

# Analisis Sensitivitas Parameter Model EFCM Berbasis BERT untuk Pendekslsian Topik = Parameter Sensitivity Analysis of BERT-based EFCM Model for Topic Detection

Yudhistira Jinawi Agung, author

Deskripsi Lengkap: <https://lib.ui.ac.id/detail?id=9999920535888&lokasi=lokal>

---

## Abstrak

Pendekslsian topik adalah suatu proses untuk mendapatkan pokok bahasan atau topik pada suatu dokumen teks. Pada data yang besar, pendekslsian topik dapat dilakukan dengan lebih efisien menggunakan metode *<i>machine learning</i>*. *<i>Clustering</i>* merupakan salah satu metode *<i>machine learning</i>* yang bertujuan untuk mengelompokkan data yang memiliki karakteristik serupa ke dalam suatu kelompok/*<i>cluster</i>*. Beberapa contoh metode *<i>clustering</i>* adalah *<i>K-Means</i>*, *<i>Fuzzy C-Means</i>* (FCM), dan *<i>Eigenspace-Based Fuzzy C-Means</i>* (EFCM). Metode *<i>clustering</i>* hanya memproses data numerik, oleh sebab itu diperlukan metode representasi teks. Metode representasi teks yang umum digunakan sebelumnya adalah *<i>Bag of Words</i>* (BoW) dan *<i>Term-Frequency Inversed Document Frequency</i>* (TFIDF). Namun, metode BoW dan TFIDF kurang baik dalam merepresentasikan teks secara kontekstual. Pada tahun 2018 metode representasi teks yang baru ditemukan yaitu metode *<i>Bidirectional Encoder Representation from Transformers</i>* (BERT). Model BERT dapat merepresentasikan teks secara kontekstual dan menghasilkan representasi teks berdimensi tinggi. EFCM merupakan teknik *<i>clustering</i>* yang menggunakan kombinasi teknik reduksi dimensi *<i>Truncated Singular Value Decomposition </i>(TSVD)* dengan teknik *<i>clustering</i>* FCM. Pada tahun 2022 terdapat penelitian yang mengombinasikan BERT dan EFCM untuk pendekslsian topik. Pada model kombinasi BERT dan EFCM terdapat beberapa nilai parameter yang dapat diatur, antara lain adalah pemilihan lapisan *<i>encoder</i>* BERT, dimensi EFCM, dan derajat *<i>fuzziness</i>*. Penelitian ini berfokus pada analisis sensitivitas parameter untuk melihat pengaruh dari nilai parameter terhadap kinerja model EFCM berbasis BERT untuk pendekslsian topik. Analisis sensitivitas parameter menggunakan metode Sobol untuk menentukan parameter yang tidak sensitif dan yang paling sensitif. Kinerja model dievaluasi menggunakan metrik evaluasi *<i>topic coherence</i>*, *<i>topic diversity</i>*, dan *<i>topic quality</i>*. Hasil penelitian menunjukkan bahwa parameter lapisan *<i>encoder</i>*, dimensi EFCM, dan derajat *<i>fuzziness</i>* sensitif terhadap kinerja model. Selain itu, diperoleh model optimal pada tiga *<i>dataset</i>* menggunakan *<i>parameter tuning</i>* metode *<i>grid search</i>*. Penerapan *<i>parameter tuning</i>* dapat meningkatkan performa model pada ketiga *<i>dataset</i>* berdasarkan nilai *<i>topic quality</i>*.

.....

Topic detection is a process to get the subject matter or topic in a text document. In large data, topic detection can be done more efficiently using machine learning methods. Clustering is a machine learning method aiming to group data with similar characteristics into a group/cluster. Some examples of clustering methods are K-Means, Fuzzy C-Means (FCM), and Eigenspace-Based Fuzzy C-Means (EFCM). The clustering method only processes numeric data; therefore, a text representation method is needed. Previously used text representation methods were Bag of Words (BoW) and Term-Frequency Inverse Document Frequency (TFIDF). However, the BoW and TFIDF methods are not good at representing text contextually.

In 2018 a new text representation method was discovered, namely the Bidirectional Encoder Representation from Transformers (BERT) method. The BERT model can contextually represent text and produce high-dimensional text representations. EFCM is a clustering technique that combines the Truncated Singular Value Decomposition (TSVD) dimension reduction technique with the FCM clustering technique. In 2022 there will be research that combines BERT and EFCM for topic detection. In the BERT and EFCM combination model, there are several parameter values that can be set, including the selection of the BERT encoder layer, EFCM dimensions, and the degree of fuzziness. This study focuses on parameter sensitivity analysis to see the effect of parameter values on the performance of the BERT-based EFCM model for topic detection. Parameter sensitivity analysis uses the Sobol method to determine which parameters are insensitive and the most sensitive. Model performance was evaluated using evaluation metrics of topic coherence, topic diversity, and topic quality. The results showed that the parameters of the encoder layer, EFCM dimensions, and degree of fuzziness were sensitive to model performance. In addition, the optimal model was obtained for three datasets using the grid search method parameter tuning. Parameter tuning can improve the model performance on the three datasets based on topic quality values.