

Name: Adarsha Poudel

Perm: 3347028

Email: adarsha@ucsb.edu

Architecture:

I have one struct named "Label" and one class called "NaiveBayesClassifier" with one parametrized constructor and four functions.

struct Label

This struct Label is used to represent the features/words of the document with their positive review count and negative review count. E.g. Label l ("love", 5, 6) means l stores the feature "love" which has 5 positive reviews and 6 negative review.

NaiveBayesClassifier: This class uses hash table named "entries" to store features and its positive/negative count. It then uses different functions and constructor to train and test data.

Constructor: NaiveBayesClassifier(std::string);

This constructor is used to train data. It takes in the filename as a string. Then it opens the file, splits each line of document into words(features) and their label. After that, it uses "insert" function to insert the word and their label into the hash table.

void insert(std::string, int);

This function is called by the constructor to help train the data. It gets the feature as string and its class(review) as an integer (0 or 1). If the feature already exists, it just increments the count of review of that word. If it does not, it creates new struct object and inserts the feature into the hash table.

int hash(std::string);

This function returns the index of my hash table given the feature as string by moding the product of ascii value of the characters in that string by the size of hashtable.

double test(std::string);

This function takes the filename of testing data as string, opens that file, and tests the file. It uses "return Probability" function to get the probability of certain feature belonging to some class, i.e, $P(\text{word}_i | \text{positive/negative review})$. It then manipulates the probabilities, compares them, and return the accuracy of our model as a double.

double returnProbability(std::string, int label);

This function is helper function for "test" which returns probability of a word belonging to certain label, i.e., $P(\text{word}_i | \text{positive/negative review})$.

Preprocessing:

I represented each document as bag-of-words model where features are the words in the document d_i . Keeping in mind this model loses order-specific information, I tried to improve the bag of words model introducing bigram to my features. So, each feature represents adjacent two words without the space between them. For example, if the document is "I love you": "Ilove", "loveyou" are the features. Preprocessing is done in my constructor of NaiveBayesClassifier class.

Model Building:

For each features/bigrams word that I extracted in preprocessing, I stored the positive review count and negative review count of that feature. So, when testing, for each feature in the testing text, I calculated the $P(\text{word} | +ve/-ve \text{ review})$ using this formula:

$$P(\text{word} | +ve/-ve \text{ review}) = \frac{\text{Number of } +\frac{+ve}{-ve} \text{ review count of that word} + \alpha}{\text{Total words in } +\frac{+ve}{-ve} \text{ review} + \alpha * \text{vocabulary}}$$
 where α is

arbitrary number and vocabulary is number of unique words in the training text. Doing this I prevented the situation where the unseen feature produces 0 probability and affects the MAP. I tried different alphas; however, $\alpha = 1.0$ gave me the best accuracy.

After calculating conditional probabilities of each feature given certain class (positive or negative review), I added logs of probabilities of each feature instead of multiplying each probability in one document to prevent overflow. I then compared the value of $P(\text{Word1, Word2, ..., WordN} | +ve \text{ review})$ and $P(\text{Word1, Word2, ..., WordN} | -ve \text{ review})$, and classified that document accordingly.

Challenges:

After implementing hash table on my own without any library, the huge problem for me was to determine a good hash table size so that hashtable perform the most efficient. So, I read up one paper, where they described that the size of hashtable should always be prime and be about $1.3 * \text{total keys you are inserting}$. Hence, I found out the total number of unique words in training text to be around 5000. So, I choose the prime 6577 to meet the criteria.

Implementing the simple naïve model using bag of words model only gave me 85 percent accuracy. To improve the accuracy, I incorporated bigram feature and was able to boost my accuracy to 0.88.

My next challenge was the overflow by multiplying many probabilities of features. I overcame this struggle using logarithmic property of addition: $\log(xy) = \log(x) + \log(y)$

Weaknesses:

Since I am using bags of words model, I lose order-specific information of the features, I also limit the context of my words in the document. While I did use bigram to value order-specific information of words, I think the mixture of trigram/bigram would have made an accurate model. I could not implement trigram, because it took longer than 10 minutes for me to execute my program. Now for the context-based information of the words, I could have used stop-word removal to remove words that does not add weight to the reviews.

Results:

I have defined important features as the words which has the highest positive/negative ratio for positive review and highest negative to positive ration for the negative review. This gives us the word with more predictive power.

Notice each word here represents a bigram: two words without spaces.

```
27 seconds (training)
28 seconds (labeling)
0.980 (training)
0.886 (testing)
The top ten important positive review words are:
onlybad: Positive to Negative Ratio = 31.000
asteal: Positive to Negative Ratio = 39.000
allages: Positive to Negative Ratio = 45.500
highlyrecommended: Positive to Negative Ratio = 29.167
definitelyrecommend: Positive to Negative Ratio = 28.333
greatreplay: Positive to Negative Ratio = 31.000
cantwait: Positive to Negative Ratio = 28.000
onlydownside: Positive to Negative Ratio = 50.000
justfun: Positive to Negative Ratio = 28.000
excellentproduct: Positive to Negative Ratio = 68.000

The top ten important negative review words are:
allcosts: Negative to Positive Ratio = 70.000
wastedmy: Negative to Positive Ratio = 49.000
dontbuy: Negative to Positive Ratio = 48.000
notwaste: Negative to Positive Ratio = 70.000
worstgame: Negative to Positive Ratio = 47.000
techsupport: Negative to Positive Ratio = 81.000
icontacted: Negative to Positive Ratio = 62.000
notpurchase: Negative to Positive Ratio = 64.000
ofjunk: Negative to Positive Ratio = 45.000
arefund: Negative to Positive Ratio = 45.000
beashamed: Negative to Positive Ratio = 51.000
```