# Lecture 6: Naive Bayes classifier
## Introduction to Machine Learning

Sophie Robert

L3 MIASHS — Semestre 2

2022-2023

Lecture 6:
Naive Bayes
classifier

Sophie Robert

Probabilistic
classifiers

Principle of
Bayes classifier

Mathematical
framework

Example: Dog
breed
prediction

Hyperparameters

Advantages
and limits

Further
algorithms

## 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

Lecture 6:
Naive Bayes
classifier

Sophie Robert

Probabilistic
classifiers

Principle of
Bayes classifier

Mathematical
framework

Example: Dog
breed
prediction

Hyperparameters

Advantages
and limits

Further
algorithms

## 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 $\mathbf{x} = (x_1, x_2, ..., x_n) \in \mathbb{R}^n$ 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 $\mathbf{x}$).

### Question

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

# Probabilistic classifiers

Lecture 6:
Naive Bayes
classifier

Sophie Robert

Probabilistic
classifiers

Principle of
Bayes classifier

Mathematical
framework

Example: Dog
breed
prediction

Hyperparameters

Advantages
and limits

Further
algorithms

## 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 . . . ).

## The naïve Bayes classifier algorithm

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

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

### Question

Can anyone remind me of the Bayes theorem ?

## Question

Can anyone remind me of the Bayes theorem ?

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

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 ?*) .

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 $\mathbb{P}(y_i, \mathbf{x})$.

Using the definition of conditional probabilities iteratively,

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

# Mathematical framework

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

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

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

$$\hat{y} = argmax_{i=1,...,k}(\mathbb{P}(y_i) \prod_{j=1}^{n} \mathbb{P}(x_j|y_i))$$

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

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 $\sigma_i$

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

# 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 $\sigma_i$

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

- **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 Bernouilli ($p_{ij}$ being the frequency of event for variable $x_j$ within class $i$):

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

# Example

**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:**

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

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

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

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

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

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

# Example: solution

Lecture 6:
Naive Bayes
classifier

Sophie Robert

Probabilistic
classifiers

Principle of
Bayes classifier

Mathematical
framework

Example: Dog
breed
prediction

Hyperparameters

Advantages
and limits

Further
algorithms

**Estimate:**

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

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

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

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

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

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

# Hyperparameters

## Hyperparameters

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

**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)

**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

# Other probabilistic classifiers

Lecture 6:
Naive Bayes
classifier

Sophie Robert

Probabilistic
classifiers

Principle of
Bayes classifier

Mathematical
framework

Example: Dog
breed
prediction

Hyperparameters

Advantages
and limits

Further
algorithms

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 ?