

Business Analytics & Machine Learning

Tutorial sheet 4: Naïve Bayes – Solution

Prof. Dr. Martin Bichler, Prof. Dr. Jalal Etesami
Julius Durmann, Markus Ewert, Johannes Knörr, Yutong Chao

Exercise T4.1 *States of a car*

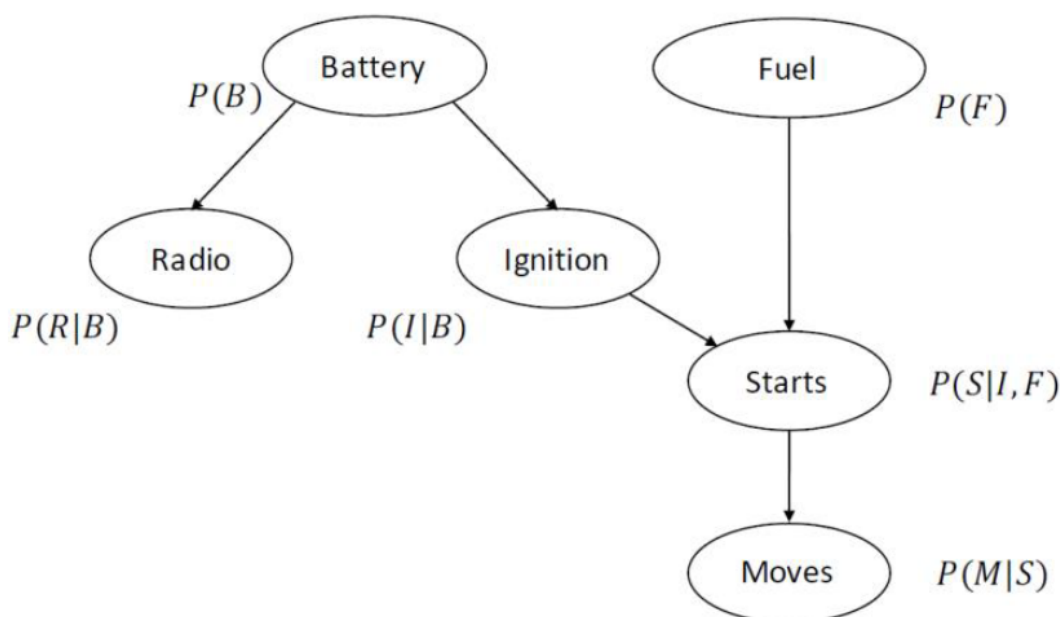
Consider the following Boolean random variables related to the state of a given car.

- Battery (B): is the battery charged?
- Fuel (F): is the fuel tank filled?
- Ignition (I): does the ignition system work?
- Moves (M): does the car move?
- Radio (R): is the radio working?
- Starts (S): does the engine work?

- Represent the joint probability density function using a Bayesian Network.
- Express the joint probability $\mathbb{P}(B, F, I, M, R, S)$ as a product of conditional probabilities using the chain rule.

Solution

- In a usual modern car, the (independent) energy sources are the Battery (B) and the Fuel Tank (F). Electronics such as the Radio (R) and the Ignition System (I) depend on the state of the Battery (B). Whether the car Starts (S) comes down to Fuel (F) being present while Ignition (I) is activated. Finally, whether the car Moves (M) is depending on whether it previously Starts (S). Combining these facts, the corresponding Bayesian Network is:



- b) Topologically ordering the variables from the Bayesian Network from a) as B, F, R, I, S, M, the chain rule yields the following factorization:

$$\mathbb{P}(B, F, R, I, S, M) = \mathbb{P}(B) \cdot \mathbb{P}(F|B) \cdot \mathbb{P}(R|B, F) \cdot \mathbb{P}(I|B, F, R) \cdot \mathbb{P}(S|B, F, R, I) \cdot \mathbb{P}(M|B, F, R, I, S).$$

Using conditional independences, the term can be simplified to:

$$\mathbb{P}(B, F, R, I, S, M) = \mathbb{P}(B) \cdot \mathbb{P}(F) \cdot \mathbb{P}(R|B) \cdot \mathbb{P}(I|B) \cdot \mathbb{P}(S|F, I) \cdot \mathbb{P}(M|S).$$

Exercise T4.2 *Cookie factory*

You are the manager of a company which produces cookies and you want to introduce a new product. Your R&D department has proposed and developed the following two alternatives:

1. Unicorn cookies (UC)
2. Vanilla-chip cookies (VC).

As part of your market research, you are interested in predicting whether certain customers are likely to buy one of the new products. For that, you have already collected data from a large number of test persons. In particular, you asked them to fill out a query with the following questions:

1. How old are you? (variable *age*)
2. What do you think is the most fascinating: Rainbows, Black holes or Cats? (variable *preferences*)
3. How much money do you spend on cookies per month? (variable *money*)
4. Which of our cookies would you buy? (variable *product*)

Note: The variable *product* can also take on the value "No product" (NP).

You can find the data in *cookie-factory.csv*. We recommend using the provided notebook template to solve sub-tasks b) - e).

- a) For each of the questions 1-4, decide
 - (i) whether the answers are continuous or discrete outcomes,
 - (ii) which range the outcomes could have,
 - (iii) to which scale of measurement (nominal, ordinal, interval, ratio) the outcomes belong to.
- b) To infer which products new customers are likely to buy, you set up a probabilistic model. You assume that the answers to questions 1 - 3 are conditionally independent (Naive Bayes) given *product* and model the dependencies as follows:

$$f(\text{age}, \text{preferences}, \text{money}, \text{product}) = \mathbb{P}(\text{age} | \text{product}) \cdot \mathbb{P}(\text{preferences} | \text{product}) \cdot f_{\text{money}}(\text{money} | \text{product}) \cdot \mathbb{P}(\text{product})$$

Estimate the parameters of your categorical prior $\mathbb{P}(\text{product})$ by using maximum likelihood:

$$\mathbb{P}(\text{product} = UC) = p_{UC} \quad \mathbb{P}(\text{product} = VC) = p_{VC} \quad \mathbb{P}(\text{product} = NP) = p_{NP}$$

Hint: The maximum likelihood estimate of the parameters for categorically distributed variables is simply the fraction of samples from a category.

c) Based on your observations in a), you decide to model the likelihoods as follows:

1. *age* follows a Poisson distribution where the parameter $\lambda_{product}$ depends on the product the customers would buy ($\lambda_{product} = \lambda_{UC}$, $\lambda_{product} = \lambda_{VC}$, or $\lambda_{product} = \lambda_{NP}$):

$$\mathbb{P}(age = k | product) = \frac{\lambda_{product}^k}{k!} e^{-\lambda_{product}}$$

2. *preferences* follows a Categorical distribution where the parameters depend on the product the customers would buy:

$$\mathbb{P}(preferences = \text{"Rainbows"} \mid prod. = UC) = \pi_{UC}^R$$

$$\mathbb{P}(preferences = \text{"Black holes"} \mid prod. = UC) = \pi_{UC}^B$$

$$\mathbb{P}(preferences = \text{"Cats"} \mid prod. = UC) = \pi_{UC}^C$$

$$\mathbb{P}(preferences = \text{"Rainbows"} \mid prod. = VC) = \pi_{VC}^R$$

$$\mathbb{P}(preferences = \text{"Black holes"} \mid prod. = VC) = \pi_{VC}^B$$

$$\mathbb{P}(preferences = \text{"Cats"} \mid prod. = VC) = \pi_{VC}^C$$

$$\mathbb{P}(preferences = \text{"Rainbows"} \mid prod. = NP) = \pi_{NP}^R$$

$$\mathbb{P}(preferences = \text{"Black holes"} \mid prod. = NP) = \pi_{NP}^B$$

$$\mathbb{P}(preferences = \text{"Cats"} \mid prod. = NP) = \pi_{NP}^C$$

3. *money* follows an exponential distribution where the parameter $\eta_{product}$ depends on the product the customers would buy ($\eta_{product} = \eta_{UC}$, $\eta_{product} = \eta_{VC}$ or $\eta_{product} = \eta_{NP}$):

$$f_{money}(m | product) = \begin{cases} \eta_{product} \cdot e^{-\eta_{product} \cdot m} & m \geq 0 \\ 0 & \text{else} \end{cases}$$

Intuitively, your model describes the profile (*age*, *preferences*, *money*) of a customer if you already know which product they would buy (*product*).

Using the data, derive maximum likelihood estimates for all parameters.

Hint: The maximum likelihood estimate of the parameters for Poisson distributed variables is simply the sample mean: \bar{x} .

Hint: The maximum likelihood estimate of the parameters for exponentially distributed variables is the inverse of their sample mean: \bar{x}^{-1} .

d) You now have access to a joint density over your data:

$$f(age, preferences, money, product) = \mathbb{P}(age \mid product) \cdot \mathbb{P}(preferences \mid product) \cdot f_{money}(money \mid product) \cdot \mathbb{P}(product)$$

With the fitted model, predict the (posterior) probability

$$\mathbb{P}(product \mid age, preferences, money)$$

that the customers below buy a unicorn cookie, a vanilla-chip cookie or no cookie at all:

Customer	<i>preferences</i>	<i>age</i>	<i>money</i>
Anna	Cats	81	53.10 €
Ben	Rainbows	15	2.30 €
Caroline	Black holes	42	10.25 €

- e) From a fourth customer, you only know that they like rainbows. Predict the probability that they buy unicorn cookies.

Solution

The solution below is meant as a reference. Please also have a look at the provided solution notebook which should give more insight into the computation.

- a) The table below shows our suggestion. Some of the answers could surely be discussed. For instance, the variable *money* could be considered discrete since cookie prices can only differ up to a unit of 1 cent. Also, it is unlikely that a person becomes much older than 150 years.

Variable	Discrete?	Range	Scale of measurement
<i>age</i>	discrete	$[0, \infty)$	Ratio
<i>money</i>	continuous	$[0, \infty)$	Ratio
<i>preferences</i>	discrete	$\{ \text{"rainbows"}, \text{"black holes"}, \text{"cats"} \}$	Nominal
<i>product</i>	discrete	$\{ \text{"UC"}, \text{"VC"}, \text{"NP"} \}$	Nominal

- b) The optimal parameters are as follows:

p_{UC}	0.622
p_{VC}	0.208
p_{NP}	0.170

- c) The optimal parameters are as follows:

1.	λ_{UC}	14.91
	λ_{VC}	56.08
	λ_{NP}	29.80

2.	π_{UC}^R	0.579	π_{UC}^B	0.305	π_{UC}^C	0.116
	π_{VC}^R	0.356	π_{VC}^B	0.245	π_{VC}^C	0.399
	π_{NP}^R	0.441	π_{NP}^B	0.418	π_{NP}^C	0.141

3.	η_{UC}	0.092
	η_{VC}	0.066
	η_{NP}	0.109

d) With the joint distribution we can derive the posterior:

$$\begin{aligned}
& \mathbb{P}(\text{product} \mid \text{age}, \text{money}, \text{preferences}) \\
&= \frac{f(\text{age}, \text{preferences}, \text{money}, \text{product})}{f(\text{age}, \text{preferences}, \text{money})} \\
&= \frac{1}{Z} \cdot f(\text{age}, \text{preferences}, \text{money}, \text{product}) \\
&= \frac{1}{Z} \cdot \mathbb{P}(\text{age} \mid \text{product}) \cdot \mathbb{P}(\text{preferences} \mid \text{product}) \cdot f_{\text{money}}(\text{money} \mid \text{product}) \cdot \mathbb{P}(\text{product})
\end{aligned}$$

By plugging in the values from b) and c), we can compute the individual terms. Once this is achieved for all products, the normalizing constant Z is such that the posterior probabilities over all products (UC, VC, NP) sum up to 1.

For the customer Anna, we need to compute:

$$\begin{aligned}
& \tilde{\mathbb{P}}(\text{product} = UC \mid \text{age} = 81, \text{money} = 53.10\text{€}, \text{preferences} = \text{"Cats"}) \\
&= \frac{1}{Z} \cdot \mathbb{P}(\text{age} = 81 \mid UC) \cdot \mathbb{P}(\text{preferences} = \text{"Cats"} \mid UC) \cdot f_{\text{money}}(53.10 \mid UC) \cdot \mathbb{P}(\text{product} = UC) \\
& \tilde{\mathbb{P}}(\text{product} = VC \mid \text{age} = 81, \text{money} = 53.10\text{€}, \text{preferences} = \text{"Cats"}) \\
&= \frac{1}{Z} \cdot \mathbb{P}(\text{age} = 81 \mid VC) \cdot \mathbb{P}(\text{preferences} = \text{"Cats"} \mid VC) \cdot f_{\text{money}}(53.10 \mid VC) \cdot \mathbb{P}(\text{product} = VC) \\
& \tilde{\mathbb{P}}(\text{product} = NP \mid \text{age} = 81, \text{money} = 53.10\text{€}, \text{preferences} = \text{"Cats"}) \\
&= \frac{1}{Z} \cdot \mathbb{P}(\text{age} = 81 \mid NP) \cdot \mathbb{P}(\text{preferences} = \text{"Cats"} \mid NP) \cdot f_{\text{money}}(53.10 \mid NP) \cdot \mathbb{P}(\text{product} = NP)
\end{aligned}$$

We find the normalizing constant Z as

$$\begin{aligned}
Z = & \tilde{\mathbb{P}}(\text{product} = UC \mid \text{age} = 81, \text{money} = 53.10\text{€}, \text{preferences} = \text{"Cats"}) \\
& + \tilde{\mathbb{P}}(\text{product} = VC \mid \text{age} = 81, \text{money} = 53.10\text{€}, \text{preferences} = \text{"Cats"}) \\
& + \tilde{\mathbb{P}}(\text{product} = NP \mid \text{age} = 81, \text{money} = 53.10\text{€}, \text{preferences} = \text{"Cats"})
\end{aligned}$$

Finally, we can use Z to normalize, and obtain:

$$\begin{aligned}
& \mathbb{P}(\text{product} = UC \mid \text{age} = 81, \text{money} = 53.10\text{€}, \text{preferences} = \text{"Cats"}) \\
&= \frac{1}{Z} \cdot \tilde{\mathbb{P}}(\text{product} = UC \mid \text{age} = 81, \text{money} = 53.10\text{€}, \text{preferences} = \text{"Cats"}) \\
& \mathbb{P}(\text{product} = VC \mid \text{age} = 81, \text{money} = 53.10\text{€}, \text{preferences} = \text{"Cats"}) \\
&= \frac{1}{Z} \cdot \tilde{\mathbb{P}}(\text{product} = VC \mid \text{age} = 81, \text{money} = 53.10\text{€}, \text{preferences} = \text{"Cats"}) \\
& \mathbb{P}(\text{product} = NP \mid \text{age} = 81, \text{money} = 53.10\text{€}, \text{preferences} = \text{"Cats"}) \\
&= \frac{1}{Z} \cdot \tilde{\mathbb{P}}(\text{product} = NP \mid \text{age} = 81, \text{money} = 53.10\text{€}, \text{preferences} = \text{"Cats"})
\end{aligned}$$

For the other customer profiles, the formulas are similar.

The computation could be done by hand. Since this is tedious, we rely on Python and suggest you have a look at the provided solution notebook.

The resulting (rounded) probabilities are as follows:

	No product	Unicorn cookie	Vanilla cookie
Anna	0.0	0.0	1.0
Ben	0.003	0.997	0.0
Caroline	0.530	0.0	0.470

e) If we have only limited information, we must consider the marginal joint distribution

$$\mathbb{P}(\text{preferences}, \text{product}) = \mathbb{P}(\text{preferences} \mid \text{product}) \cdot \mathbb{P}(\text{product})$$

The posterior becomes

$$\mathbb{P}(\text{product} \mid \text{preferences}) = \frac{1}{Z} \cdot \mathbb{P}(\text{preferences} \mid \text{product}) \cdot \mathbb{P}(\text{product})$$

Plugging in the maximum-likelihood values and normalizing gives

	No product	Unicorn cookie	Vanilla cookie
Unknown customer	0.147	0.707	0.145