

# Assignment 6: Naïve Bayes Classifier

Assignment written by Sara Sood (while at Pomona College),  
Northwestern University with changes by Nathan Shelly

You may work alone or in a pair. If you work in a pair, both students must submit the file on Canvas to get credit. Additionally, please put both students' NetIDs at the top of the file as a comment so that there's no academic integrity issues. For example:

```
# aaa0000, zzz9999
```

The purpose of this assignment is to apply what you've learned about supervised Machine Learning and practice using python dictionaries and the pickle library.

In particular, you will implement a Naïve Bayes Classifier that analyzes the *sentiment* conveyed in text. When given a selection of text as input, your system will return whether the text is positive, or negative. The system will use a corpus of movie reviews as training data, using the number of stars assigned to each review by its author as the truth data (1 star is negative, 5 stars is positive).

## Overview

Here's a high level overview of a Bayes Classifier. The **Part 3: Start Classifying** section further below provides a more in-depth, mathier explanation.

Let's say we have a movie review called `mystery-?.txt`. The text of our review is **amazing movie**. Let's also say that for every word we know the probability of it existing in the positive reviews and the probability of it existing in the negative reviews. For example, we might know that the word **amazing** shows up 8 times in positive reviews out of 1032 total words (for a probability of  $\frac{8}{1032} = .0077$ ) and 2 times in negative reviews out of 743 reviews (for a probability of  $\frac{2}{743} = .0027$ ). See the training sections below for an explanation of how we calculate this.

Okay, now we're ready to classify our document. We're going to calculate two probabilities, the probability that our document is positive and the probability that it is negative. Whichever probability is higher will be the class we output. To calculate the probability of the document being positive we will combine the probability of each individual word being positive. For now we'll just multiply the probabilities together, below we mention issues with this and a better way to combine them. For example here's a calculation of positive (made up data):

$$P(\text{positive} \mid \text{mystery} - ?.txt) = P(\text{"amazing"} \mid \text{positive}) * P(\text{"movie"} \mid \text{positive}) = .0077 * .0016$$

For anyone who hasn't seen the above notation,  $|$  means given, so  $P(\text{some\_event} \mid \text{some\_data})$  means the probability of *some\_event* given *some\_data*. Here we want to calculate the probability of a positive review given the document. To do so we assume the document is positive and multiply the probability of each word given that knowledge. We would then repeat the same calculation for negative. Finally, we would output whichever class had the higher probability.

Important to note that we've simplified one or two things above that we'll explain more in the **Part 3: Start Classifying** section. Also important, Naïve Bayes classifiers are complicated. Even for people with a strong grasp of probability, the concept can be confusing. Please feel free to ask lots of questions and we'll try and help as best we can!

# The BayesClassifier Class

This class has the following attributes:

- **pos\_freqs** - dictionary of frequencies of positive words. This is a dictionary of **str**->**int** where keys are words (tokens) and values are the count of that word across the training texts. For example if the positive reviews had 8 instances of the word "Tomate" then **pos\_freqs["Tomate"]** would equal 8.
- **neg\_freqs** - dictionary of frequencies of negative words. See **pos\_freqs** above except for negative words.
- **pos\_filename** - name of positive dictionary cache (pickled) file
- **neg\_filename** - name of negative dictionary cache (pickled) file
- **training\_data\_directory** - relative path to training directory
- **neg\_file\_prefix** - prefix of negative reviews
- **pos\_file\_prefix** - prefix of positive reviews

It should also contain the following methods (a few have already been provided, those are marked below):

- **\_\_init\_\_(self)** (*provided*): Initializes various attributes described again and begins training (by calling **train**) if it does not detect **pos.dat** and **neg.dat** data files
- **train(self)**: Trains the classifier by generating **pos\_freqs** and **neg\_freqs** dictionaries on all files in the folder given by **training\_data\_directory**.
- **classify(self, text)**: Provides a classification for given text. We'd hope **some\_classifier.classify("I love Tomate")** would return **positive**.
- **load\_file(self, filepath)** (*provided*): Loads file given by filepath, returning the text of the file.
- **save\_dict(self, dict, filepath)** (*provided*): Saves (by pickling) given dictionary to the provided filepath.
- **load\_dict(self, filepath)** (*provided*): Loads a pickled dictionary at the given filepath, returning the dictionary.
- **tokenize(self, text)** (*provided*): Splits given string into a list of individual words (tokens) in the text (all lowercase).
- **update\_dict(self, words, freqs)**: Updates the frequency of words in the **freqs** dictionary with the given **words** list. This involves incrementing the count of each word in **words** in the **freqs** dictionary or adding words with a starting count of 1 if they don't already exist in **freqs**.

## Your Work

### Part 1: The Training Data

To train your system, you will be using data from a corpus of movie and product reviews collected from RateItAll. The reviews come from a variety of domains, e.g., movies, music, books, and politicians and are stored in individual text files where the file name is of the form *domain-number\_of\_stars-id.txt*. For example, the file "movies-1-32.txt" contains a movie review that was assigned one star by its author, and has an id of 32. Please note that we will not use the id number in this assignment so feel free to ignore it.

We will only be training using reviews that had either one or five stars, where the reviews with one star are the "negative" training data and the reviews with five stars are the "positive" training data.

Two provided functions **load\_file** and **tokenize** will be helpful here. Familiarize yourself with these functions by using them at the python interpreter to read reviews from the training data files. Take a moment to explore the code provided in the **a6.py** file so that you don't create extra work for yourself. If you would like to edit the provided code, please ask on Piazza first to make sure it won't affect grading.

### Part 2: Train the System

Now that you're familiar with the training data, our next task will be to train the classifier (filling in the body of the **train** method).

Given a **target document** (a document to be classified), Naïve Bayes Classifiers calculate the probability of each **feature** (i.e. each individual word) in the target document occurring in a document of each **document class** (i.e. positive and negative) in order to find the class that is most probable. For example, in our review from before called “mystery-?.txt” with the text “amazing movie” the features would be “amazing” and “movie”.

In order to make these calculations, the system must store the number of times that each feature occurs in each document class in the training data. For example, given a word like “good,” it may occur 329 times in the positive documents, but only 53 times in the negative documents. There are two different types of tallies that are commonly used: **presence** refers to the number of documents in the training data that the feature occurs in and **frequency** refers to the number of times that the feature occurs in the training set (i.e. if it occurs three times in one document, that adds three to the count). These instructions use **frequency** but feel free to experiment with both.

Naïve Bayes Classifiers typically use a database to store all the counts. However, since you are writing this code in Python, we’ll instead use a handy Python utility called **pickle**. Using **pickle**, one can write a data structure to a file, and then load it back into memory later. In particular, you’ll have two dictionaries, one holding all the words that occur in positive documents keyed to the frequencies at which they occurred, and the corresponding dictionary for negative documents. If the positive reviews had 300 unique words the positive dictionary would have 300 keys, with the corresponding value equal to the number of times that word appeared across all positive documents. Since these dictionaries are time consuming to construct, it is useful to only calculate them once, **pickle** them, and then load them into memory the next time you need to use them.

You’re now ready to start implementing training. Before filling out the **train** method, implement the **update\_dict** method as you will want to use this in your implementation of **train**. **update\_dict** takes a list of words and a dictionary and updates the counts in the dictionary with the new words in the given words list (see the docstring in the code for more details).

Once you’ve implemented **update\_dict** you’re now ready to fill in the body of the **train** method. Notice that we’ve started it for you by giving you a list of the training files. Here is a loose overview of the steps you’ll need to take:

- For each file name, parse the file name and determine if it’s a positive or negative review. **Hint:** the **startswith** method on strings or the keyword **in** may be helpful here.
- For each word in the file, update the frequency for that word in the appropriate dictionary. **Hint:** **update\_dict** might come in handy here.
- Save these dictionaries using **pickle** so that the system only has to calculate them once. We’ve provided you with a **save\_dict** function to help with this step (the complimentary function **load\_dict** is used in **\_\_init\_\_**).

More detailed instructions are provided in the comments in the provided code. Be sure to notice the attributes provided in the **\_\_init\_\_** method. In particular, **self.pos\_freqs** and **self.neg\_freqs** are the dictionaries you’ll use to store the data.

After writing the **train** method, you can test it by creating an instance of a **BayesClassifier** because the constructor (**\_\_init\_\_**) method calls the **train** method. First, run your code. Then, in the IDLE shell, you can type:

```
b = BayesClassifier()
```

If you make any changes to your code, you must rerun your code before calling the constructor from the Shell window. Additionally, if you change your code and want it to retrain, you’ll have to delete the **.dat** files that have been produced/pickled (**pos.dat** and **neg.dat**).

## Part 3: Start Classifying

At this point, you now have all the data you need stored in memory, and can start classifying text! Given a document with some text, the goal is for your system to output the correct document class (i.e. positive or negative). In particular, you'll calculate the conditional probability of each document class given the features in the target document (e.g.  $P(\text{positive} \mid \text{word}_1, \text{word}_2, \dots)$ ) and return the document class with the highest probability.

**IMPORTANT:** In the **Overview** section above we made a few simplifications. Your classifier **must** classify according to the instruction in the following paragraphs, not the simplified version.

The first important simplification we made was using the product of conditional probabilities. This causes an issue called underflow. Because we are multiplying so many fractions, the product becomes too small to be represented. The standard solution to this problem (which we're going to implement) is to calculate the sum of the logs (base e) of the probabilities, as opposed to the product of the probabilities themselves. Our earlier example with `mystery-?.txt` would now look like this:

$$P(\text{positive} \mid \text{mystery} - ?.txt) = \\ \log(P(\text{"amazing"} \mid \text{positive})) + \log(P(\text{"movie"} \mid \text{positive})) = \log(.0077) + \log(.0016) = -2.11 + -2.80$$

**Note:** this results in negative probabilities. We will still output whichever class has the larger probability, in this case, the least negative one.

Another important piece is called smoothing (specifically we'll be doing "add one smoothing"). This problem is more subtle as it won't produce a bug, but will throw off your calculations. Suppose a feature  $f$  (say "amazing") occurs some number of times in the training data for class *positive* and not at all in the training data for class *negative*. If the document we're trying to classify contains feature  $f$ , the probability for class *negative* (assuming all else equal) will be higher than the probability for class *positive*, which intuitively should not be the case. Additionally, we run into issues with trying to take the log of 0. To solve this problem we pretend that the feature has a count of 1 in the class where it doesn't exist in the training data. To account for the fact that we added 1 to one class we need to also add 1 to the count in each other class. We'll discuss this in class more; you can also read more online by searching "add one smoothing".

In case you prefer equations, assume we are trying to calculate the probability that a target document  $d$  is positive given the set of  $n$  features of  $d$ , which we'll call  $f$ :

$$P(\text{positive} \mid f) = \sum_{i=1}^n \log(P(f_i \mid \text{positive}))$$

The term on the right side of the equation is the sum of the conditional probabilities of each of the features occurring given that a document is positive.

Here is an expanded version of this calculation (for our "amazing movie" example):

$$P(\text{positive} \mid f) = \log(P(\text{"amazing"} \mid \text{positive})) + \log(P(\text{"movie"} \mid \text{positive}))$$

**IMPORTANT:** The equations above don't show "add one smoothing" but you *must* implement it for an accurate classifier.

Fill in the body of the `classify` method which takes as input a string of text - like "amazing movie" (not a file name). It should return a string of the classification - either "positive" or "negative". More details are provided in comments inside the method.

Once you've implemented it you should now have a classifier that can be used in the following manner:

```
bc = BayesClassifier()
print(bc.classify("I love python!")) # should print positive
print(bc.classify("I hate python!")) # should print negative
```

## Possible Extensions (optional)

Congratulations! Once you reach this point, you will have built your first Naïve Bayes Classifier! In the previous part, you probably realized some ways in which you can improve upon your system. In this part, try and implement these ideas and build *your best classifier*. For this part, it might make sense to create a copy of your code and store it in a separate file so that you don't break your existing classifier.

The most important problem with respect to Naïve Bayes Classifiers is the choice of features used in classification. Thus far, you've only used "unigrams" (i.e. single words) as features. Another common feature is "bigrams" (i.e. two word phrases). Another feature might be the length of the document, or amount of capitalization or punctuation used. You could also remove frequently occurring words (`stop_words`) that are often thought of as noise. See `sorted_stoplist.txt` (provided) or store the frequencies of the stems of the words instead of the words themselves.

Further, you could attempt to write some code to evaluate your system measuring accuracy, precision, recall and f-measure. You could also change your classify function to attempt to add a "neutral" class.

Additionally, you could also account for any imbalances in the number of documents in each class. In an extreme example, if we had 1,000,000 positive reviews and 1 negative review we would want to bias our classifier towards a positive classification. To do so, we could add the prior probability of the document being a given class to each class calculation. That is, when calculating the probability that our text is positive we would add the probability of a document being positive based on the number of documents, in this case  $1,000,000/1,000,001 = 0.9999$ . To the negative probability we'd add  $1/1,000,001 = 0.00001$ . In practice, the prior probability can be somewhat artificial, which is why we did not include it in the calculation.