

**Writeup Question 1.1:** There are two other ways we could have handled the square brackets. We could have deleted those words completely, or we could have removed the square brackets but left the words untagged. Why did we do it the way we did (with tags)? What are some pros and cons of the three different ways of handling the square brackets? Give at least one pro or one con for each way.

Answer: by tagging the words, we are making sure that we will not lose information in the piece of text, and we can explicitly mark the origin of the words.

Pros:

Tag all words: can keep all the information.

Deleted completely: easy to achieve.

Removed brackets: retain all information, and also easy to achieve.

Cons:

Tag all words: it is time consuming.

Deleted completely: we will lose some information by doing this.

Removed brackets: we will not be able to distinguish which part were not from the author.

**Writeup Question 1.2:** What data type did you use to save the negation-ending tokens? Explain your choice.

Answer: I used a set to store the negation-ending tokens. Because a set can provide me with  $O(1)$  complexity of verifying whether a token belongs to it.

**Writeup Question 3.1:** Looking at the formulae for precision, recall, and f-measure, what does each of them measure? Why do we need all three of them?

Answer: **Precision** means among all the labels that are predicted true, the ratio of the labels that are actually true. **Recall** means among all the labels that are actually true, the ratio of the labels that are also predicted true. **F-measure** is the measure of a test's accuracy and is defined as the weighted harmonic mean of the precision and recall of the test. We need all three of them so we can make a generic evaluation of the model, plus, the calculation of **f-measure** also requires the calculation of **precision** and **recall**.

**Writeup Question 3.2:** What is the performance of the GaussianNB model?

Answer: The results are below, as we can see, the model is very precise.

**Writeup Question 3.3:** What is the performance of the Logistic-Regression model? Discuss which model performed better on this data and why you think that might be.

Answer:

From the results we can see that Logistic-Regression model performs better, this is because NB models and Logistic-Regression treats the data in different patterns, and for data with large size, Logistic-Regression model will usually have better performance.

**Writeup Question 3.4:** Report the top 10 features of your LogisticRegression model. Why do we sort by the absolute value of weight, rather than the actual value? Explain why it is important to do so and what it means in terms of the features.

Answer: The top 10 features of the Logistic-Regression model are below:

The reason we use the absolute value instead of the actual value is because the coefficients in a Logistic-Regression model are log odds ratios. The sign of each number indicates that the direction of the effect. The absolute value of the number indicates the scale of the effect. Thus, we need to use the absolute value instead of the actual value.

**Writeup Question 4.1:** What is the performance of the new model with extra features? What are its top 10 features? How does it compare with the old model without extra features? Why do you think that might be?

Answer:

We can see from the results that the precision and F-measure of the new model decreased but the recall increased. Overall, the changes are small. The top 10 features are same as the previous model, except that the orders are slightly different. The reason is that the newly added feature can increase the information for the model to predict the snippets, thus increasing the overall recall rate, but the precision might not be affected by the added information.

**Writeup Question 5.1:** How long did this homework take you to complete (not counting extra credit)?

Answer: it took me around 5 hours to write the code and writeup questions, which the time has been distributed in 3 days.

**Writeup Question 5.2:** Did you discuss this homework with anyone?

Answer: I discussed the homework with Joel Yin.

**Writeup Question 6.1:** Train, test, and examine the feature weights of a new LogisticRegression model using the new version of score snippet(). What is the performance of the new model with extra features? What are its top 10 features? How does it compare with the model that used the old version of score snippet()?

Answer:

From the results, we can see that, after adding the WordNet, we increased the precision of the model, nevertheless, the recall has been decreased and the F-measure remained roughly the same as the previous model.

The top 10 features are also changed, the last two features in this model '*and*' and '*entertaining*', whereas in previous model, the last two were '*enjoyable*' and '*fun*'.