

Pendeteksian berita palsu menggunakan RoBERTa dengan Optimalisasi Word Embedding

Adisaputra Nur Arminta¹, Yuliant Sibaroni²

^{1,2}Fakultas Informatika, Universitas Telkom, Bandung

¹adisputraa@student.telkomuniversity.ac.id, ²yuliantsibaroni@telkomuniversity.ac.id

Abstrak

Penyebaran berita palsu (*hoax*) telah menjadi permasalahan serius yang mempengaruhi opini publik dan menciptakan polarisasi di masyarakat. Penelitian ini bertujuan untuk mendeteksi berita palsu menggunakan model *RoBERTa* yang dioptimalkan dengan tiga teknik *word embedding*. *Word embedding* yang digunakan adalah *RoBERTa*, *Word2Vec*, dan *GloVe*. Dataset yang digunakan adalah "Indonesian fact and hoax political news" yang diambil dari Kaggle, Dataset ini memerlukan tahap pre-processing untuk membersihkan ketidakkonsistenan data, seperti mengubah singkatan menjadi kata lengkap dan menghapus tanda baca. Selanjutnya, dilakukan representasi teks menggunakan tiga metode word embedding yaitu *Word2Vec*, *GloVe*, dan *RoBERTa*. Proses pelatihan model dilakukan dengan validasi silang K-Fold untuk meningkatkan generalisasi model. Hasil penelitian menunjukkan bahwa embedding *RoBERTa* mencapai akurasi terbaik 96%, sedangkan word embedding *Word2Vec* mendapatkan akurasi 94%. *Word Embedding GloVe* menunjukkan performa paling rendah dengan akurasi 51%. Penelitian ini membuktikan bahwa pemilihan teknik word embedding yang tidak tepat untuk model *RoBERTa* dapat mengurangi akurasi dan efektivitas model dalam mendeteksi berita palsu. Diharapkan bahwa temuan dalam penelitian ini dapat memberikan kontribusi terhadap peningkatan sistem deteksi berita palsu di masa mendatang.

Kata kunci: *hoax*, *RoBERTa*, *GloVe*, *Word2Vec*

Abstract

The dissemination of misinformation (hoaxes) has emerged as a significant issue, shaping public perception and contributing to societal division. This research aims to detect fake news using the *RoBERTa* model which is optimized with three word embedding techniques. The word embeddings used are *RoBERTa*, *Word2Vec*, and *GloVe*. The dataset used is "Indonesian fact and hoax political news" taken from Kaggle. This dataset requires a pre-processing stage to clean up data inconsistencies, such as changing abbreviations into complete words and removing punctuation. Next, text representation is carried out using three word embedding methods: *Word2Vec*, *GloVe*, and *RoBERTa*. The training phase of the model utilizes the K-Fold cross-validation approach to enhance its ability to generalize across different data distributions. The research results show that *RoBERTa* embedding achieves the best accuracy of 96%, while *Word2Vec* word embedding achieves 94% accuracy. *Word Embedding GloVe* shows the lowest performance with 51% accuracy. This study demonstrates that an unsuitable word embedding method for the *RoBERTa* model can negatively impact both accuracy and efficiency in identifying fake news. The findings from this research are expected to support advancements in future fake news detection systems.

Keywords: *hoax*, *RoBERTa*, *GloVe*, *Word2Vec*

1. Pendahuluan

1.1 Latar Belakang

Hoax berasal dari *hocus* yaitu untuk menipu sering kali muncul pada topik yang sedang hangat dibicarakan. Tujuannya adalah untuk membujuk atau memanipulasi orang dan kemudian melakukan tindakan yang telah ditetapkan sebelumnya, biasanya menggunakan ancaman atau membuat mereka mempercayai hal-hal yang tidak nyata [1]. *Hoax* adalah informasi palsu yang beredar di berbagai platform dan saluran komunikasi, contohnya seperti media sosial dan situs web di internet.

Dengan berkembangnya teknologi digital, pengguna internet dapat dengan mudah membagikan informasi dan berinteraksi sesama pengguna, namun di sisi lain, informasi tersebut juga bisa disalahgunakan untuk menyebarkan berita palsu yang dapat memberikan dampak negatif kepada masyarakat [2]. Oleh karena itu, penyebaran hoaks di dunia maya perlu segera ditangani, karena dapat memicu ketegangan sosial, meningkatkan rasa permusuhan, dan bahkan menyebabkan konflik antar kelompok [3]. Dengan demikian, pengembangan sistem deteksi hoaks otomatis yang mampu mengidentifikasi serta menangani berita palsu secara cepat menjadi hal yang krusial sebelum informasi tersebut menyebar lebih luas.

Dalam upaya mengatasi penyebaran *hoax*, penggunaan salah satu metode yang menjanjikan untuk mengklasifikasi berita *hoax* adalah menggunakan *transformer* [4]. Salah satu metode *transformer* yaitu adalah *RoBERTa* (*Robustly Optimized BERT Approach*), *RoBERTa* adalah modifikasi dari *BERT* yang sederhana namun efektif [19]. Modifikasi tersebut meliputi melatih model lebih lama dengan *batch* yang lebih besar, menggunakan lebih banyak data, menghapus objektif prediksi kalimat berikutnya, melatih pada urutan yang lebih panjang, dan secara dinamis mengubah pola *masking* yang diterapkan pada data latih. Melalui modifikasi ini, *RoBERTa* dapat menghasilkan hasil terbaik pada berbagai tugas pemrosesan bahasa seperti *GLUE*, *RACE*, dan *SQuAD* [19].