

# Assignment for Week 3 - Naive-Bayes

## Task 1 - Bayesian Classification

**Important Note:** This exercise is not a programming exercise, it is a math exercise to help reinforce the math behind a Bayesian Classification. You can fill free to complete this in any method you feel is appropriate (ie: pencil/paper (you will need to scan your work to submit), Excel workbook, markdown text with a jupyter notebook, etc)

Please show all your work.

1. In a study of pleas and prison sentences, it is reported that 42% of the subjects were sent to prison. Among those sent to prison, 38% pleaded guilty. Among those not sent to prison, 50% pleaded guilty.
- a) If a subject is randomly selected, what is the probability of getting a person who was not sent to prison?
- b) If a subject is randomly selected, and it is known that the subject entered a guilty plea, what is the probability that this subject was not sent to prison?
- c) If a subject is randomly selected, what is the probability of getting someone who was sent to prison?
- d) If a subject is randomly selected, and it is known that the subject entered a guilty plea, what is the probability that this person was sent to prison?

## Bayes Theorem:

The probability of event A, given that event B has subsequently occurred can be defined as

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

**Note:**  $P(\neg A)$  denotes  $P(\text{not } A)$  and the same applies for  $P(\neg B)$ .

## What we know at this point

Let A be the event of **sent to prison**.

$P(A) = 42\%$   
Let B be the event of **entering a guilty plea**.

A) 58% of the subjects was not sent to prison. B) 50% of people who pleaded guilty was not sent to prison. C) 38% of people who pleaded guilty was sent to prison.

To make it easier to calculate, I going to assume that the number of subjects is 100,000.

If the percentage of subjects sent to prison is 42%, then  $0.42 \cdot 100,000 = 42,000$  and if the percentage of subjects who sent to prison and pleaded guilty is 38%, then  $0.38 \cdot 42,000 = 15,960$ , which is the true positive. 15,960% of the subjects who pleaded guilty was actually sent to prison.

The number of subject who pleaded not guilty and sent to prison is  $42,000 - 15,960 = 26,040 = 26.04\%$  of subjects sent to prison not guilty

subject was not sent to prison =  $58\% \cdot 100,000 = 58,000$ .

50% of subject who didn't plea guilty and not sent to prison =  $50\% \cdot 0.5 \cdot 58,000 = 33,640 = 33\%$  of all subjects (True negative) the number of subjects who pleaded guilty and sent to jail =  $58,000 - 33,640 = 24,360 = 24.36\%$  of all subjects (False negative)

number of subjects who pleaded guilty is  $24,360 + 15,960 = 40,320 = 40.32\%$  of all subjects. number of subject who pleaded not guilty is  $100,000 - 40,320 = 59,680 = 59.68\%$  of all subjects

A:  $1 - 0.42 = 0.58 \cdot 100 = 58\%$  of the subject was not sent to jail. B:  $24,360 / 40,320 = 60.4\%$  probability someone who entered a guilty plea was not sent to prison. C:  $42\%$  probability of selecting someone who was sent to prison. D:  $24.36\%$  probability of someone who entered a guilty plea was sent to prison.

*Continue defining the probabilities and compute b-d*

1. Given the following table:

Customer ID	Gender	Car Type	Shirt Size	Class
1	M	Family	Small	C0
2	M	Sports	Medium	C0
3	M	Sports	Medium	C0
4	M	Sports	Large	C0
5	M	Sports	Extra Large	C0
6	M	Family	Extra Large	C0
7	F	Sports	Small	C0
8	F	Sports	Small	C0
9	F	Sports	Medium	C0
10	F	Luxury	Large	C0
11	M	Family	Large	C1
12	M	Family	Extra Large	C1
13	M	Family	Medium	C1
14	M	Family	Extra Large	C1
15	F	Luxury	Small	C1
16	F	Luxury	Small	C1
17	F	Luxury	Medium	C1
18	F	Luxury	Medium	C1
19	F	Luxury	Medium	C1
20	F	Luxury	Large	C1

a) What is the value of each of the following probabilities?

- $P(\text{Gender} = \text{M} \mid \text{Class} = \text{C0})$
- $P(\text{Gender} = \text{F} \mid \text{Class} = \text{C1})$
- $P(\text{Car Type} = \text{Family} \mid \text{Class} = \text{C0})$
- $P(\text{Car Type} = \text{Family} \mid \text{Class} = \text{C1})$
- $P(\text{Shirt Size} = \text{Medium} \mid \text{Class} = \text{C0})$
- $P(\text{Shirt Size} = \text{Medium} \mid \text{Class} = \text{C1})$

50% of subjects are M 50% of subjects are C0 50% of subjects are F 50% of subjects are C1

A)  $P(\text{Gender} = \text{M} \mid \text{Class} = \text{C0}) = \frac{P(\text{Gender} = \text{M} \cap \text{Class} = \text{C0})}{P(\text{Class} = \text{C0})} = \frac{6/10}{0.6} = \frac{P(\text{Gender} = \text{M}) \cdot P(\text{Class} = \text{C0})}{P(\text{Class} = \text{C0})} = \frac{6/10}{0.6} = 0.5$   
 $P(\text{Gender} = \text{M} \mid \text{Class} = \text{C0}) = 0.5$

B)  $P(\text{Gender} = \text{F} \mid \text{Class} = \text{C1}) = \frac{P(\text{Gender} = \text{F} \cap \text{Class} = \text{C1})}{P(\text{Class} = \text{C1})} = \frac{6/10}{0.6} = \frac{P(\text{Gender} = \text{F}) \cdot P(\text{Class} = \text{C1})}{P(\text{Class} = \text{C1})} = \frac{6/10}{0.6} = 0.5$   
 $P(\text{Gender} = \text{F} \mid \text{Class} = \text{C1}) = 0.5$

C)  $P(\text{Car Type} = \text{Family} \mid \text{Class} = \text{C0}) = \frac{P(\text{Car Type} = \text{Family} \cap \text{Class} = \text{C0})}{P(\text{Class} = \text{C0})} = \frac{4/10}{0.6} = \frac{P(\text{Car Type} = \text{Family}) \cdot P(\text{Class} = \text{C0})}{P(\text{Class} = \text{C0})} = \frac{4/10}{0.6} = 0.6667$   
 $P(\text{Car Type} = \text{Family} \mid \text{Class} = \text{C0}) = 0.6667$

D)  $P(\text{Car Type} = \text{Family} \mid \text{Class} = \text{C1}) = \frac{P(\text{Car Type} = \text{Family} \cap \text{Class} = \text{C1})}{P(\text{Class} = \text{C1})} = \frac{4/10}{0.6} = \frac{P(\text{Car Type} = \text{Family}) \cdot P(\text{Class} = \text{C1})}{P(\text{Class} = \text{C1})} = \frac{4/10}{0.6} = 0.6667$   
 $P(\text{Car Type} = \text{Family} \mid \text{Class} = \text{C1}) = 0.6667$

E)  $P(\text{Shirt Size} = \text{Medium} \mid \text{Class} = \text{C0}) = \frac{P(\text{Shirt Size} = \text{Medium} \cap \text{Class} = \text{C0})}{P(\text{Class} = \text{C0})} = \frac{3/10}{0.6} = \frac{P(\text{Shirt Size} = \text{Medium}) \cdot P(\text{Class} = \text{C0})}{P(\text{Class} = \text{C0})} = \frac{3/10}{0.6} = 0.5$   
 $P(\text{Shirt Size} = \text{Medium} \mid \text{Class} = \text{C0}) = 0.5$

F)  $P(\text{Shirt Size} = \text{Medium} \mid \text{Class} = \text{C1}) = \frac{P(\text{Shirt Size} = \text{Medium} \cap \text{Class} = \text{C1})}{P(\text{Class} = \text{C1})} = \frac{4/10}{0.6} = \frac{P(\text{Shirt Size} = \text{Medium}) \cdot P(\text{Class} = \text{C1})}{P(\text{Class} = \text{C1})} = \frac{4/10}{0.6} = 0.6667$   
 $P(\text{Shirt Size} = \text{Medium} \mid \text{Class} = \text{C1}) = 0.6667$

b) Use Naive Bayes Classifier to find the class of  $P(\text{Gender} = \text{F} \mid \text{Car Type} = \text{Family}, \text{Shirt Size} = \text{Medium})$

$P(\text{C1} \mid \text{F}, \text{Family}, \text{Medium}) = 0.5 \cdot 0.6 \cdot 0.3 \cdot 0.4 = 0.036$

0.036% of  $P(\text{F} \mid \text{Family} \mid \text{Medium})$

**NOTE:** Julio helped me going through this exercise through Zoom.

## Task 2 - Text Classification

**Data Set:** spam.csv located at <https://www.kaggle.com/uciml/sms-spam-collection-dataset/version/1>  
**Note:** you may want to use `encoding = 'latin-1'` when loading this file (<https://www.kaggle.com/bervozza/spam-classification>)

**Objective:** to classify SMS message as spam or not spam (ham).

From the given data set, use Naive Bayes to classify the SMS message. The framework for text classification is briefly summarized here:

- Transformation of your dataset (change to lower case, remove numbers, remove punctuation, stop words, white space, word stemming, etc)
- Document-Term-Matrix creation - matrix of word counts for each individual document in the matrix (e.g. documents as rows, words as columns or vice versa)
- Text Analysis (e.g. word counts, visualizations using wordclouds)

**Helpful links:**  
<https://machinelearningmastery.com/clean-text-machine-learning-python/>  
<http://textminingonline.com/dive-into-nltk-part-4-word-stemming-and-lemmatization>  
<https://machinelearningmastery.com/prepare-text-data-machine-learning>

**Analysis Questions:**

- What is the accuracy of the model? Report your finding with corresponding graphs.
- Print the 5 most frequent words in each class, and their posterior probability generated by the model.
- How would you improve the model performance?
- If the data set is bigger, do you think the accuracy increases? Discuss.

## Deliverables:

Upload your work from task 1 and your notebook from task 2

**Important:** Make sure your provide complete and thorough explanations for all of your analysis. You need to defend your thought processes and reasoning.

```
# Loading required libraries
import pandas as pd
import seaborn as sns
import re
import nltk
from nltk.corpus import stopwords
from nltk.stem import WordNetLemmatizer
from sklearn.feature_extraction.text import CountVecorizer
from sklearn.model_selection import train_test_split
from sklearn.feature_extraction.text import TfidfVectorizer
from sklearn.metrics import confusion_matrix, accuracy_score, classification_report
from nltk import RegexpTokenizer

import matplotlib.pyplot as plt
%matplotlib inline
```

UsageError: line magic function '%' not found.

Downloading stop-words files from nltk package

```
# nltk downloading stopwords files
nltk.download('stopwords')
nltk.download('wordnet')
```

```
[nltk_data] downloading package stopwords to
[nltk_data] C:\Users\Ahmad\AppData\Roaming\nltk_data...
[nltk_data] Package stopwords is already up-to-date!
[nltk_data] downloading package wordnet to
[nltk_data] C:\Users\Ahmad\AppData\Roaming\nltk_data...
[nltk_data] Package wordnet is already up-to-date!
```

```
True
```

```
spam = pd.read_csv('spam.csv', encoding='latin-1')
spam.head(10)
```

	v1	v2	Unnamed: 2	Unnamed: 3	Unnamed: 4
0	ham	Go until jurong point, crazy. Available only in...	NaN	NaN	NaN
1	ham	Ok lar... Joking wif u oni...	NaN	NaN	NaN
2	spam	Free entry in 2 a wkly comp to win FA Cup fina...	NaN	NaN	NaN
3	ham	U dun say so early hor... U c already then say...	NaN	NaN	NaN
4	ham	Nah I don't think he goes to usf, he lives aro...	NaN	NaN	NaN
5	spam	FreeMsg Hey there darling it's been 3 week's n...	NaN	NaN	NaN
6	ham	Even my brother is not like to speak with me ...	NaN	NaN	NaN
7	ham	As per your request 'Melle Melle (Dru Minnamin...	NaN	NaN	NaN
8	spam	WINNER!! As a valued network customer you have...	NaN	NaN	NaN
9	spam	Had your mobile 11 months or more? U R entitle...	NaN	NaN	NaN

spam.info()

```
<class 'pandas.core.frame.DataFrame'>
RangeIndex: 5572 entries, 0 to 5571
Data columns (total 5 columns):
 #   Column  Non-Null Count  Dtype
---  --
 0   v1      5572 non-null     object
 1   v2      5572 non-null     object
 2   Unnamed: 2    50 non-null      object
 3   Unnamed: 3    12 non-null      object
 4   Unnamed: 4    6 non-null       object
dtypes: object (5)
memory usage: 217.8+ KB
```

Dataset has 5 columns. We only need V1 and V2 all other columns are null and not useful. So, I dropped them and renamed V1 as Class, and V2 as text

```
# drop unused columns
spam = spam.drop(["Unnamed: 2", "Unnamed: 3", "Unnamed: 4"], axis=1)
# rename columns
spam = spam.rename(columns={"v1": "Class", "v2": "text"})
spam.head(10)
```

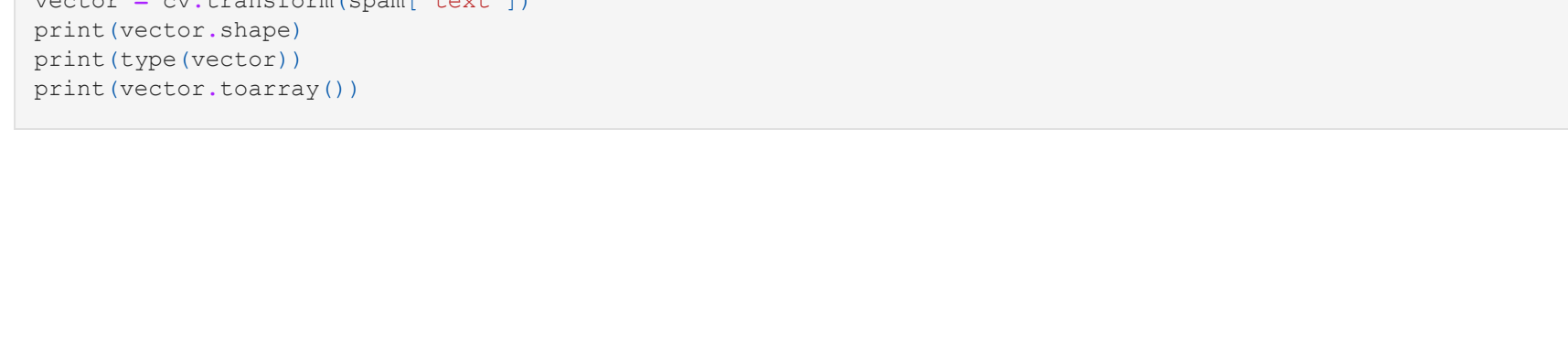
	class	text
0	ham	Go until jurong point, crazy. Available only in...
1	ham	Ok lar... Joking wif u oni...
2	spam	Free entry in 2 a wkly comp to win FA Cup fina...
3	ham	U dun say so early hor... U c already then say...
4	ham	Nah I don't think he goes to usf, he lives aro...
5	spam	FreeMsg Hey there darling it's been 3 week's n...
6	ham	Even my brother is not like to speak with me ...
7	ham	As per your request 'Melle Melle (Dru Minnamin...
8	spam	WINNER!! As a valued network customer you have...
9	spam	Had your mobile 11 months or more? U R entitle...

```
# column that store text length
len_text=[]
for i in spam['text']:
    len_text.append(len(i))
spam['text_length']=len_text
```

spam.head()

	class	text	text_length
0	ham	Go until jurong point, crazy. Available only in...	111
1	ham	Ok lar... Joking wif u oni...	29
2	spam	Free entry in 2 a wkly comp to win FA Cup fina...	155
3	ham	U dun say so early hor... U c already then say...	49
4	ham	Nah I don't think he goes to usf, he lives aro...	61

```
plt.figure(figsize=(12,5))
spam[spam['class']=='spam']['text_length'].plot(bins=35,kind='hist',color='blue',label='spam',alpha=0.5)
plt.legend()
plt.xlabel('message length')
plt.show()
```



```
plt.figure(figsize=(12,5))
spam[spam['class']=='ham']['text_length'].plot(bins=35,kind='hist',color='red',label='ham',alpha=0.5)
plt.legend()
plt.xlabel('message length')
plt.show()
```



We can see that most of spam messages' length is between 100 and 150 while ham is between 40 and 70.

The function below is to remove stop-words, punctuations, and cnvert upper-case t lower-case.

```
def process(sentence):
    text = re.sub('[^a-zA-Z]', ' ', sentence).split()
    words = [x.lower() for x in text if x not in stopwords.words('english')]
    lemma = WordNetLemmatizer()
    word = [lemma.lemmatize(word,'v') for word in words]
    word = ' '.join(word)
    return word
```

```
spam['text'] = spam['text'].apply(process)
```

spam.head()

	class	text	text_length
0	ham	go jurong point crazy available bugis n great ...	111
1	ham	ok lar joke wif u oni	29
2	spam	free entry 2 wkly comp win fa cup final tkts 2...	155
3	ham	u dun say early hor u already say	49
4	ham	nah i think go usf five around though	61

The data now is clean. Stopwords, punctuations removed, and upper-cases converted to lower-cases.

```
cv = CountVecorizer()
cv.fit(spam['text'])
print(cv.vocabulary_)
vector = cv.transform(spam['text'])
print(vector.shape)
print(type(vector))
print(vector.toarray())
```



[illegible]







[illegible]



