

Text Classification with Naive Bayes

20221009 GKS Kim Eunhong

1. Introduction

Text classification is a natural language processing (NLP) task that involves assigning predefined categories or labels to text documents. In this implementation, I utilize the Naive Bayes classifier, a probabilistic algorithm based on Bayes' theorem, which assumes that the features (words in this case) are conditionally independent. The Naive Bayes algorithm is commonly used for text classification tasks due to its simplicity and efficiency. This text classification systems also lies the Naive Bayes algorithm, a probabilistic model rooted in Bayes' theorem.

Bayesian Foundation: Bayesian methods are grounded in probability theory, particularly Bayes' theorem, which calculates the probability of a hypothesis given observed evidence. In the context of text classification, the hypothesis is the category to which a document belongs, and the evidence is the content of the document itself.

Naive Bayes Assumption: The "naive" in Naive Bayes arises from the assumption of conditional independence among features, in this case, the words or tokens in a document. Despite this simplification, Naive Bayes has proven remarkably effective for text classification tasks, demonstrating its resilience and efficiency in the face of large and diverse datasets.

2. Implementation:

2.1 Data Fetching:

I used the `fetch_20newsgroups` dataset from scikit-learn, which contains a collection of approximately 20,000 newsgroup documents across 20 different categories.

2.2 Vocabulary Creation:

I constructed a vocabulary set by tokenizing and lowercasing words from the training dataset. The unique words in the vocabulary set were used to create dictionaries for storing word counts and log probabilities.

2.3 Model Training:

For each category, Laplace smoothing was applied to the word counts, and the probabilities were computed. Log probabilities were then calculated for each word, providing a more stable measure for very small probabilities.

2.4 Model Testing:

The trained model was applied to the test dataset. For each test sample, log probabilities for each category were computed, and the category with the maximum log probability was selected as the predicted category.

3. Performance Evaluation

3.1 Accuracy:

The accuracy of the model on the test dataset was calculated and printed. The accuracy is a measure of the proportion of correctly predicted labels among all test samples. The measured model accuracy was about 76%.

3.2 Confusion Matrix Visualization:

A confusion matrix was generated to provide a detailed breakdown of correct and incorrect predictions across different categories. The confusion matrix was visualized as a heatmap using seaborn.

