

Sprachverarbeitung: Übung

SoSe 24

Janis Pagel

Department for Digital Humanities, University of Cologne

2024-06-11

For this exercise, you need to both submit manual calculations as well as Python code. Please submit two files in Ilias, one a PDF with your calculations and one a file containing your Python code (either Jupyter Notebook or Python script). You can also combine both files into a zip-archive and submit only the archive. You can either solve the calculations by hand on a sheet of paper, scan it and submit as a PDF file or use the capabilities to write mathematical equations of tools like MS Word / LibreOffice / LaTeX, etc. to write down your calculations digitally.

Exercise 1.

For this part of the exercise, you need to do manual calculations and submit your solution in a PDF file.

Given are the following nine sentences in different languages, Dutch (nl), English (en) and French (fr):

| Sentence | Class |
|------------------------------|-------|
| een man rookt | nl |
| er loopt een dier | nl |
| de man eet pizza | nl |
| the man snorted | en |
| never received the product | en |
| never received my package | en |
| apprendre de manière ludique | fr |
| ça ne marche absolument pas | fr |
| pas de version en français | fr |

Table 1: Language identification dataset

Given this training data, create a Naïve Bayes model by hand by calculating all probabilities for a certain feature to occur/not occur in a sentence given a certain class (language). Apply “Add-One” smoothing to all probabilities. Only use the following seven features: “een”, “man”, “de”, “never”, “pas”, “version” and “en”.

Using the calculated probabilities, determine the class of the following three sentences by calculating the probability of a certain class given the presence or absence of the features in the sentences:

| Sentence |
|-----------------------------|
| een man en een vrouw praten |
| never ordered this version |
| je n'ai pas les mots |

Solution 1.

Tables of feature occurrences in the data.

Feature value=1:

| | nl | en | fr |
|----------|----|----|----|
| een | 2 | 0 | 0 |
| man | 2 | 1 | 0 |
| de | 1 | 0 | 2 |
| never | 0 | 2 | 0 |
| pas | 0 | 0 | 2 |
| version | 0 | 0 | 1 |
| en | 0 | 0 | 1 |
| Σ | 5 | 3 | 6 |

Feature value=0:

| | nl | en | fr |
|----------|----|----|----|
| een | 1 | 3 | 2 |
| man | 1 | 2 | 3 |
| de | 2 | 3 | 1 |
| never | 3 | 1 | 3 |
| pas | 3 | 3 | 1 |
| version | 3 | 3 | 2 |
| en | 3 | 3 | 2 |
| Σ | 16 | 18 | 14 |

Class probabilities

$$p(nl) = p(en) = p(fr) = \frac{3}{9} = \frac{1}{3}$$

Probability for a feature given class=nl

$$\begin{aligned} p(een = 1|nl) &= \frac{2+1}{5+1} &&= 0.5 \\ p(man = 1|nl) &= \frac{2+1}{5+1} &&= 0.5 \\ p(de = 1|nl) &= \frac{1+1}{5+1} &&= \frac{1}{3} \\ p(never = 1|nl) &= \frac{0+1}{5+1} = \frac{1}{6} &&\approx 0.16 \\ p(pas = 1|nl) &= \frac{0+1}{5+1} &&= \frac{1}{6} \\ p(version = 1|nl) &= \frac{0+1}{5+1} &&= \frac{1}{6} \\ p(en = 1|nl) &= \frac{0+1}{5+1} &&= \frac{1}{6} \\ p(een = 0|nl) &= \frac{1+1}{16+1} = \frac{2}{17} &&\approx 0.12 \\ p(man = 0|nl) &= \frac{1+1}{16+1} &&= \frac{2}{17} \\ p(de = 0|nl) &= \frac{2+1}{16+1} = \frac{3}{17} &&\approx 0.18 \\ p(never = 0|nl) &= \frac{3+1}{16+1} = \frac{4}{17} &&\approx 0.24 \\ p(pas = 0|nl) &= \frac{3+1}{16+1} &&= \frac{4}{17} \\ p(version = 0|nl) &= \frac{3+1}{16+1} &&= \frac{4}{17} \\ p(en = 0|nl) &= \frac{3+1}{16+1} &&= \frac{4}{17} \end{aligned}$$

Probability for a feature given class=en

$$\begin{aligned}p(een = 1|en) &= \frac{0 + 1}{3 + 1} = 0.25 \\p(man = 1|en) &= \frac{1 + 1}{3 + 1} = 0.5 \\p(de = 1|en) &= \frac{0 + 1}{3 + 1} = 0.25 \\p(never = 1|en) &= \frac{2 + 1}{3 + 1} = 0.75 \\p(pas = 1|en) &= \frac{0 + 1}{3 + 1} = 0.25 \\p(version = 1|en) &= \frac{0 + 1}{3 + 1} = 0.25 \\p(en = 1|en) &= \frac{0 + 1}{3 + 1} = 0.25 \\p(een = 0|en) &= \frac{3 + 1}{18 + 1} = \frac{4}{19} \approx 0.21 \\p(man = 0|en) &= \frac{2 + 1}{18 + 1} = \frac{3}{19} \approx 0.16 \\p(de = 0|en) &= \frac{3 + 1}{18 + 1} = \frac{4}{19} \\p(never = 0|en) &= \frac{1 + 1}{18 + 1} = \frac{2}{19} \approx 0.11 \\p(pas = 0|en) &= \frac{3 + 1}{18 + 1} = \frac{4}{19} \\p(version = 0|en) &= \frac{3 + 1}{18 + 1} = \frac{4}{19} \\p(en = 0|en) &= \frac{3 + 1}{18 + 1} = \frac{4}{19}\end{aligned}$$

Probability for a feature given class=fr

$$\begin{aligned}
 p(een = 1|fr) &= \frac{0+1}{6+1} = \frac{1}{7} \approx 0.14 \\
 p(man = 1|fr) &= \frac{0+1}{6+1} = \frac{1}{7} \\
 p(de = 1|fr) &= \frac{2+1}{6+1} = \frac{3}{7} \approx 0.43 \\
 p(never = 1|fr) &= \frac{0+1}{6+1} = \frac{1}{7} \\
 p(pas = 1|fr) &= \frac{2+1}{6+1} = \frac{3}{7} \\
 p(version = 1|fr) &= \frac{1+1}{6+1} = \frac{2}{7} \approx 0.29 \\
 p(en = 1|fr) &= \frac{1+1}{6+1} = \frac{2}{7} \\
 p(een = 0|fr) &= \frac{2+1}{14+1} = 0.2 \\
 p(man = 0|fr) &= \frac{3+1}{14+1} = \frac{4}{15} \approx 0.26 \\
 p(de = 0|fr) &= \frac{1+1}{14+1} = \frac{2}{15} \approx 0.13 \\
 p(never = 0|fr) &= \frac{3+1}{14+1} = \frac{4}{15} \\
 p(pas = 0|fr) &= \frac{1+1}{14+1} = \frac{2}{15} \\
 p(version = 0|fr) &= \frac{2+1}{14+1} = 0.2 \\
 p(en = 0|fr) &= \frac{2+1}{14+1} = 0.2
 \end{aligned}$$

Use these probabilities to calculate the most likely class the test sentences, given their distribution of features: For the sentence “een man en een vrouw praten”

$$\begin{aligned}
 & p(nl|een = 1, man = 1, de = 0, never = 0, pas = 0, version = 0, en = 1) \\
 \propto & p(een = 1|nl) \times p(man = 1|nl) \times p(de = 0|nl) \times p(never = 0|nl) \times p(pas = 0|nl) \times p(version = 0|nl) \times p(en = 1|nl) \times p(nl) \\
 & = 0.5 \times 0.5 \times \frac{3}{17} \times \frac{4}{17} \times \frac{4}{17} \times \frac{4}{17} \times \frac{1}{6} \times \frac{1}{3} \approx 0.000032 \\
 & p(en|een = 1, man = 1, de = 0, never = 0, pas = 0, version = 0, en = 1) \propto 0.25 * 0.5 * \frac{4}{19} * \frac{2}{19} * \frac{4}{19} * \frac{4}{19} * 0.25 * \frac{1}{3} \approx 0.00001 \\
 & p(fr|een = 1, man = 1, de = 0, never = 0, pas = 0, version = 0, en = 1) \propto \frac{1}{7} * \frac{1}{7} * \frac{2}{15} * \frac{4}{15} * \frac{2}{15} * 0.2 * \frac{2}{7} * \frac{1}{3} \approx 0.0000018 \\
 & \qquad \qquad \qquad 0.000032 > 0.00001 > 0.0000018
 \end{aligned}$$

This means that “nl” is the most likely class for this sentence.

For the sentence “never ordered this version”

$$\begin{aligned}p(nl|een = 0, man = 0, de = 0, never = 1, pas = 0, version = 1, en = 0) &\propto \frac{2}{17} \times \frac{2}{17} \times \frac{3}{17} \times \frac{1}{6} \times \frac{4}{17} \times \frac{1}{6} \times \frac{4}{17} \times \frac{1}{3} \approx 0.0000013 \\p(en|een = 0, man = 0, de = 0, never = 1, pas = 0, version = 1, en = 0) &\propto \frac{4}{19} * \frac{3}{19} * \frac{4}{19} * 0.75 * \frac{4}{19} * 0.25 * \frac{4}{19} * \frac{1}{3} \approx 0.000019 \\p(fr|een = 0, man = 0, de = 0, never = 1, pas = 0, version = 1, en = 0) &\propto 0.2 * \frac{4}{15} * \frac{2}{15} * \frac{1}{7} * \frac{2}{15} * \frac{2}{7} * 0.2 * \frac{1}{3} \approx 0.0000026 \\&0.000019 > 0.0000026 > 0.0000013\end{aligned}$$

“en” is the most likely class for this sentence.

For the sentence “je n’ai pas les mots”

$$\begin{aligned}p(nl|een = 0, man = 0, de = 0, never = 0, pas = 1, version = 0, en = 0) &\propto \frac{2}{17} \times \frac{2}{17} \times \frac{3}{17} \times \frac{4}{17} \times \frac{1}{6} \times \frac{4}{17} \times \frac{4}{17} \times \frac{1}{3} \approx 0.0000018 \\p(en|een = 0, man = 0, de = 0, never = 0, pas = 1, version = 0, en = 0) &\propto \frac{4}{19} * \frac{3}{19} * \frac{4}{19} * \frac{2}{19} * 0.25 * \frac{4}{19} * \frac{4}{19} * \frac{1}{3} \approx 0.0000027 \\p(fr|een = 0, man = 0, de = 0, never = 0, pas = 1, version = 0, en = 0) &\propto 0.2 * \frac{4}{15} * \frac{2}{15} * \frac{4}{15} * \frac{3}{7} * 0.2 * 0.2 * \frac{1}{3} \approx 0.000011 \\&0.000011 > 0.0000027 > 0.0000018\end{aligned}$$

“fr” is the most likely class for this sentence.

Exercise 2.

For this part of the exercise, you need to write Python code.

The dataset in table 1 is part of a larger dataset from <https://huggingface.co/datasets/papluca/language-identification>. From the course website, download the files https://lehre.idh.uni-koeln.de/site/assets/files/5151/languageidentification_train.tsv and https://lehre.idh.uni-koeln.de/site/assets/files/5151/languageidentification_test.tsv, which contain a train and test split for this dataset, with each column representing a token feature. The column containing the language classes is called “labels”. Train a sklearn Bernoulli Naïve Bayes model on the train set and let it predict on the test set. Use sklearn’s `classification_report` function to obtain several evaluation metrics for your prediction. Also create a confusion matrix of the test classes and the classes predicted by the model using sklearn’s `confusion_matrix` function. Afterwards, plot the confusion matrix using seaborn’s heatmap function.

Solution 2.

```
from sklearn.naive_bayes import BernoulliNB
from sklearn.metrics import classification_report, confusion_matrix
import pandas as pd
import seaborn as sns
import matplotlib.pyplot as plt

train_df = pd.read_csv("language_identification_train.csv")
X_train, y_train = train_df["text"], train_df["labels"]
test_df = pd.read_csv("language_identification_test.csv")
X_test, y_test = test_df["text"], test_df["labels"]
```

```
clf = BernoulliNB()
clf.fit(X_train, y_train)
y_pred = clf.predict(X_test)

print(classification_report(y_test, y_pred))

conf_matrix = confusion_matrix(y_test, y_pred, labels=clf.classes_)
conf_matrix = pd.DataFrame(conf_matrix, index=clf.classes_, columns=clf.
                           classes_)

plt.subplots(figsize=(12,12))
sns.heatmap(conf_matrix, annot=True, fmt="g")
```