

# Analysis of Topic Modeling for Text Classification

Xinrui Yan  
yanxinrui2021@ia.ac.cn

March 2025

## Abstract

This paper presents an analysis of the performance of topic modeling-based text classification. The influence of different parameters, including the number of topics ( $T$ ), the choice of word-level versus character-level tokenization, and varying values of the number of topics ( $K$ ), is explored. The results demonstrate how these factors affect the classification accuracy using a Naive Bayes classifier, and provide insights into the trade-offs involved in optimizing topic modeling parameters for text classification tasks.

## 1 Introduction

Topic modeling is a popular approach in natural language processing (NLP) for uncovering the latent topics that exist within a corpus of text. Latent Dirichlet Allocation (LDA) is one of the most commonly used algorithms for topic modeling. This study investigates the impact of several parameters on the classification performance using LDA, specifically exploring the impact of the number of topics ( $T$ ), tokenization units (words versus characters), and the number of topics per document ( $K$ ). The classification is performed using a Naive Bayes classifier, a widely used probabilistic classifier known for its simplicity and effectiveness, especially in text classification tasks.

## 2 Methodology

The analysis involves the following steps:

- Preprocessing the text data, which involves reading and cleaning the text files from the dataset.
- Experimenting with different values for the number of topics ( $T$ ) in LDA, such as 5, 10, 20, and 30.
- Comparing word-level and character-level tokenization for text representation.
- Running text classification using LDA with different values for the number of topics ( $K$ ) on both short and long texts.
- Evaluating the classification accuracy using a Naive Bayes classifier and 10-fold cross-validation.

The Naive Bayes classifier assumes that the features (words or characters) are conditionally independent given the class, which simplifies the calculation of class probabilities. The classification accuracy is computed using the formula:

$$\text{Accuracy} = \frac{\text{Number of Correct Predictions}}{\text{Total Number of Predictions}} \times 100 \quad (1)$$

## 3 Experimental Results

The following results were obtained from the experiments on the dataset. The performance was measured using accuracy for different values of the number of topics ( $T$ ), tokenization units (word-level vs. character-level), and the number of topics in each document ( $K$ ).

### 3.1 Effect of Different Number of Topics ( $T$ )

The classification accuracy was measured for various values of  $T$ , and the results are shown in Table 1.

From the results, we can observe that the classification accuracy fluctuates as the number of topics changes. For smaller values of  $T$  (such as 5 and

Number of Topics ( $T$ )	Average Accuracy
5	0.2440
10	0.2390
20	0.2470
30	0.1950

Table 1: Classification accuracy for different values of the number of topics ( $T$ )

10), the accuracy is relatively stable, but as  $T$  increases beyond 20, the performance begins to decrease, suggesting that too many topics can introduce noise and make the model less effective.

### 3.2 Effect of Tokenization Units (Word-level vs. Character-level)

We compared the classification results using word-level tokenization and character-level tokenization. The accuracy for both approaches is summarized in Table 2.

Tokenization Method	Average Accuracy
Word-level	0.2440
Character-level	0.2070

Table 2: Comparison of classification accuracy using word-level and character-level tokenization

The results show that word-level tokenization outperforms character-level tokenization. Word-level tokenization captures more meaningful semantic information, which is crucial for the classification task, whereas character-level tokenization tends to break down the text into more granular components that are less informative for the model.

### 3.3 Effect of the Number of Topics per Document ( $K$ ) for Short and Long Texts

The classification accuracy was also tested with different values of  $K$  for both short and long texts. Table 3 summarizes the results.

Number of Topics ( $K$ )	Short Text Accuracy	Long Text Accuracy
20	0.2440	0.2670
100	0.2070	0.2280
500	0.2400	0.2420
1000	0.2670	0.2550
3000	0.2260	0.2460

Table 3: Classification accuracy for different values of  $K$  for short and long texts

The results indicate that the number of topics ( $K$ ) affects the performance differently for short and long texts. For short texts, a smaller value of  $K$  (such as 20) yields better accuracy, while for long texts, the performance increases as  $K$  becomes larger (such as 1000). This suggests that long texts provide more information, allowing the model to benefit from a larger number of topics.

## 4 Conclusion

This study demonstrates the impact of various parameters on the performance of topic modeling-based text classification. The number of topics ( $T$ ), the choice of tokenization unit (word-level vs. character-level), and the number of topics per document ( $K$ ) all influence the classification accuracy. The Naive Bayes classifier, used for classification in this study, was found to be effective for text classification tasks. Key findings include the optimal choice of tokenization (word-level), the effect of topic number on performance, and how text length affects model performance. Future work can explore further optimizations, such as hyperparameter tuning and the use of more sophisticated models.

## 5 References

### References

- [1] Baum, L. E. (1972). An inequality and associated maximization technique in statistical estimation for probabilistic functions of Markov processes.

- In *Inequalities III: Proceedings of the 3rd Symposium on Inequalities* (pp. 1-8). Academic Press.
- [2] Baum, L. E., & J. A. Eagon. (1967). An inequality with applications to statistical estimation for probabilistic functions of Markov processes and to a model for ecology. *Bulletin of the American Mathematical Society*, 73(3), 360–363.
  - [3] Baum, L. E., & T. Petrie. (1966). Statistical inference for probabilistic functions of finite-state Markov chains. *Annals of Mathematical Statistics*, 37(6), 1554–1563.
  - [4] Dempster, A. P., N. M. Laird, & D. B. Rubin. (1977). Maximum likelihood from incomplete data via the EM algorithm. *Journal of the Royal Statistical Society*, 39(1), 1-21.
  - [5] Eisner, J. (2002). An interactive spreadsheet for teaching the forward-backward algorithm. In *Proceedings of the ACL Workshop on Effective Tools and Methodologies for Teaching NLP and CL*.
  - [6] Forney, Jr., G. D. (1973). The Viterbi algorithm. *Proceedings of the IEEE*, 61(3), 268–278.
  - [7] Jelinek, F. (1997). *Statistical Methods for Speech Recognition*. MIT Press.
  - [8] Kruskal, J. B. (1983). An overview of sequence comparison. In D. Sankoff and J. B. Kruskal (Eds.), *Time Warps, String Edits, and Macromolecules: The Theory and Practice of Sequence Comparison* (pp. 1–44). Addison-Wesley.
  - [9] Markov, A. A. (1913). Essai d’une recherche statistique sur le texte du roman “Eugene Onegin” illustrant la liaison des epreuve en chain. *Izvestia Imperatorskoi Akademii Nauk*, 7, 153–162.
  - [10] Markov, A. A. (2006). Classical text in translation: A. A. Markov, an example of statistical investigation of the text Eugene Onegin concerning the connection of samples in chains. *Science in Context*, 19(4), 591–600. Translated by David Link.
  - [11] Needleman, S. B., & C. D. Wunsch. (1970). A general method applicable to the search for similarities in the amino-acid sequence of two proteins. *Journal of Molecular Biology*, 48, 443–453.

- [12] Rabiner, L. R. (1989). A tutorial on hidden Markov models and selected applications in speech recognition. *Proceedings of the IEEE*, 77(2), 257–286.
- [13] Rabiner, L. R., & B. H. Juang. (1993). *Fundamentals of Speech Recognition*. Prentice Hall.
- [14] Reichert, T. A., D. N. Cohen, & A. K. C. Wong. (1973). An application of information theory to genetic mutations and the matching of polypeptide sequences. *Journal of Theoretical Biology*, 42, 245–261.
- [15] Sakoe, H., & S. Chiba. (1971). A dynamic programming approach to continuous speech recognition. In *Proceedings of the Seventh International Congress on Acoustics*, Volume 3, Akadémiai Kiadó.