

Report on Sentiment analysis

Development Environment:

The environment I used for this assignment on sentiment analysis was Jupyter notebook running on python 3 and the packages I used are NumPy, os, nltk. In nltk I used word_tokenize and sent_tokenize for feature extraction purposes. The os package was used to read the given files.

Features:

In this task, a total of eight features were used for sentiment analysis, including the number of tokens in each review, the number of positive words, the number of negative words, the number of exclamation marks, the number of question marks, the number of sentences in the last paragraph, the length of the longest sentence, and the number of unique words in each review. These features were chosen based on their potential relevance to sentiment analysis.

The number of tokens and sentences in a review provide basic structural information about each review, which can be used to understand its length and complexity. The number of positive and negative words, based on the provided lexicons, provide sentiment information and can help determine whether the review has a positive or negative sentiment. The number of exclamation and question marks can indicate the level of emotion in the review. The number of unique words in each review can provide information about the diversity of language used in the review, and the length of the longest sentence can provide information about the complexity of the review's syntax.

By including a mix of syntactic and semantic information in the features, the classifier can better understand the sentiment expressed in the review. Each feature was chosen to provide a different perspective on the review, and together they form a broad range of information that can be used to classify its sentiment.

Overall, the features used in this task were carefully selected to provide a comprehensive understanding of each review, including its structural, sentiment, and syntactic features.

Hyperparameters:

Hyperparameters are parameters that are set before training the model and can have a significant impact on the model's performance. The hyperparameters used in this logistic regression classifier are the number of epochs, batch size, learning rate, and epsilon for cross-entropy loss.

The number of epochs refers to the number of times the model is trained on the entire dataset. A higher number of epochs may improve the model's accuracy, but may also lead to overfitting. In this case, the number of epochs was set to 100.

Batch size refers to the number of samples used in each iteration of training. A smaller batch size can lead to more noise in the updates of the model's weights, but may allow the model to converge faster. In this case, the batch size was set to 32.

Learning rate refers to the step size used in each update of the model's weights during training. A higher learning rate can cause the model to converge faster, but may also cause the model to overshoot the optimal weights and lead to instability. In this case, the learning rate was set to 0.01.

Epsilon for cross-entropy loss is a small value used to avoid taking the logarithm of zero in the calculation of the loss function. In this case, the epsilon value was set to $1e-12$.

The weights of the model were initialized with respect to the number of features and the batch size. The bias was initially set to zero.

I tried different combinations of hyperparameters to find the best combination that would yield the highest accuracy. After tuning the hyperparameters, the best results were obtained with a batch size of 32, learning rate of 0.01, epsilon of $1e-12$, and 100 epochs. Other hyperparameters that were tried did not yield as good of an accuracy.

Regression equation for logistic regression classifier with full features:

The logistic regression equation can be written as:

$$P(y=1|x) = \text{sigmoid}(w_0 + w_1x_1 + w_2x_2 + \dots + w_nx_n)$$

where $P(y=1|x)$ is the probability of the review having a positive sentiment given the features x , w_0 is the bias term, and w_1, w_2, \dots, w_n are the weights of the features x_1, x_2, \dots, x_n . The sigmoid function is used to map the output to a probability between 0 and 1.

The weights are given as follows:

$$W_1 = 0.4630$$

$$W_2 = 2.0419$$

$$W_3 = -2.0068$$

$$W_4 = -0.1999$$

$$W_5 = -0.2439$$

$$W_6 = 0.0111$$

$$W_7 = -1.0435$$

$$W_8 = -1.8966$$

And the bias is: -0.4834

The equation consists of two main parts: the linear combination of the features with their respective weights, and the sigmoid function. The sigmoid function is used to map the output of the linear combination to a probability value between 0 and 1.

In this particular implementation, the logistic regression equation with full features is given by:

$P(y=1|x) = \text{sigmoid} (-0.4834 + 0.4630x_1 + 2.0419x_2 - 2.0068x_3 - 0.1999x_4 - 0.2439x_5 + 0.0111x_6 - 1.0435x_7 - 1.8966x_8)$

where x_1 corresponds to the number of tokens in the review, x_2 to the number of positive words, x_3 to the number of negative words, x_4 to the number of exclamation marks, x_5 to the number of question marks, x_6 to the number of sentences in the last paragraph, x_7 to the length of the longest sentence, and x_8 to the number of unique words.

The weights of each feature are given by the values W_1 to W_8 , and the bias is given by -0.4834 . These weights were obtained after training the logistic regression model with the full set of features using the hyperparameters previously mentioned.

Performance results for Task 1 and 2:

Task 1 (Baseline lexicon-based classifier):

Accuracy: 0.71

Precision: 0.75

Recall: 0.64

F1 Score: 0.69

Task 2 (Logistic Regression Classifier):

Accuracy: 0.5

Precision: 0.5

Recall: 0.5625

F1 Score: 0.5294

Analysis:

Task 1 used a lexicon-based approach which relied solely on the presence of positive and negative words in the reviews to classify their sentiment. The accuracy of this approach was 0.71, which means that it correctly classified 71% of the reviews. The precision was 0.75, which means that out of all the reviews classified as positive, 75% were actually positive. The recall was 0.64, which means that out of all the positive reviews in the dataset, the classifier correctly identified 64% of them. Finally, the F1 score was 0.69, which is the harmonic mean of precision and recall.

Task 2 used a logistic regression classifier with the features described earlier. The accuracy dropped to 0.5, which is significantly worse than the baseline approach. The precision was 0.5, which means that out of all the reviews classified as positive, only 50% were actually positive. The recall was 0.5625, which means that out of all the positive reviews in the dataset, the classifier correctly identified only 56.25% of them. The F1 score was 0.5294, which is lower than the baseline approach.

Overall, it can be concluded that the chosen features were not effective for sentiment classification and a different set of features could have been chosen. The logistic regression classifier did not perform better than the lexicon-based baseline, which suggests that the additional information provided by the features did not contribute to the classification task in a meaningful way. It is also possible that the hyperparameters were not properly tuned, which could have affected the performance of the model.

Ablation study:

An ablation study was conducted to evaluate the impact of each feature on the performance of the logistic regression classifier. The following results were obtained:

Ablation study for num_tokens:
Accuracy: 0.5
Precision: 0.5
Recall: 0.34375
F1 Score: 0.4074074074074074

Ablation study for num_negative:
Accuracy: 0.5
Precision: 0.5
Recall: 0.4375
F1 Score: 0.4666666666666667

Ablation study for num_question:
Accuracy: 0.5
Precision: 0.5
Recall: 0.53125
F1 Score: 0.5151515151515151

Ablation study for num_longest_sentence:
Accuracy: 0.5
Precision: 0.5
Recall: 0.53125
F1 Score: 0.5151515151515151

Ablation study for num_positive:
Accuracy: 0.5
Precision: 0.5
Recall: 0.53125
F1 Score: 0.5151515151515151

Ablation study for num_exclamation:
Accuracy: 0.5
Precision: 0.5
Recall: 0.625
F1 Score: 0.5555555555555556

Ablation study for num_last_paragraph:
Accuracy: 0.5
Precision: 0.5
Recall: 0.374921875
F1 Score: 0.42852040360746496

Ablation study for num_unique_words:
Accuracy: 0.5
Precision: 0.5
Recall: 0.445390625
F1 Score: 0.47111808941409805

The results of the ablation study showed that the most important feature among the set of features was the number of tokens in a review. When this feature was removed, the F1 score and recall were significantly reduced compared to the original evaluation results. This indicates that the number of tokens is a critical feature in determining the sentiment of a review.

On the other hand, the least important feature was found to be the number of exclamation marks. When this feature was removed, the recall and F1 score increased compared to the original evaluation results. This suggests that the number of exclamation marks is not a strong indicator of sentiment in reviews.

Another interesting finding from the ablation study was the difference in importance between the number of negative and positive words features. The study showed that the number of negative words was more important than the number of positive words in determining the sentiment of a review. This is surprising because one would expect that the number of positive words would be more important in determining the sentiment of a positive review.

Overall, the ablation study provides useful insights into the importance of each feature in the logistic regression classifier for sentiment analysis. The results can be used to improve the performance of the model by selecting the most important features and discarding the less important ones.

General Reflections on the Assignment:

This assignment provided me with a great opportunity to apply my knowledge of natural language processing and machine learning techniques to a practical problem. It helped me understand the process of building a sentiment analysis classifier and the importance of feature engineering in improving the performance of the model.

What I Learned from this Assignment:

I learned how to pre-process text data, extract relevant features, and train a classifier for sentiment analysis. I also learned about the importance of hyperparameter tuning and the trade-off between model complexity and performance.

Difficulty Level:

I found this assignment to be moderately challenging. It required a good understanding of natural language processing concepts and programming skills. However, the detailed instructions and guidance provided in the assignment made it easier to complete.

Difficulties Encountered:

I did not encounter any major difficulties while completing this assignment. However, I had to spend some time researching the different types of feature extraction techniques and hyperparameter optimization strategies to select the best ones for this task.

Approach for Next Time:

If I were to approach this assignment again, I would spend more time exploring different feature extraction techniques and experimenting with more advanced machine learning algorithms. I would also try to incorporate deep learning techniques to improve the performance of the sentiment analysis classifier.