

# Lecture 6: Naive Bayes classifier

## Introduction to Machine Learning

Sophie Robert

L3 MIASHS | Semestre 2

2023-2024

- 1 Probabilistic classifiers
- 2 Principle of Bayes classifier
- 3 Mathematical framework
- 4 Example: Dog breed prediction
- 5 Hyperparameters
- 6 Advantages and limits
- 7 Further algorithms

## Probabilistic classifiers

# Probabilistic classifiers

## Probabilistic classifiers

Probabilistic classifiers are classifiers that predict, given an observation of an input, a **probability distribution over a set of classes** (instead of simply the class like standard classifiers).

# Probabilistic classifiers

## Probabilistic classifiers

Probabilistic classifiers are classifiers that predict, given an observation of an input, a **probability distribution over a set of classes** (instead of simply the class like standard classifiers).

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  $\mathbf{x}$ ).

# Probabilistic classifiers

## Question

What is in your opinion one of the strength of working with a probabilistic model rather than a classification function ?

# Probabilistic classifiers

## Question

What is in your opinion one of the strength of working with a probabilistic model rather than a classification function ?

Possible native models are:

- Logistic regression
- Subtypes of neural networks
- Bayes classifiers

# Probabilistic classifiers

## Question

What is in your opinion one of the strength of working with a probabilistic model rather than a classification function ?

Possible native models are:

- Logistic regression
- Subtypes of neural networks
- Bayes classifiers

Non-probabilistic models can also be turned into a probabilistic classifier (SVM, trees ...).

## Principle of Bayes classifier

# Main idea

## The naïve Bayes classifier algorithm

The naïve Bayes classifier algorithm is a probabilistic classifier based on applying Bayes' theorem with independence assumptions between the features.

# Main idea

## The naïve Bayes classifier algorithm

The naïve Bayes classifier algorithm is a probabilistic classifier based on applying Bayes' theorem with independence assumptions between the features.

Each feature **contributes to the class probability**

**independently**: for example, the size and the weight of the dog contribute independently to its breed.

## Mathematical framework

## Mathematical framework

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 ?*) .

## Mathematical framework

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.

## Mathematical framework

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:

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

$P(\mathbf{x})$  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

Using the definition of conditional probabilities iteratively,

# Mathematical framework

Using the definition of conditional probabilities iteratively,

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

## Mathematical framework

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

### Conditional independence

**Conditional independence** describes situations where an observation is irrelevant: the probability of the hypothesis given the uninformative observation is equal to the probability without. If A is the hypothesis, B and C the observations,

$$P(A | B, C) = P(A | C)$$

## Mathematical framework

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

### Conditional independence

**Conditional independence** describes situations where an observation is irrelevant: the probability of the hypothesis given the uninformative observation is equal to the probability without.

If A is the hypothesis, B and C the observations,

$$P(A | B, C) = P(A | C)$$

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

# Mathematical framework

We now have:

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

and :

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

# Mathematical framework

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

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

# Mathematical framework

Can you guess why this algorithm can be called naive ?

# Mathematical framework

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  $\sigma_i^2$

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

# 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  $\sigma_i^2$

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

- If  $X_j$  is a **categorial variable** ( $x_j \in \{0, \dots, K\}$ ), the probability can be estimated as the proportion of values within class:

$$\mathbb{P}(x = j | y_i) = \frac{N_{ji}}{N_i}$$

Example: Dog breed prediction

# Example

## Training dataset:

Height	Weight	Tail	Label
45	30	0	Legendary
30	25	1	Legendary
40	35	1	Legendary
20	15	0	Non legendary
22	18	1	Non legendary
25	20	1	Non legendary

## Individual to classify

Height	Weight	Tail	Label
25	31	1	?

## Example: training the model

### Training the model

Training the model consists in computing **statistical estimators over the population**.

For the legendary population:

	Height	Weight
Mean	38.33	30
Var	38.89	16.67

For the non-legendary population:

	Height	Weight
Mean	22.33	17.66
Var	4.22	4.22

## Example: solution

**Estimate:**

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

$$\propto \mathbb{P}(\text{legendary}) \times \mathbb{P}(\text{height} = 25 \mid \text{legendary})$$

$$\times \mathbb{P}(\text{weight} = 31 \mid \text{legendary})$$

$$\times \mathbb{P}(\text{tail} = 1 \mid \text{legendary})$$

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

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

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

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

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

## Example: solution

### Estimate:

$$\mathbb{P}(\text{non-legends} | \text{height} = 25, \text{weight} = 31, \text{tail} = 1)$$

$$\propto \mathbb{P}(\text{non-legends}) \times \mathbb{P}(\text{height} = 25 | \text{non-legends})$$

$$\times \mathbb{P}(\text{weight} = 31 | \text{non-legends})$$

$$\times \mathbb{P}(\text{tail} = 1 | \text{non-legends})$$

$$\mathbb{P}(\text{non-legends}) = \frac{1}{2}$$

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

$$\mathbb{P}(\text{weight} = 31 | \text{non-legends}) = \frac{1}{\sqrt{2\pi \times 4.22}} e^{-\frac{(31 - 17.66)^2}{2 \times 4.22}} = 1.31e - 11$$

$$\mathbb{P}(\text{tail} = 1 | \text{non-legends}) = \frac{2}{3}$$

$$\mathbb{P}(\text{non-legends} | \text{height} = 25, \text{weight} = 31, \text{tail} = 1) \propto 0.00$$

# Hyperparameters

# Hyperparameters

## Hyperparameters

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

## Advantages and limits

# Advantages and limits

## Limits:

# Advantages and limits

## Limits:

- Strong independance 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 and limits

## Limits:

- Strong independance 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:

# Advantages and limits

## Limits:

- Strong independance 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

## Further algorithms

## Other probabilistic classifiers

One of the most famous probabilistic classifier is **logistic regression**: probability distribution is expressed as the logit of the linear combination of features.

## Other probabilistic classifiers

One of the most famous probabilistic classifier is **logistic regression**: probability distribution is expressed as the logit of the linear combination of features.

Using **softmax function** as activation layer in **neural networks** transforms output into a probability distribution consisting of K probabilities proportional to the exponentials of the input numbers.

# Questions

Questions ?