

# Lecture 6: Naive Bayes classifier

## Introduction to Machine Learning

Sophie Robert

L3 MIASHS — Semestre 2

2022-2023

## 1 Principle

## 2 Mathematical framework

## 3 Example: Dog breed prediction

## 4 Hyperparameters

## 5 Advantages and limits

# Reminders on Bayes theorem

Lecture 6:  
Naive Bayes  
classifier

Sophie Robert

Principle

Mathematical  
framework

Example: Dog  
breed  
prediction

Hyperparameters

Advantages  
and limits

## Question

Can anyone remind me of the Bayes theorem ?

$$\mathbb{P}(A|B) = \frac{\mathbb{P}(B|A) \times \mathbb{P}(A)}{\mathbb{P}(B)}$$

# Main idea

Lecture 6:  
Naïve Bayes  
classifier

Sophie Robert

Principle

Mathematical  
framework

Example: Dog  
breed  
prediction

Hyperparameters

Advantages  
and limits

## The naïve Bayes classifier algorithm

The naïve Bayes classifier algorithm is a probabilistic classifier which consists in applying to each record the class which is the most probable according to a probability model.

Given a record  $\mathbf{x} = (x_1, x_2, \dots, x_n) \in \mathbb{R}^n$  and a set of labels  $y_i \in \mathcal{Y}$ , provide an estimation of  $\mathbb{P}(y_i|\mathbf{x})$   $y_i \in \mathcal{Y}$  (and assign most likely label to  $\mathbf{x}$ ).

# Mathematical framework

Lecture 6:  
Naive Bayes  
classifier

Sophie Robert

Principle

Mathematical  
framework

Example: Dog  
breed  
prediction

Hyperparameters

Advantages  
and limits

We want to estimate for each label  $y_i \in \mathcal{Y}$   $\mathbb{P}(y_i|\mathbf{x})$  (*what is the probability of being label  $y_i$  given the data records ?*) .

However, if  $n$  is large, the computation is infeasible.

Using the definition of conditional probabilities:

$$\mathbb{P}(y_i|\mathbf{x}) = \frac{\mathbb{P}(y_i, \mathbf{x})}{\mathbb{P}(\mathbf{x})}$$

$\mathbb{P}(\mathbf{x})$  is a constant because  $\mathbf{x}$  is given so we only need to find the value of  $\mathbb{P}(y_i, \mathbf{x})$ .

# Mathematical framework

Lecture 6:  
Naive Bayes  
classifier

Sophie Robert

Principle

Mathematical  
framework

Example: Dog  
breed  
prediction

Hyperparameters

Advantages  
and limits

Using the definition of conditional probabilities iteratively,

$$\begin{aligned}\mathbb{P}(y_i, \mathbf{x}) &= \mathbb{P}(x_1, x_2, \dots, x_n, y_i) \\ &= \mathbb{P}(x_1 | x_2, x_3, \dots, y_i) \times \mathbb{P}(x_2, x_4, \dots, y_i) \\ &= \mathbb{P}(x_1 | x_2, x_3, \dots, y_i) \times \mathbb{P}(x_2 | x_3, x_4, \dots, y_i) \times \mathbb{P}(x_3, x_4, \dots, y_i) \\ &= \dots \\ &= \mathbb{P}(x_1 | x_2, x_3, \dots, y_i) \times \mathbb{P}(x_2 | x_3, x_4, \dots, y_i) \times \mathbb{P}(x_n | y_i) \times \mathbb{P}(y_i)\end{aligned}$$

# Mathematical framework

Lecture 6:  
Naive Bayes  
classifier

Sophie Robert

Principle

Mathematical  
framework

Example: Dog  
breed  
prediction

Hyperparameters

Advantages  
and limits

We now make the hypothesis that each feature  $x_i$  is independant (and only depends on the label  $y_i$ ):

$$\mathbb{P}(x_1|x_2, x_3, \dots, y_i) = \mathbb{P}(x_1|y_i)$$

We now have:

$$\mathbb{P}(y_i, \mathbf{x}) = \mathbb{P}(y_i) \prod_{j=1}^n \mathbb{P}(x_j|y_i)$$

and :

$$\mathbb{P}(y_i|\mathbf{x}) \propto \mathbb{P}(y_i) \prod_{j=1}^n \mathbb{P}(x_j|y_i)$$

We then select the most probable class

$$\hat{y} = \underset{i=1, \dots, k}{\operatorname{argmax}} (\mathbb{P}(y_i) \prod_{j=1}^n \mathbb{P}(x_j|y_i))$$

# Mathematical framework

Lecture 6:  
Naive Bayes  
classifier

Sophie Robert

Principle

Mathematical  
framework

Example: Dog  
breed  
prediction

Hyperparameters

Advantages  
and limits

Can you guess why this algorithm can be called naive ?

We have two terms to estimate:

- $\mathbb{P}(y_i)$ : either assume class equiprobability or estimate using the frequency in training dataset
- $\mathbb{P}(x_j|y_i)$ : we need to decide on a conditional law



# Mathematical framework

Possible assumptions include:

- If  $X_j$  is a **continuous variable** ( $x_j \in \mathbb{R}$ ), the continuous values associated within class  $i$  are distributed according to a Gaussian distribution parametrized with mean  $\mu_i$  and variance  $\sigma_i$

$$f(v | y_i) = \frac{1}{\sqrt{2\pi\sigma_i^2}} e^{-\frac{(v-\mu_i)^2}{2\sigma_i^2}}$$

- If  $X_j$  is a **binary variable** ( $x_j \in \{0, 1\}^n$ ), the proportion of binary values observed within class  $y_i$  can be treated as a multivariate Bernoulli ( $p_{ij}$  being the frequency of event for variable  $x_j$  within class  $i$ ):

$$\mathbb{P}(x_j | y_i) = \prod_{j=1}^n p_{ij}^{x_j} (1 - p_{ij})^{(1-x_j)}$$

# Example

Lecture 6:  
Naive Bayes  
classifier

Sophie Robert

Principle

Mathematical  
framework

Example: Dog  
breed  
prediction

Hyperparameters

Advantages  
and limits

## Training dataset:

Height	Weight	Tail	Label
45	30	0	Labradoodle
30	25	1	Labradoodle
40	35	1	Labradoodle
20	15	0	English cocker
22	18	1	English cocker
25	20	1	English cocker

## Individual to classify

Height	Weight	Tail	Label
25	31	1	?

# Example: solution

## Estimate:

$$\begin{aligned} & \mathbb{P}(\text{labradoodle} \mid \text{height} = 25, \text{weight} = 31, \text{tail} = 1) \\ & \propto \mathbb{P}(\text{labradoodle}) \times \mathbb{P}(\text{height} = 25 \mid \text{labradoodle}) \\ & \times \mathbb{P}(\text{weight} = 31 \mid \text{labradoodle}) \\ & \times \mathbb{P}(\text{tail} = 1 \mid \text{labradoodle}) \end{aligned}$$

$$\mathbb{P}(\text{labradoodle}) = \frac{1}{2}$$

$$\mathbb{P}(\text{height} = 25 \mid \text{labradoodle}) = \frac{1}{\sqrt{2\pi \times 38.89}} e^{-\frac{(25-38.33)^2}{2 \times 38.89}} = 0.006$$

$$\mathbb{P}(\text{weight} = 31 \mid \text{labradoodle}) = \frac{1}{\sqrt{2\pi \times 16.67}} e^{-\frac{(31-30)^2}{2 \times 16.67}} = 0.09$$

$$\mathbb{P}(\text{tail} = 1 \mid \text{labradoodle}) = \frac{2}{3}$$

$$\mathbb{P}(\text{labradoodle} \mid \text{height} = 25, \text{weight} = 31, \text{tail} = 1) = 0.00017$$

# Example: solution

## Estimate:

$$\begin{aligned}\mathbb{P}(\text{cocker} \mid \text{height} = 25, \text{weight} = 31, \text{tail} = 1) \\ \propto \mathbb{P}(\text{cocker}) \times \mathbb{P}(\text{height} = 25 \mid \text{cocker}) \\ \times \mathbb{P}(\text{weight} = 31 \mid \text{cocker}) \\ \times \mathbb{P}(\text{tail} = 1 \mid \text{cocker})\end{aligned}$$

$$\mathbb{P}(\text{cocker}) = \frac{1}{2}$$

$$\mathbb{P}(\text{height} = 25 \mid \text{cocker}) = \frac{1}{\sqrt{2\pi \times 4.22}} e^{-\frac{(25-22.33)^2}{2 \times 4.22}} = 0.08$$

$$\mathbb{P}(\text{weight} = 31 \mid \text{cocker}) = \frac{1}{\sqrt{2\pi \times 16.67}} e^{-\frac{(31-30)^2}{2 \times 16.67}} = 1.39e - 10$$

$$\mathbb{P}(\text{tail} = 1 \mid \text{cocker}) = \frac{2}{3}$$

$$\mathbb{P}(\text{cocker} \mid \text{height} = 25, \text{weight} = 31, \text{tail} = 1) = 0.00$$

# Hyperparameters

Lecture 6:  
Naive Bayes  
classifier

Sophie Robert

Principle

Mathematical  
framework

Example: Dog  
breed  
prediction

Hyperparameter

Advantages  
and limits

## Hyperparameters

What **hyperparameters\*** do the naive Bayes classifier require ?

# Advantages and limits

Lecture 6:  
Naive Bayes  
classifier

Sophie Robert

Principle

Mathematical  
framework

Example: Dog  
breed  
prediction

Hyperparameters

Advantages  
and limits

## Limits:

- Strong independence hypothesis (but in practice, naive bayes behave rather well)
- Unable to classify unknown classes that do not show in training set (always sets it to 0, except in the case of artificial equiprobability)

## Advantages:

- Extends naturally to multi-class
- Naturally deals with categorical variables

# Questions

Lecture 6:  
Naive Bayes  
classifier

Sophie Robert

Principle

Mathematical  
framework

Example: Dog  
breed  
prediction

Hyperparameters

Advantages  
and limits

Questions ?