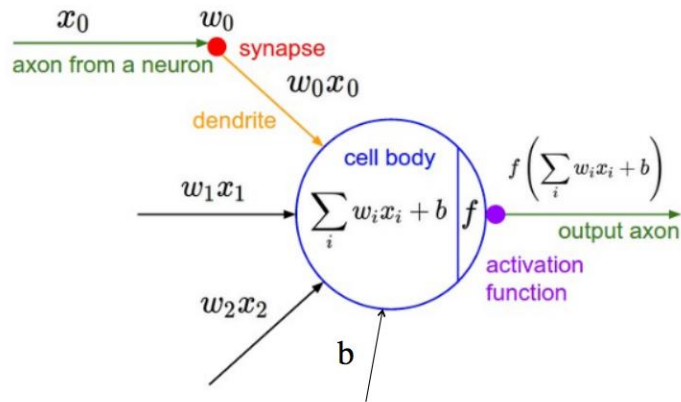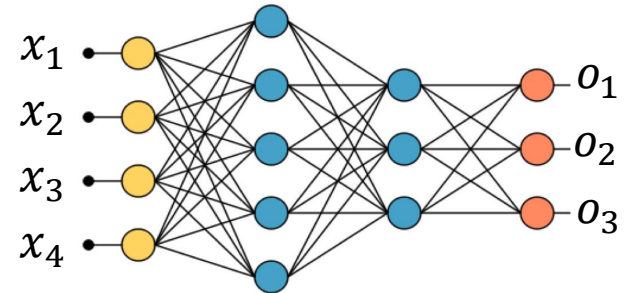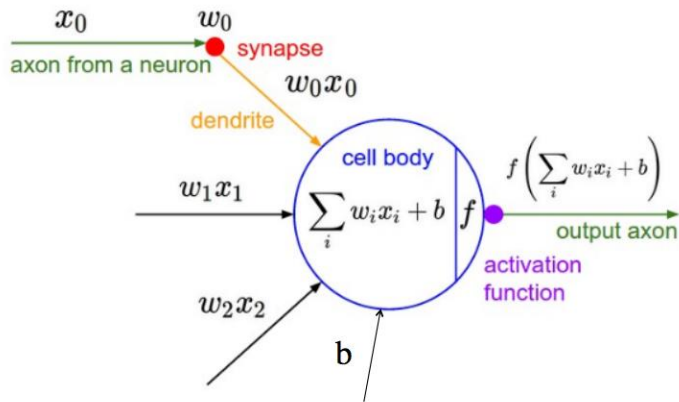
# Neural networks

- Neuron



- Neural Network



The union, in reticular form, of neurons is **neural network**.

It is divided into **layers**, having an **input layer**, one or more hidden layers and an output layer.

They are used as many neurons as variables in the input and as many as classes in the output.

- Usual activation functions



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

**Sigmoid**

$$\tanh(x) = \frac{2}{1 + e^{-2x}} - 1$$

**Hyperbolic tangent**

$$f(x) = \begin{cases} 0 & \text{for} \quad x < 0 \\ x & \text{for} \quad x \geq 0 \end{cases}$$

**ReLu**

# Neural networks

- Shallow Learning



Input 1

Input 2

Input 3

Output
Value

Input
Layer

Hidden
Layer

Output
Layer

# Neural networks

- ## Shallow Learning



Input 1
Input 2
Input 3

Output Value

Input Layer
Hidden Layer
Output Layer

- ## Deep Learning



Input 1
Input 2
Input 3

Output Value

Input Layer
Hidden Layer
Output Layer

# Ensemble learning

- Combination of several predictions for a given paradigm

# Ensemble learning

- Combination of several predictions for a given paradigm

- Different configurations and/or training data subsets are employed

# Ensemble learning

- Combination of several predictions for a given paradigm

- Different configurations and/or training data subsets are employed

- Results aggregation. Several ways depending of the type of problem or data

# Ensemble learning

- Combination of several predictions for a given paradigm

- Different configurations and/or training data subsets are employed

- Results aggregation. Several ways depending of the type of problem or data

- Variance reduction keeping bias

# Ensemble learning

- Combination of several predictions for a given paradigm

- Different configurations and/or training data subsets are employed

- Results aggregation. Several ways depending of the type of problem or data

- Variance reduction keeping bias

# Boosting & AdaBoost

- <u>Question</u>: Can a set of weak learners create a single strong learner? (Kearns & Valiant)

# Boosting & AdaBoost

- <u>Question</u>: Can a set of weak learners create a single strong learner? (Kearns & Valiant)
- <u>Answer</u>: Yes. Think of an analogy with a set of experts with different levels of expertise taking decisions vs a single oracle.

# Boosting & AdaBoost

- <u>Question</u>: Can a set of weak learners create a single strong learner? (Kearns & Valiant)
- <u>Answer</u>: Yes. Think of an analogy with a set of experts with different levels of expertise taking decisions vs a single oracle.

- Boosting steps

# Boosting & AdaBoost

- <u>Question</u>: Can a set of weak learners create a single strong learner? (Kearns & Valiant)

- <u>Answer</u>: Yes. Think of an analogy with a set of experts with different levels of expertise taking decisions vs a single oracle.

- Boosting steps:
  - Train the weak learners iteratively

# Boosting & AdaBoost

- <u>Question</u>: Can a set of weak learners create a single strong learner? (Kearns & Valiant)

- <u>Answer</u>: Yes. Think of an analogy with a set of experts with different levels of expertise taking decisions vs a single oracle.

- Boosting steps:
    - Train the weak learners iteratively
    - Aggregate their predictions using a weighted voting process, where the weights are adjusted using the weak kearners performance

# Boosting & AdaBoost

- <u>Question</u>: Can a set of weak learners create a single strong learner? (Kearns & Valiant)

- <u>Answer</u>: Yes. Think of an analogy with a set of experts with different levels of expertise taking decisions vs a single oracle.

- Boosting steps:
    - Train the weak learners iteratively
    - Aggregate their predictions using a weighted voting process, where the weights are adjusted using the weak kearners performance
    - Update the weights every time a new weak learner is trained

# Boosting & AdaBoost

- <u>Question</u>: Can a set of weak learners create a single strong learner? (Kearns & Valiant)

- <u>Answer</u>: Yes. Think of an analogy with a set of experts with different levels of expertise taking decisions vs a single oracle.

- Boosting steps:

  – Train the weak learners iteratively

  – Aggregate their predictions using a weighted voting process, where the weights are adjusted using the weak kearners performance

  – Update the weights every time a new weak learner is trained

- **Adaboost** stands for **adaptive boosting**

# Boosting & AdaBoost

- <u>Question</u>: Can a set of weak learners create a single strong learner? (Kearns & Valiant)

- <u>Answer</u>: Yes. Think of an analogy with a set of experts with different levels of expertise taking decisions vs a single oracle.

- Boosting steps:
  - Train the weak learners iteratively
  - Aggregate their predictions using a weighted voting process, where the weights are adjusted using the weak kearners performance
  - Update the weights every time a new weak learner is trained

- **Adaboost** stands for **adaptive boosting**

- It is a boosting algorithm in which the sequential aggregation of weak learners is guided by an **adaptive process**

# Boosting & AdaBoost

- <u>Question</u>: Can a set of weak learners create a single strong learner? (Kearns & Valiant)

- <u>Answer</u>: Yes. Think of an analogy with a set of experts with different levels of expertise taking decisions vs a single oracle.

- Boosting steps:
  – Train the weak learners iteratively
  – Aggregate their predictions using a weighted voting process, where the weights are adjusted using the weak kearners performance
  – Update the weights every time a new weak learner is trained

- **Adaboost** stands for **adaptive boosting**

- It is a boosting algorithm in which the sequential aggregation of weak learners is guided by an **adaptive process**

- Such process considers the error of each weak learner on each sample in order to **adapt the weights for the weak learners** to be aggregated, so that **the overall error decreases as much as possible**.

# Bagging

- **Bagging** stands for **boostrap aggregating**

# Bagging

- **Bagging** stands for **boostrap aggregating**

- Process

# Bagging

- **Bagging** stands for **boostrap aggregating**

- Process:

  Given a dataset $A$ with $n$ samples, a number of samples $n' \leq n$, and a number of models $m$,

# Bagging

- **Bagging** stands for **boostrap aggregating**

- Process:

   Given a dataset $A$ with $n$ samples, a number of samples $n' \leq n$, and a number of models $m$,

   – Sample $m$ datasets, with $n'$ samples each, by sampling with replacement

# Bagging

- **Bagging** stands for **boostrap aggregating**

- Process:

  Given a dataset $A$ with $n$ samples, a number of samples $n' \leq n$, and a number of models $m$,

  – Sample $m$ datasets, with $n'$ samples each, by sampling with replacement

  – For each dataset, create a predictive model

# Bagging

- **Bagging** stands for **boostrap aggregating**

- Process:

  Given a dataset $A$ with $n$ samples, a number of samples $n' \leq n$, and a number of models $m$,

  - Sample $m$ datasets, with $n'$ samples each, by sampling with replacement

  - For each dataset, create a predictive model

  - Aggregate the output of all $m$ models by voting

# Bagging

- **Bagging** stands for **boostrap aggregating**

- Process:

  Given a dataset $A$ with $n$ samples, a number of samples $n' \leq n$, and a number of models $m$,

  - Sample $m$ datasets, with $n'$ samples each, by sampling with replacement

  - For each dataset, create a predictive model

  - Aggregate the output of all $m$ models by voting

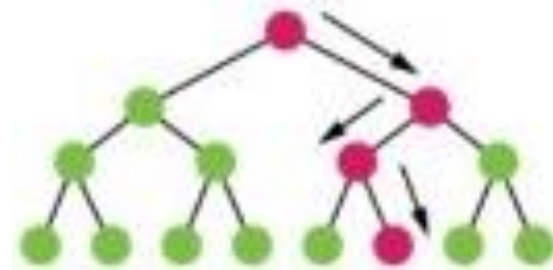- The most famous bagging-type algorithm is **Random Forests**

# Bagging

- **Bagging** stands for **boostrap aggregating**

- Process:

  Given a dataset $A$ with $n$ samples, a number of samples $n' \leq n$, and a number of models $m$,

  – Sample $m$ datasets, with $n'$ samples each, by sampling with replacement

  – For each dataset, create a predictive model

  – Aggregate the output of all $m$ models by voting

- The most famous bagging-type algorithm is **Random Forests**

- More recently, new approaches arised, such as **Gradient Boosting classifiers**, being **Extreme Gradient Boosting** (XGBoost) the new ML rockstar

# Random Forests
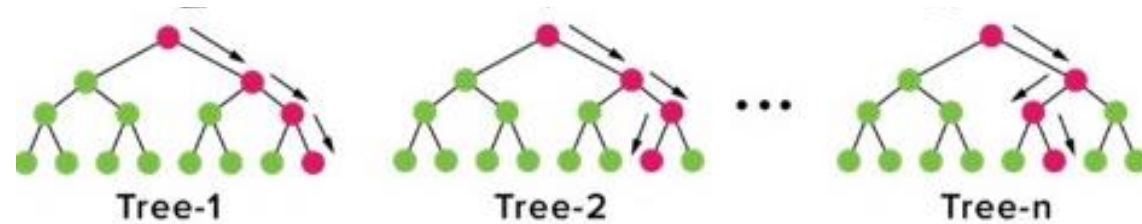
- Ensemble learning for decision trees

# Random Forests

- Ensemble learning for decision trees

- The number $n$ of trees is chosen

# Random Forests

- Ensemble learning for decision trees

- The number $n$ of trees is chosen

- For each tree



Tree-1       Tree-2  •••  Tree-n

# Random Forests

- Ensemble learning for decision trees

- The number $n$ of trees is chosen

- For each tree:

  - A **subset of $L$ variables** are randomly chosen. By default, $L = \left\lceil \sqrt{M} \right\rceil$

Tree-1    Tree-2    · · ·    Tree-n

# Random Forests

- Ensemble learning for decision trees

- The number $n$ of trees is chosen

- For each tree:

  - A **subset of $L$ variables** are randomly chosen. By default, $L = \lceil \sqrt{M} \rceil$

  - $K$ samples are chosen by **random sampling with reposition**. By default, $K = N$



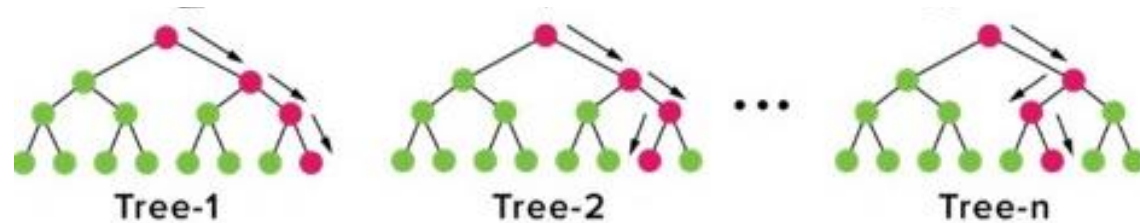Tree-1          Tree-2          Tree-n

# Random Forests

- Ensemble learning for decision trees

- The number $n$ of trees is chosen

- For each tree:

  – A **subset of $L$ variables** are randomly chosen. By default, $L = \lceil \sqrt{M} \rceil$

  – $K$ samples are chosen by **random sampling with reposition**. By default, $K = N$

- All trees are trained



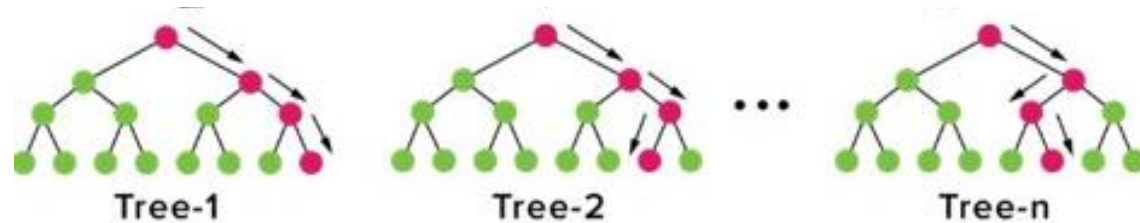Tree-1          Tree-2    ...    Tree-n

# Random Forests

- Ensemble learning for decision trees

- The number $n$ of trees is chosen

- For each tree:

  - A **subset of $L$ variables** are randomly chosen. By default, $L = \lceil \sqrt{M} \rceil$

  - $K$ samples are chosen by **random sampling with reposition**. By default, $K = N$

- All trees are trained

- For each sample, $n$ predictions are obtained



24.10.19

# Random Forests

- Ensemble learning for decision trees

- The number $n$ of trees is chosen

- For each tree:

    - A **subset of $L$ variables** are randomly chosen. By default, $L = \lceil \sqrt{M} \rceil$

    - $K$ samples are chosen by **random sampling with reposition**. By default, $K = N$
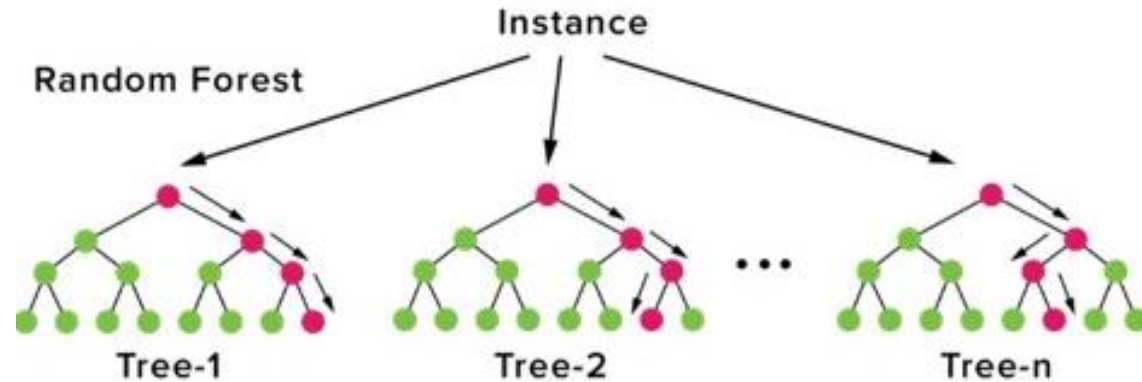
- All trees are trained

- For each sample, $n$ predictions are obtained

# Random Forests

- Ensemble learning for decision trees

- The number $n$ of trees is chosen

- For each tree:

  - A **subset of $L$ variables** are randomly chosen. By default, $L = \lceil \sqrt{M} \rceil$

  - $K$ samples are chosen by **random sampling with reposition**. By default, $K = N$
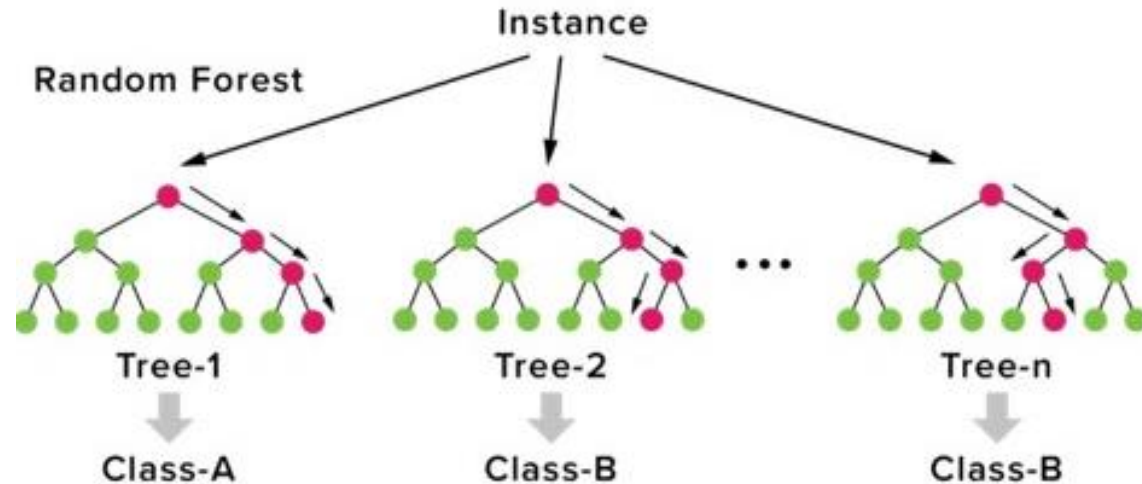
- All trees are trained

- For each sample, $n$ predictions are obtained

- Predictions are **aggregated** by **majority voting**. Different possible options

# Random Forests

- Ensemble learning for decision trees

- The number $n$ of trees is chosen

- For each tree:

    - A **subset of $L$ variables** are randomly chosen. By default, $L = \lceil \sqrt{M} \rceil$

    - $K$ samples are chosen by **random sampling with reposition**. By default, $K = N$

- All trees are trained
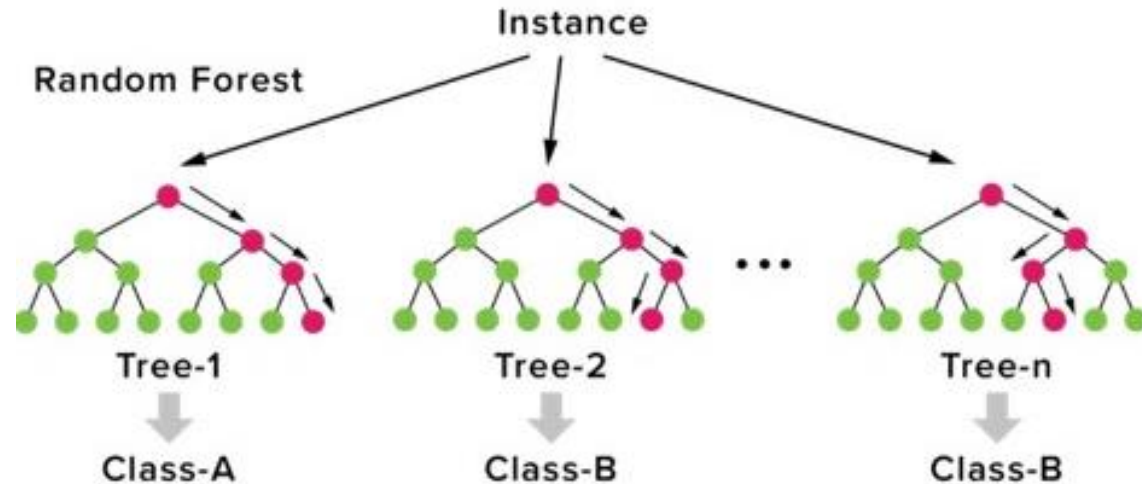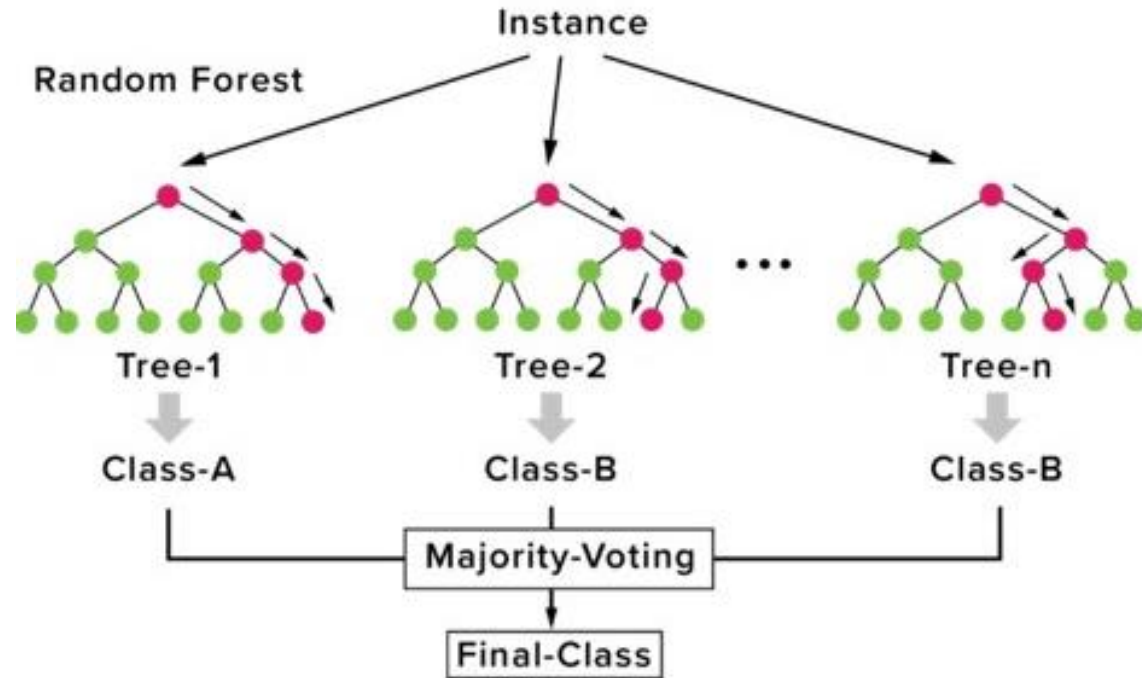
- For each sample, $n$ predictions are obtained

- Predictions are **aggregated** by **majority voting**. Different possible options

# Validation

- We must distinguish binary from multiclass classification

# Validation

- We must distinguish binary from multiclass classification

- All classification metrics seen in FAA are valid for the binary case

# Validation

- We must distinguish binary from multiclass classification

- All classification metrics seen in FAA are valid for the binary case

- Only a few are directly applicable to multiclass situations

# Validation

- We must distinguish binary from multiclass classification

- All classification metrics seen in FAA are valid for the binary case

- Only a few are directly applicable to multiclass situations

|  | Predicted | | | |
|---|---|---|---|---|
|  | Iris-setosa | Iris-versicolor | Iris-virginica | Σ |
| Iris-setosa | 100.0 % | 0.0 % | 0.0 % | 50 |
| Iris-versicolor | 0.0 % | 88.7 % | 6.4 % | 50 |
| Iris-virginica | 0.0 % | 11.3 % | 93.6 % | 50 |
| Σ | 50 | 53 | 47 | 150 |

# Validation

- We must distinguish binary from multiclass classification

- All classification metrics seen in FAA are valid for the binary case

- Only a few are directly applicable to multiclass situations

|  | Predicted | | | |
|---|---|---|---|---|
|  | Iris-setosa | Iris-versicolor | Iris-virginica | Σ |
| Iris-setosa | 100.0 % | 0.0 % | 0.0 % | 50 |
| Iris-versicolor | 0.0 % | 88.7 % | 6.4 % | 50 |
| Iris-virginica | 0.0 % | 11.3 % | 93.6 % | 50 |
| Σ | 50 | 53 | 47 | 150 |

(Actual)

- All scores are usable by means of 2 possible simplifications:

# Validation

- We must distinguish binary from multiclass classification

- All classification metrics seen in FAA are valid for the binary case

- Only a few are directly applicable to multiclass situations

|  | Predicted | | | |
|---|---|---|---|---|
| | Iris-setosa | Iris-versicolor | Iris-virginica | Σ |
| Iris-setosa | 100.0 % | 0.0 % | 0.0 % | 50 |
| Iris-versicolor | 0.0 % | 88.7 % | 6.4 % | 50 |
| Iris-virginica | 0.0 % | 11.3 % | 93.6 % | 50 |
| Σ | 50 | 53 | 47 | 150 |

- All scores are usable by means of 2 possible simplifications:
  - <u>One-vs-All</u>: If we have k classes, we train k binary model predicting each class against the rest iteratively

# Validation

- We must distinguish binary from multiclass classification

- All classification metrics seen in FAA are valid for the binary case

- Only a few are directly applicable to multiclass situations

| | | Predicted | | |
| --- | --- | --- | --- | --- |
| | Iris-setosa | Iris-versicolor | Iris-virginica | Σ |
| Iris-setosa | 100.0 % | 0.0 % | 0.0 % | 50 |
| Iris-versicolor | 0.0 % | 88.7 % | 6.4 % | 50 |
| Iris-virginica | 0.0 % | 11.3 % | 93.6 % | 50 |
| Σ | 50 | 53 | 47 | 150 |

- All scores are usable by means of 2 possible simplifications:
  - One-vs-All: If we have k classes, we train k binary model predicting each class against the rest iteratively
  - One-vs-One: If we have k classes, we train all pairwise possible binary models. In total, the total number of binary models is k(k-1)/2

# Validation

- We must distinguish binary from multiclass classification

- All classification metrics seen in FAA are valid for the binary case

- Only a few are directly applicable to multiclass situations

|  | Predicted | | | |
| --- | --- | --- | --- | --- |
|  | Iris-setosa | Iris-versicolor | Iris-virginica | Σ |
| Iris-setosa | 100.0 % | 0.0 % | 0.0 % | 50 |
| Iris-versicolor | 0.0 % | 88.7 % | 6.4 % | 50 |
| Iris-virginica | 0.0 % | 11.3 % | 93.6 % | 50 |
| Σ | 50 | 53 | 47 | 150 |

(Actual)

- All scores are usable by means of 2 possible simplifications:
  - One-vs-All: If we have k classes, we train k binary model predicting each class against the rest iteratively
  - One-vs-One: If we have k classes, we train all pairwise possible binary models. In total, the total number of binary models is k(k-1)/2

- In both cases, numerical scores are **aggregated by averaging**

# Python functions

| Algoritmo | Entorno | Función |
| --- | --- | --- |
| K-Nearest Neighbors | Scikit-learn | KNeighborsClassifier |
| Decision tree | Scikit-learn | DecisionTreeClassifier |
| Logistic Regression | Scikit-learn | LogisticRegression |
| Naïve Bayes | Scikit-learn | GaussianNB<br>MultinomialNB<br>ComplementNB<br>BernoulliNB |
| Support Vector Machines | Scikit-learn | SVC<br>NuSVC<br>LinearSVC |
| Neural Networks | Scikit-learn | MLPClassifier |
| AdaBoost | Scikit-learn | AdaBoostClassifier |
| Bagging | Scikit-learn | BaggingClassifier |
| Random Forests | Scikit-learn | RandomForestClassifier |
| Gradient Boosting | Scikit-learn | GradientBoostingClassifier<br><br>HistGradientBoostingClassifier |

Eskerrik asko
Muchas gracias
Thank you

**Carlos Cernuda**

ccernuda@mondragon.edu

MGEP
Goiru, 2
20500 Arrasate – Mondragon
Tlf. 662420414