

## BAB II

### TINJAUAN PUSTAKA

#### 2.1 Penelitian Terdahulu

Dalam pengembangan metode pemrosesan sistem untuk penelitian ini, beberapa penelitian terkait telah menjadi sumber referensi yang penting. Referensi ini membantu memperkuat dasar teoritis dan metodologis penelitian, serta memberikan wawasan yang berharga dalam merancang pendekatan yang efektif dan inovatif untuk mencapai tujuan penelitian yang ditetapkan. Referensi ini antara lain:

##### 2.1.1 Penelitian Terkait Judi Online

- Penelitian dengan judul “*Detecting Online Gambling Promotions on Indonesian Twitter Using Text Mining Algorithm*” [11] yang dilakukan oleh Reza Bayu Perdana, Ardin, Indra Budi, Aris Budi Santoso, Amanah Ramadiah, dan Prabu Kresna Putra menunjukkan bahwa promosi judi online di indonesia telah meningkat di salah satu platform media sosial, yaitu Twitter. Penelitian juga menunjukkan pentingnya penggunaan algoritma machine learning seperti text mining untuk mendeteksi promosi judi secara otomatis dan diharapkan adanya eksplorasi algoritma deep learning lain, seperti menggabungkan model hybrid untuk meningkatkan performa dan efisiensi. Penelitian juga mengevaluasi berbagai percobaan skenario model, yaitu Random Forest, Logistic Regression, dan Convolutional Neural Networks dalam mendeteksi konten judi. Dataset yang digunakan sebanyak 6038 berbahasa indonesia dengan memanfaatkan kata kunci seperti "slot" "gacor" "untung" dan "cuan". Kata kunci yang sering muncul pada tweet promosi judi meliputi "link," "situs," "prediksi," "jackpot," "maxwin," dan "togel”.
- Penelitian dengan judul “*Perbandingan Model IndoBERT Dan Model Hybrid Pada Analisis Sentimen Opini Masyarakat Terhadap Judi Online*” [12] yang dilakukan oleh Ghufran Afham Asnawi

menunjukkan performa INDOBERT dalam menganalisis sentimen opini masyarakat terhadap judi online di Indonesia. Penelitian memberikan wawasan dan pandangan bahwa menangani judi online di Indonesia dapat ditingkatkan dengan adanya deteksi dan analisa judi online pada media sosial berbasis machine learning. Dataset yang digunakan diambil dari platform media sosial Instagram dengan total 1811 yang menunjukkan komentar mengenai judi online cukup banyak di media sosial. Pengembangan model terdiri dari 3 skenario, yaitu IndoBERT standar, Hybrid IndoBERT-CNN, dan Hybrid IndoBERT-RNN. Hasil evaluasi menunjukkan Model IndoBERT murni mencapai akurasi 60%, model IndoBERT CNN menunjukkan kinerja terbaik dengan akurasi 75%, dan model IndoBERT-RNN dengan akurasi 72%. Ini membuktikan bahwa model hybrid lebih unggul dalam analisis sentimen teks berbahasa Indonesia.

- Penelitian dengan judul “*Penegakan Hukum Perjudian Online di Indonesia, Tantangan dan Solusi*” [27] yang dilakukan oleh Reza Ditya Kesuma menyoroti pentingnya deteksi dan integrasi teknologi AI sebagai bagian dari strategi hukum preventif untuk judi online. Dengan adanya deteksi dini, maka pencegahan judi online oleh anonimitas pelaku dapat lebih efektif.
- Penelitian dengan judul “*Pengaruh Terpaan Iklan Judi Online di Media Sosial terhadap Minat Bermain Judi Online*” [28] yang dilakukan oleh Muhammad Yusuf Akbar menunjukkan bahwa terpaan konten iklan, terutama dalam bentuk teks di media sosial memiliki pengaruh kuat terhadap minat berjudi, khususnya pada pengguna muda usia 17 – 25 tahun. Instagram, Facebook, dan Youtube disebut sebagai platform utama yang digunakan untuk menyebarkan iklan promosi judi, sehingga butuh adanya pemblokiran sistematis terhadap konten judi.

### 2.1.2 Penelitian Terkait BERT

- Penelitian dengan judul “*BERT Models for Arabic Text Classification: A Systematic Review*” [5] yang dilakukan oleh Ali Saleh Alammery membahas tentang penggunaan model algoritma BERT, seperti AraBERT, MARBERT, QARiB dalam mengklasifikasi teks arab. Sebagian besar dataset yang digunakan berasal dari media sosial dengan teks yang relatif pendek. Hasil penelitian menunjukkan bahwa model BERT bersifat multilingual yang mendukung 104 bahasa dan memperoleh nilai F1 yang tinggi, bahkan melebihi kinerja model BERT bahasa Inggris dalam beberapa kasus.
- Penelitian dengan judul “*Discovering a tourism destination with social media data: BERT based sentiment analysis*” [6] yang dilakukan oleh Marlon Santiago, Vinan-Ludena, dan Luis M. de Campos membahas mengenai menentukan destinasi wisata dengan menggunakan teknik sentiment analysis berbasis BERT dengan data dari media sosial Twitter dan Instagram untuk membandingkan beberapa arsitektur seperti stacked BiLSTM, multi-convnets, dan model BERT dengan varian BETO untuk bahasa Spanyol. Penelitian mengumpulkan 90.725 post Instagram dan 235.755 tweet, dengan fokus pada data berbahasa Spanyol dan Inggris. Dilakukan juga tokenisasi, penghapusan stopwords, dan text cleaning. Hasil evaluasi menggunakan metrik akurasi, precision, recall, dan F1 score menunjukkan Model Spanish-BERT memiliki performa unggul jika dibandingkan dengan arsitektur deep learning lain dengan akurasi sekitar 75% untuk bahasa spanyol. Untuk bahasa inggris, model juga memberikan hasil terbaik dengan akurasi dan F1 score di atas 75%.
- Penelitian dengan judul “*Transformer Models for Text-based Emotion Detection: A Review of BERT-based Approaches*” [7] yang dilakukan oleh Francisca Adoma Acheampong, Henry Nunoo-Mensah, dan Wenyu Chen menunjukkan bahwa algoritma BERT

dapat digunakan dalam mendeteksi emosi seseorang dari sebuah teks. Penelitian ini menyoroti pentingnya penggunaan NLP untuk klasifikasi teks dan penggunaan encoder-decoder untuk menguraikan komponen utama dataset. BERT dalam penelitian ini juga menggunakan pendekatan Masked Language Modeling (MLM) dan Next Sentence Prediction (NSP) untuk menangkap konteks dua arah.

- Penelitian dengan judul “*Machine Learning Untuk Deteksi Berita Hoax Menggunakan BERT*” [8] yang dilakukan oleh Isnaeni Imroatus Sholikhah, Aris Tri Jaka Harjanta, dan Khoiriya Latifah menunjukkan penggunaan IndoBERT dengan fine-tuning dalam kasus yang terjadi di Indonesia, yaitu berita hoax di media sosial dan internet. Penelitian menekankan bahwa Pemilihan dataset yang berkualitas menjadi tantangan utama untuk keberhasilan model, sehingga perlu adanya pre-processing dengan tokenisasi untuk menghapus karakter khusus dan angka yang tidak relevan. Hasil evaluasi model IndoBERT menunjukkan akurasi yang cukup yaitu 67% dan dianggap lebih baik dibandingkan model tradisional lainnya. Recall untuk label positif adalah 0.80, yang berarti model berhasil menemukan sekitar 80% dari keseluruhan kasus positif yang ada. F1 Score untuk label positif adalah 0.67, yang merupakan ukuran rata-rata dari precision dan recall. Precision untuk label positif adalah 0.57, yang berarti sekitar 57% dari prediksi positif model benar-benar positif. Peneliti juga menyebutkan bahwa masih terdapat ruang untuk meningkatkan performa model IndoBERT dengan mengoptimalkan parameter atau mengkombinasikannya dengan metode hybrid lainnya.

### 2.1.3 Penelitian Terkait Fuzzy Matching

- Penelitian dengan judul “*Implementasi Fuzzy Search Untuk Pendeteksi Kata Asing Pada Dokumen Microsoft Word*” [9] oleh Ichsan Taufik, Izma Dewi Aishia, dan Jumadi menunjukkan penggunaan *Fuzzy Matching* dalam mendeteksi kata asing agar dalam penulisan karya ilmiah semua kata asing ditulis dalam huruf miring.

Penelitian ini menekankan penggunaan kata kunci dalam bentuk database untuk meningkatkan akurasi dalam mendeteksi kata asing. Perbedaan antara kata yang diketik dan kata dalam database bahasa Indonesia dan mengkategorikannya menjadi Tidak Mirip, Kurang Mirip, Cukup Mirip, Mirip, dan Sangat Mirip. Hasil evaluasi menunjukkan bahwa Fuzzy Search memiliki akurasi rata-rata 89,6% dalam mendeteksi kata asing dan membantu memastikan kepatuhan terhadap kaidah bahasa Indonesia.

- Penelitian dengan judul “*Handwritten Word Recognition Using Fuzzy Matching Degrees*” [10] yang dilakukan oleh Michał Wróbel, Janusz T. Starczewski<sup>1</sup>, Justyna Fijałkowska, Agnieszka Siwocha, dan Christian Napoli menunjukkan adanya penggunaan *Fuzzy Matching* untuk membuat sistem pengenalan tulisan tangan untuk menginterpretasi teks tulisan tangan berdasarkan citra statis dan tantangan seperti variasi gaya tulis dan huruf yang menyambung. Dalam penelitian ini tulisan tangan diuraikan menjadi stroke, yang kemudian diaproksimasi dengan polinomial. Perbandingan antara stroke dari dokumen dengan pola huruf tidak bersifat biner, melainkan diukur dengan fuzzy degrees yang menggambarkan derajat kemiripan bentuk. Hasil eksperimen pada pengenalan kata "tex" menunjukkan bahwa integrasi *fuzzy matching* dengan analisis bigram dapat menghasilkan nilai kecocokan yang membantu memilih hasil pengenalan terbaik.
- Penelitian dengan judul “*A Survey of Text-Matching Techniques*” [25] yang dilakukan oleh Jiang menunjukkan bahwa integrasi BERT dan Fuzzy Matching untuk deteksi spam menghasilkan performa tinggi dalam klasifikasi teks manipulatif. Dataset berupa komentar diambil dari sosial media Reddit yang bertipe spam. Hasil metrik menunjukkan akurasi 0.93, Precision 0.92, Recall 0.91, dan F1-Score 0.915. Kombinasi model menunjukkan kemampuan untuk menangani teks dengan sinonim atau ejaan tidak baku.

## 2.2 Tinjauan Teori

### 2.2.1 *Text Classification*

*Text Classification* adalah salah satu pendekatan NLP yang secara otomatis mengkategorikan teks ke dalam beberapa kategori yang telah ditentukan. Klasifikasi ini memiliki berbagai macam aplikasi seperti, deteksi spam, analisis sentimen, kategorisasi berita, klasifikasi maksud pengguna, moderasi konten, dan sebagainya [5].

Data teks dapat berasal dari berbagai sumber, termasuk data web, email, obrolan, media sosial, tiket, klaim asuransi, ulasan pengguna, dan pertanyaan serta jawaban dari layanan pelanggan, dan masih banyak lagi. Teks adalah sumber informasi yang sangat kaya. Namun, mengekstraksi wawasan dari teks dapat menjadi tantangan dan memakan waktu, karena sifatnya yang tidak terstruktur. Klasifikasi teks dapat dilakukan baik melalui anotasi manual maupun dengan pelabelan otomatis. Dengan semakin besarnya skala data teks dalam aplikasi industri, klasifikasi teks otomatis menjadi semakin penting [13].

### 2.2.2 *NLP*

*Natural Language Processing (NLP)* adalah bidang penelitian yang mengeksplorasi bagaimana komputer dapat digunakan untuk memahami dan memanipulasi teks atau ucapan bahasa alami. Fondasi dari NLP adalah linguistik komputasi, pemodelan berbasis aturan bahasa manusia, dan pemodelan statistik [14].

### 2.2.3 *Fuzzy Matching*

*Fuzzy Matching* atau dalam penelitian ini *Fuzzy String Matching* adalah metode pencarian kata yang menggunakan proses pendekatan terhadap pola dari kata yang dicari. Kunci metode ini adalah bagaimana memutuskan bahwa sebuah kata yang dicari memiliki kesamaan dengan kata yang tertampung, meskipun tidak sama persis dalam susunan karakternya.

Ukuran kuantitatif seperti peresentase untuk mengukur kemiripan antar suatu kata dengan kata yang lain dihitung dengan algoritma Levenshtein Distance atau sering disebut juga algoritma Edit Distance. Algoritma ini menghitung jumlah operasi penghapusan, penyisipan, dan penukaran yang harus dilakukan terhadap suatu kata agar sama dengan kata lain yang menjadi pembanding. Sebagai contoh, kata “komputer” dan “computer” memiliki distance 1 karena hanya perlu dilakukan satu operasi saja untuk mengubah satu karakter “c” menjadi “k” [15].

#### **2.2.4 BERT**

*Bidirectional Encoder Representations (BERT)* adalah salah satu framework NLP yang dikembangkan oleh Google yang menggunakan arsitektur jaringan saraf dalam berdasarkan model transformator canggih. Arsitektur model BERT didasarkan pada jaringan saraf dalam yang disebut transformator. Berbeda dari model NLP tradisional yang memproses teks satu kata dalam satu waktu, BERT dapat memproses seluruh masukan teks sekaligus dan membantu menangkap hubungan antara kata dan frasa dengan lebih efektif. Artinya, BERT dapat mempertimbangkan keseluruhan konteks setiap kata dalam kalimat dan membantu memahami makna teks lebih baik [16]. Salah satu varian lokal dari BERT adalah IndoBERT yang dilatih khusus untuk Bahasa Indonesia. Dengan 3 sumber utama, yaitu Wikipedia Indonesia sebanyak 74 juta kata, artikel Kompas, Tempo, Liputan6 sebanyak 55 juta kata, dan Web Corpus Indonesia sebanyak 90 juta kata [29].

UNIVERSITAS  
MULTIMEDIA  
NUSANTARA