

## Mohammad Firas Sada - A20401140

### CS 481 Spring 2023 Programming Assignment #02

Due: **Sunday, April 2, 2023 at 11:59 PM CST**

Points: **100**

#### Instructions:

1. Place **all your deliverables (as described below) into a single ZIP** file named:

`LastName_FirstName_CS481_Programming02.zip`

2. Submit it to Blackboard Assignments section before the due date **[presentation slides can be added AFTER you presented]**. **No late submissions will be accepted.**

#### Objectives:

1. (100 points) Implement and evaluate a Naïve Bayes classifier algorithm.

#### Task:

Your task is to implement, train, and test a Naïve Bayes classifier using a publicly available data set.

#### Data set:

Pick a publicly available data (**follow the guidelines provided in Blackboard**) set first and do an initial exploratory data analysis. Note

#### Deliverables:

Your submission should include:

- Python code file(s). Your py file should be named:

`cs481_P02_XXXXXXXXX.py`

where XXXXXXXXXX is your IIT A number (**this is REQUIRED!**). If your solution uses multiple files, makes sure that the main (the one that will be run to solve the problem) is named that way and others include your IIT A number in their names as well.

- Presentation slides in PPTX or PDF format. **Slides can be added to your submission AFTER you presented your work in class [You can resubmit everything then]**. Name it:

`LastName_FirstName_CS481_P02_Slides.pptx or pdf`

- This document with your observations and conclusions. You should rename it to:

`LastName_FirstName_CS481_P02.doc or pdf`

## Implementation:

Your task is to implement, train, and test a Naïve Bayes classifier (as outlined in class) and apply it to classify sentences entered using keyboard.

Your program should:

- Accept one (1) command line argument, i.e. so your code could be executed with

```
python cs481_P02_XXXXXXXXX.py IGNORE
```

where:

- `cs481_P02_XXXXXXXXX.py` is your python code file name,
- `IGNORE` is a YES / NO switch deciding if your implementation will skip one pre-processing step (the one selected by you in Google Spreadsheet for the assignment),

Example:

```
python cs480_P01_A11111111.py YES
```

If the number of arguments provided is NOT one (none, two or more) the argument is neither YES nor NO assume that the value for `IGNORE` is NO.

- Load and process input data set:
  - Apply any data clean-up / wrangling you consider necessary first (mention and discuss your choices in the Conclusions section below).
  - Text pre-processing:
    - ◆ treat every document in the data set as a single sentence, even if it is made of many (no segmentation needed),
    - ◆ if `IGNORE` is set to YES, skip one (selected earlier) of the steps below:
      - apply lowercasing,
      - remove all stop words (use the `stopwords` corpora for that purpose),
      - stem your data using NLTK's Porter Stemmer (NO lemmatization necessary)
- Train your classifier on your data set:
  - assume that vocabulary `V` is the set of ALL words in the data set (after pre-processing above),
  - divide your data set into:
    - ◆ training set: FIRST (as appearing in the data set) 80% of samples / documents,
    - ◆ test set: REMAINING 20 % of samples / documents,

- use **binary BAG OF WORDS** with “**add-1 smoothing**” representation for documents,
- train your classifier (find its parameters. HINT: use Python dictionary to store them),
- Test your classifier:
  - use the test set to test your classifier,
  - calculate (and display on screen) following metrics:
    - ◆ number of true positives,
    - ◆ number of true negatives,
    - ◆ number of false positives,
    - ◆ number of false negatives,
    - ◆ sensitivity (recall),
    - ◆ specificity,
    - ◆ precision,
    - ◆ negative predictive value,
    - ◆ accuracy,
    - ◆ F-score,
- Ask the user for keyboard input (a single sentence S):
  - use your Naïve Bayes classifier to decide (HINT: use log-space calculations to avoid underflow) which class S belongs to,
  - display classifier decision along with  $P(\text{CLASS\_A} | S)$  and  $P(\text{CLASS\_B} | S)$  values on screen

Your program output should look like this (if pre-processing step is NOT ignored, output NONE):

```
Last Name, First Name, AXXXXXXXX solution:
Ignored pre-processing step: STEMMING
```

```
Training classifier...
Testing classifier...
Test results / metrics:
```

```
Number of true positives: xxxx
Number of true negatives: xxxx
Number of false positives: xxxx
Number of false negatives: xxxx
Sensitivity (recall): xxxx
Specificity: xxxx
```

Precision: **xxxx**  
Negative predictive value: **xxxx**  
Accuracy: **xxxx**  
F-score: **xxxx**

Enter your sentence:

Sentence S:

**<entered sentence here>**

was classified as **<CLASS\_LABEL here>**.

P(**<CLASS\_A>** | S) = **xxxx**

P(**<CLASS\_B>** | S) = **xxxx**

Do you want to enter another sentence [Y/N]?

If user responds Y, classify new sentence (you should not be re-training your classifier).

where:

- **xxxx** is an actual numerical result,
- **<entered sentence here>** is actual sentence entered y the user,
- **<CLASS\_LABEL here>** is the class label decided by your classifier,
- **<CLASS\_A>**, **<CLASS\_B>** are available labels (SPAM/HAM, POSITIVE/NEGATIVE, etc.).

### Classifier testing results:

Enter your classifier performance metrics below:

With ALL pre-processing steps:	Without <b>REMOVING STOP WORDS</b> step:
Number of true positives: <b>9585</b>	Number of true positives: <b>9455</b>
Number of true negatives: <b>6247</b>	Number of true negatives: <b>6411</b>
Number of false positives: <b>1685</b>	Number of false positives: <b>1521</b>
Number of false negatives: <b>2483</b>	Number of false negatives: <b>1521</b>
Sensitivity (recall): <b>0.79</b>	Sensitivity (recall): <b>0.78</b>
Specificity: <b>0.79</b>	Specificity: <b>0.81</b>
Precision: <b>0.85</b>	Precision: <b>0.86</b>
Negative predictive value: <b>0.72</b>	Negative predictive value: <b>0.71</b>
Accuracy: <b>0.79</b>	Accuracy: <b>0.79</b>
F-score: <b>0.82</b>	F-score: <b>0.82</b>

What are your observations and conclusions? When did the algorithm perform better? a summary below

### Summary / observations / conclusions

In this assignment we are asked to implement a naive bayes classifier to perform sentiment analysis on a dataset. The task involved implementing a naive bayes classifier from scratch to perform sentiment analysis on a dataset of 4,000,000 Amazon product reviews. The code was then tested on a subset of 100,000 reviews. The classifier was trained using the bag-of-words model with add-1 smoothing and binary representation. The dataset was split into 80% training and 20% testing sets.

The results showed that the classifier was able to correctly classify the sentiment of the reviews almost 80% of the time, regardless of whether stop words were removed or not. This suggests that the classifier is effective at distinguishing between positive and negative sentiments in text data.

Here are some additional details about the implementation:

- The bag-of-words model is a simple way to represent text data as a set of words. Each word is assigned a unique ID, and the frequency of each word in the text is recorded.
- Add-1 smoothing is a technique that is used to prevent the naive Bayes classifier from overfitting the training data. It does this by adding a small constant value to the probability of each word in the vocabulary.
- Binary representation is a way to represent text data as a set of binary values. Each word is assigned a value of 1 if it appears in the text and 0 if it does not.
- The dataset was split into 80% training and 20% testing sets. This means that the classifier was trained on 80,000 reviews and tested on 20,000 reviews.
- The results showed that the classifier was able to correctly classify the sentiment of the reviews more than 79% of the time. This means that the classifier was able to correctly identify whether the reviews were positive or negative.
- Skipping removing stop words degrades the accuracy but only by a small margin. This can be because we are training our model on a fairly large dataset with an abundance of training data.