

# MADS-ML – Machine Learning

## Naive Bayes

Prof. Dr. Stephan Doerfel



**FACHHOCHSCHULE KIEL**  
University of Applied Sciences



Moodle (WiSe 2024/25)

---

“For events  $A$  and  $B$   
with  $P(B) > 0$ ,

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

holds.”

---

Bayes' Theorem

# Outline

**A (very brief) Intro to Probability Theory**

Basic Naive Bayes

Multinomial Bayes

Gaussian Bayes

Discussion

# Sample Spaces and Events

## Definition 1 (Sample Space, Event)

The **sample space**  $\Omega$  of an experiment is the set of all possible outcomes. The elements of  $\Omega$  are called **sample outcomes**. Subsets of  $\Omega$  are called **events**.

Example: Role a regular die.

- ▶  $\Omega = \{1, 2, 3, 4, 5, 6\}$
- ▶ Events: e.g.
  - ▶ “Role an even number”  $\{2, 4, 6\}$
  - ▶ “Role a 3”  $\{3\}$
  - ▶ “Role a 1 or a 4”  $\{1, 4\}$

# Probability Distribution

## Definition 2 (Probability Distribution)

Given a sample space  $\Omega$ , a function  $P : 2^\Omega \rightarrow \mathbb{R}$  is called a **probability distribution** if the following are true:

1. for  $A \subseteq \Omega : P(A) \geq 0$
2.  $P(\Omega) = 1$
3. If in a set of events  $\{A_i \mid i \in I\}$ , the  $A_i$  are pairwise disjoint, then

$$P\left(\bigcup_{i \in I} A_i\right) = \sum_{i \in I} P(A_i).$$

- ▶ Note that the index set  $I$  can be an infinite set!
- ▶ This definition is actually too narrow. A broader version leads to the mathematical theory of measures ...

# Probability Distribution – Example

Roll a (fair) die as in the example before:

- ▶  $P(\Omega) = 1$
- ▶  $P(\{2, 4, 6\}) = \frac{1}{2}$
- ▶  $P(\{3\}) = \frac{1}{6}$

Toss a (fair) coin three times:

- ▶  $\Omega = \{HHH, HHT, HTH, THH, HTT, THT, TTH, TTT\}$
- ▶  $P(3 \times H) = \frac{1}{8}$
- ▶  $P(2 \times H) = \frac{3}{8}$
- ▶  $P(\geq 2 \times H) = \frac{1}{2}$

# Independence

## Definition 3 (Independence)

Two events  $A, B \in \Omega$  are called **independent** if  $P(A \wedge B) = P(A) \cdot P(B)$ . A set of events  $\{A_i \mid i \in I\}$  is independent if

$$P\left(\bigwedge_{i \in I} A_i\right) = \prod_{i \in I} P(A_i)$$

Example: Throw a (fair) coin twice.

► Example 1:

- Event  $A$  = “first throw is heads”:  $P(A) = .5$
- Event  $B$  = “at least 1 heads”:  $P(B) = .75$
- $P(A \wedge B) = 0.5 \neq 0.5 \cdot 0.75$  (not independent)

► Example 2:

- Event  $A$  = “first throw is heads”:  $P(A) = .5$
- Event  $B$  = “second throw is heads”:  $P(B) = .5$
- $P(A \wedge B) = 0.25 = 0.5 \cdot 0.5$  (independent)

# Conditional Probability

## Definition 4 (Conditional Probability)

Given two events  $A, B \in \Omega$  with  $P(B) > 0$ . The **conditional probability** of  $A$  given  $B$  is

$$P(A \mid B) = \frac{P(A \cap B)}{P(B)}.$$

Example: Throw a (fair) coin three times.

- ▶  $A = \text{"exactly two heads"}: P(A) = \frac{3}{8}$
- ▶  $B = \text{"first throw is heads"}: P(B) = \frac{1}{2}$
- ▶  $P(A \mid B) = \frac{\frac{2}{8}}{\frac{1}{2}} = \frac{1}{2}$



# Independence and Conditional Probabilities

## Lemma 5

*If two events  $A, B \in \Omega$  are independent and  $P(B) > 0$ , then  $P(A \mid B) = P(A)$ .*

Example: Throw a (fair) coin three times.

- ▶  $A = \text{"second throw is heads"}: P(A) = \frac{1}{2}$
- ▶  $B = \text{"first throw is heads"}: P(B) = \frac{1}{2}$
- ▶  $P(A \mid B) = \frac{\frac{1}{4}}{\frac{1}{2}} = \frac{1}{2}$

# Bayes' Theorem

## Theorem 6 (Bayes' Theorem)

*Given two events  $A, B \in \Omega$  with  $P(A), P(B) > 0$ . Then*

$$P(A \mid B) = \frac{P(B \mid A) \cdot P(A)}{P(B)}.$$

# Random Variable

## Definition 7 (Random Variable)

A **random variable** is a function  $X : \Omega \rightarrow \mathbb{R}$ .

Example: Throw a (fair) coin three times. Count the number of heads.

# The Point Mass Distribution

- There exists one element  $a \in \Omega$  s.t. for  $\omega \in \Omega$  :

$$X(\omega) = \begin{cases} 1 & \text{if } \omega = a \\ 0 & \text{else.} \end{cases}$$

# The Discrete Uniform Distribution

- ▶ Let  $|X[\Omega]| = k$  ( $X$  takes  $k$  values) and  $P(X = x)$  is either 0 or  $\frac{1}{k}$ .
- ▶ Each value of  $X$  has the same probability.

# The Bernoulli Distribution

- ▶  $X$  represents the number of heads in a single (biased) coin flip.
- ▶  $P(X = 1) = p$  and  $P(X = 0) = 1 - p$ .
- ▶  $p$  is a (fix) parameter of the distribution.

# The Binomial Distribution

- ▶  $X$  represents the number of heads in  $n$  (independent) flips of a (biased) coin.



$$P(X = x) = \begin{cases} \binom{n}{x} p^x (1 - p)^{n-x} & \text{for } x \in 0, 1, \dots, n \\ 0 & \text{else.} \end{cases}$$

- ▶  $n$  and  $p$  are (fix) parameters of the distribution.
- ▶ generalization of Bernoulli distribution ( $n = 1$ )

# The Categorical Distribution

- ▶  $X$  represents the rolling of a  $k$  sided (biased) die.



$$P(X = x) = \begin{cases} p_i & \text{for } i \in \{1, \dots, k\} \\ 0 & \text{else.} \end{cases}$$

- ▶  $k$  and the  $p_i$  are (fix) parameters of the distribution with  $\sum_{i=1}^k p_i = 1$ .
- ▶ generalization of Bernoulli distribution ( $k = 2$ )



# Outline

A (very brief) Intro to Probability Theory

**Basic Naive Bayes**

Multinomial Bayes

Gaussian Bayes

Discussion

# Naive Bayes – Basics

Which class has the maximum likelihood,  
given the observed feature values?

- ▶ Given an instance (feature tuple)  $\mathbf{x}$ , which class  $c \in C$  has the highest likelihood to fit  $\mathbf{x}$ , i.e, for which class is it the most plausible that we would see  $\mathbf{x}$ ?
- ▶ Compare conditional probabilities:  $P(c \mid \mathbf{x})$
- ▶ Predict the class with the highest likelihood, i.e.

$$\underset{c \in C}{\operatorname{argmax}} P(c \mid \mathbf{x})$$

How can we estimate these conditional probabilities for arbitrary  $\mathbf{x}$ ?

# Estimating Likelihoods

- ▶ According to Bayes Theorem:

$$P(c \mid \mathbf{x}) = \frac{P(\mathbf{x} \mid c) \cdot P(c)}{P(\mathbf{x})}$$

- ▶ the likelihood of class  $c$  is computed from
  - ▶ the evidence  $P(\mathbf{x})$ ,
  - ▶ the prior probability of  $c$  and
  - ▶ the conditional probability of  $\mathbf{x}$  given  $c$ , which can be interpreted as the plausibility (likelihood) of  $c$  given  $\mathbf{x}$ .

# Naive Bayes – Classification

- ▶ Given an instance  $\mathbf{x}$ ,  $P(\mathbf{x})$  is constant. Thus

$$P(c \mid \mathbf{x}) = \frac{P(\mathbf{x} \mid c) \cdot P(c)}{P(\mathbf{x})} \propto P(\mathbf{x} \mid c) \cdot P(c).$$

- ▶ Thus the decision rule becomes

$$\underset{c \in C}{\operatorname{argmax}} P(\mathbf{x} \mid c) \cdot P(c).$$

- ▶  $P(c)$  can be estimated by counting the relative class frequency in the training data.

How to determine  $P(\mathbf{x} \mid c)$ ?

# Example 1: A Single Feature

day	wind	tennis?
1	weak	no
2	strong	no
3	weak	yes
4	weak	yes
5	weak	yes
6	strong	no

Classification task: Play tennis when wind is weak?

- ▶  $P(\text{no}|\text{weak}) \sim P(\text{weak}|\text{no}) \cdot P(\text{no}) = \frac{1}{3} \cdot \frac{3}{6} = \frac{1}{6}$
- ▶  $P(\text{yes}|\text{weak}) \sim P(\text{weak}|\text{yes}) \cdot \text{Prob}(\text{yes}) = \frac{3}{3} \cdot \frac{3}{6} = \frac{1}{2}$
- ▶ Decision: yes

Probability estimate: count frequencies in the training data.

## Example 2: Multiple Attributes

Attribute space is more sparsely covered (much bigger space through combinatorial explosion)

day	outlook	temperature	humidity	wind	tennis?
1	sunny	hot	high	weak	no
2	sunny	hot	high	strong	no
3	clouded	hot	high	weak	yes
4	rain	mild	high	weak	yes
5	rain	cool	normal	weak	yes
6	rain	cool	normal	strong	no

Classification task: Play tennis on a clouded, mild day with normal humidity and weak wind?

$$P(\text{clouded} \wedge \text{mild} \wedge \text{normal} \wedge \text{weak} \mid \text{no}) = ?$$

# Computing the Probabilities 1/2

- ▶ Task: Estimate  $P(\mathbf{x} \mid c)$ , for  $\mathbf{x} = (x_1, \dots, x_d)$
- ▶ **Naive Independence Assumption**: for each class  $c$ , the attribute values are independent, i.e.

$$P(\mathbf{x} \mid c) = \prod_{i=1}^d P(x_i \mid c)$$

- ▶ Estimate  $P(x_i \mid c)$  for  $i = 1, \dots, d$  by counting relative frequencies.
- ▶ Decision rule of Naive Bayes:

$$\underset{c \in C}{\operatorname{argmax}} \left( P(c) \cdot \prod_{i=1}^d P(x_i \mid c) \right)$$

---

“All models are wrong,  
but some are useful.”

---

George Box (statistician)



# Computing the Probabilities 2/2

## Naiveté:

- ▶ The naive independence assumption is usually never correct.
- ▶ Implicit assumption: All attributes are equally important (equally contribute to the product).

## Restriction

Works only for categorical attributes.

## Zero-Probabilities

- ▶ If for one  $x_i/c$ -combination  $P(x_i | c) = 0$ , then  $\prod_{i=1}^d P(x_i | c) = 0$ .
- ▶ Remedy: Smoothing. Introduce hyperparameter  $\alpha$ :

$$P(x_i | c; \alpha) = \frac{|\{t \in T \mid t_i = x_i, t \in c\}| + \alpha}{|\{t \in T \mid t \in c\}| + \alpha |\{t_i \mid t \in T\}|},$$

# Naive Bayes in Python

`sklearn.naive_bayes.CategoricalNB`

- ▶ Parameters: Mainly  $\alpha$  for smoothing the probabilities
- ▶ Variables must be categorical, yet presented as consecutive integers  $0, 1, \dots, n - 1$ .
- ▶ Variables can be transformed using the `OrdinalEncoder` as their order is ignored in the algorithm.

 Notebook 08\_1\_bayes\_tennis

# Outline

A (very brief) Intro to Probability Theory

Basic Naive Bayes

**Multinomial Bayes**

Gaussian Bayes

Discussion

# Restrictions of Categorical NB

- ▶ Assumes that in each class, each feature has a categorical distribution, i.e.  $k$  values each with their own probability, summing up to 1
- ▶ This makes it expensive if features have many possible categories
- ▶ This does not work with counts (all possible values for counts would have to be in the training data).
- ▶ This does not work with continuous values.

# Multivariate Distributions

- ▶ Consider multiple random variables  $X_i : \Omega \rightarrow \mathbb{R}$  for  $i = 1, \dots, n$ .
- ▶ These random variables are **independent** if for all choices of intervals  $A_i$ :

$$P(X_1 \in A_1, \dots, X_n \in A_n) = \prod_{i=1}^n P(X_i \in A_i).$$

# The Multinomial Distribution

- ▶ is a multivariate distribution
- ▶  $X_i$  the frequency of each side  $i$  of a  $k$  sided (biased) die in  $n$  (independent) roles.
- ▶  $P(X_1 = x_1, \dots, X_k = x_k)$

$$= \begin{cases} \frac{n!}{\prod_{i=1}^k x_i!} \prod_{i=1}^k p_i^{x_i} & \text{for } x_i \in \{0, \dots, n\} \text{ and } \sum_{i=1}^k x_i = n \\ 0 & \text{else.} \end{cases}$$

- ▶  $k$  and the  $p_i$  are (fix) parameters of the distribution with  $\sum_{i=1}^k p_i = 1$ .
- ▶ generalization of categorical distribution ( $n = 1$ )
- ▶ generalization of binomial distribution  $k = 2$

# Multinomial Naive Bayes

- ▶ Assumes that per class, the features are multinomial, i.e. each class has a multinomial distribution describing the (integer) values of all the features.
  - ▶ very popular in text mining, features are words with their frequency in a document as value
    - ▶ role a (biased) die with the vocabulary as sides
    - ▶ role as many times as there are words in a document
  - ▶ smoothing the probabilities is extremely important, as many document classes do not contain all words at least once
  - ▶ in python: `sklearn.naive_bayes.MultinomialNB`
    - ▶ parameter `alpha` for smoothing
    - ▶ probabilities are computed from sums (not frequencies) of the features
    - ▶ works with non-integer data as well (e.g. tf-idf)

 Notebook 08\_2\_bayes\_20\_news\_groups

# Outline

A (very brief) Intro to Probability Theory

Basic Naive Bayes

Multinomial Bayes

**Gaussian Bayes**

Discussion



# Continuous Distributions

We cannot always attribute probability to specific individual values.

## Definition 8

A random variable  $X$  is continuous if a function  $f_X$ , called **probability density function** exists, such that

1.  $f_X(x) \geq 0$  for  $x \in \mathbb{R}$
2.  $\int_{-\infty}^{\infty} f_X(x) dx = 1$
3. for  $a \leq b$  holds

$$P(a \leq X \leq b) = \int_a^b f_X(x) dx.$$

- ▶ Note: we can only determine the probability for intervals of real number here!
- ▶ The distribution is given indirectly through the density function.
- ▶  $f(X)$  is not the same as  $P(X = x)$ . In fact,  $P(X = x) = 0$ .

# The Continuous Uniform Distribution

- ▶ probability density function:

$$f(x) = \begin{cases} \frac{1}{b-a} & \text{for } x \in [a, b] \\ 0 & \text{else} \end{cases}$$

- ▶  $a, b$  are fix parameters of the distribution
- ▶ the size of the area under the density function, in the interval  $[a, b]$  is 1
- ▶ the probability of each individual value is 0

# The Gaussian Distribution

- ▶ probability density function:

$$f(x) = \frac{1}{\sigma\sqrt{2\pi}} e^{-\frac{1}{2}\left(\frac{x-\mu}{\sigma}\right)^2}$$

- ▶  $\mu$  and  $\sigma$  are (fix) parameters of the distribution, representing mean and standard deviation of data sampled from the distribution
- ▶ Gaussians are central elements of various statistical analyses

# Gaussian Naive Bayes

- ▶ Allows continuous data (counting frequencies or summing them does not work with continuous data).
- ▶ The likelihood of  $c$  being the class of an observed instance  $\mathbf{x}$  with feature  $x_i$  is estimated as

$$P(x_i | c) = \frac{1}{\sqrt{2\pi\sigma_c^2}} \exp\left(-\frac{(x_i - \mu_c)^2}{2\sigma_c^2}\right)$$

- ▶ one Gaussian per class and feature.

# Outline

A (very brief) Intro to Probability Theory

Basic Naive Bayes

Multinomial Bayes

Gaussian Bayes

**Discussion**

# Discussion: Bayes-Classifications

## Positive:

- ▶ Optimality property: On average, no other classification algorithm with the same a-priori knowledge can reach a higher average classification quality.
- ▶ always a good baseline for evaluating other classifiers
- ▶ high accuracy on many classification problems
- ▶ incremental learning: new training instances can easily be added
- ▶ includes knowledge about the domain (data distributions)

## Negative:

- ▶ naive assumption sometimes too naive
- ▶ inefficiency with many attributes