

## BAB 2

### LANDASAN TEORI

#### 2.1 Penelitian Pendahuluan

Beberapa penelitian telah dilakukan untuk mencari dan mengembangkan metode terbaik untuk mendeteksi lirik lagu eksplisit. Setiap metode atau model yang digunakan dalam penelitian tersebut memberikan hasil yang berbeda-beda, masing-masing dengan kelebihan dan kekurangan. Komparasi dari beberapa penelitian ditunjukkan pada Tabel 2.1.

Tabel 2.1. Komparasi Penelitian Terkait Deteksi Lirik Lagu Eksplisit

Judul	Penulis	Metode	Hasil
<i>A novel approach for explicit song lyrics detection using machine and deep ensemble learning models</i>	Chen et al. (2023)	<i>Soft voting</i> menggabungkan <i>Extra Trees Classifier (ETC)</i> dan <i>LSTM</i> .	Model mencapai akurasi 96,92%, <i>precision</i> 100,00%, <i>recall</i> 92,97%, dan <i>F1-score</i> 96,79%.
<i>Explicit content detection in music lyrics using machine learning</i>	Chin et al. (2018)	<i>Bagging</i> dengan kamus kata eksplisit.	Model menghasilkan <i>precision</i> 0,96, <i>recall</i> 0,97, dan <i>F1-score</i> 0,96.

Penelitian oleh Chen et al. (2023) menggunakan pendekatan *deep ensemble* dengan menggabungkan model berbasis pohon keputusan dan jaringan saraf untuk meningkatkan performa deteksi konten eksplisit dalam lirik lagu. Sementara itu, Chin et al. (2018) memanfaatkan pendekatan berbasis *bagging* dan kamus kata kasar yang terbukti efektif dalam klasifikasi kata eksplisit. Meskipun pendekatan Chin lebih sederhana, hasil akurasi dan metrik evaluasi lainnya tetap tinggi. Penelitian-penelitian ini menjadi dasar yang kuat bagi pengembangan model klasifikasi eksplisit berbasis *voting ensemble* dalam penelitian ini.

Sementara itu, Chin et al. (2018) mengusulkan metode klasifikasi berdasarkan teknik *bagging* dan penggunaan kamus kata eksplisit untuk menandai kata-kata bermuatan seksual atau kasar dalam lirik. Meskipun tidak menggunakan teknik *deep learning*, penelitian ini tetap memberikan hasil yang sangat baik dengan *f1-score* sebesar 0,96. Hal ini menunjukkan bahwa teknik berbasis kamus dan *ensemble learning* sederhana masih relevan untuk digunakan dalam klasifikasi konten eksplisit, terutama untuk dataset yang lebih terbatas atau sistem dengan kebutuhan komputasi rendah.

## 2.2 Lirik Lagu

Lagu merupakan kumpulan alunan irama yang dinyanyikan oleh musisi. Lirik lagu merupakan karya cipta musisi yang mengandung ungkapan perasaan penulis lagu yang memiliki bentuk puisi pendek. Lirik lagu berisi emosi dan ekspresi penulis lagu yang disampaikan kepada masyarakat. Lirik lagu sendiri digunakan untuk menyampaikan isi pikiran dan perasaan penulis lagu untuk pendengar lagu [21].

Lirik lagu merupakan karya seni seorang pengarang dimana ide, perasaan, dan kreativitas dipadukan sehingga menjadi harmoni yang dapat dinikmati oleh pendengar. Lirik lagu dapat menjadi karya sastra puisi dikarenakan lirik lagu dan puisi memiliki kesamaan sebagai media untuk mengungkapkan pikiran dan perasaan seseorang dengan rima, irama, serta harmonisasi[22].

## 2.3 Kata Seksual

Kata seksual atau eksplisit termasuk ke dalam kategori *offensive language* atau *explicit content*, yang mencakup kata-kata dengan muatan seksual, kekerasan, diskriminatif, atau vulgar. Pendeteksian kata-kata seperti ini sangat penting dalam konteks penyiaran, pendidikan, dan platform digital [13]. Daftar kata seksual umumnya disusun dalam bentuk kamus (*lexicon-based*), namun sering kali tidak lengkap dan tidak kontekstual, sehingga pendekatan berbasis pembelajaran mesin diperlukan.

## 2.4 Preprocessing

Pada tahap *preprocessing*, data teks diolah agar lebih terstruktur dan siap diproses lebih lanjut. Ada beberapa hal yang dilakukan pada tahap *preprocessing*,

yaitu *data cleaning*, *tokenizing*, *stopwords removal*, dan *stemming*.

#### 1. *Data Cleaning*

Pada proses *data cleaning*, data teks akan menghapus data tidak relevan seperti angka, simbol, dan karakter. Proses ini digunakan untuk menghapus data yang tidak diperlukan serta mengurangi waktu proses penelitian karena tidak perlu membaca data yang tidak penting.

#### 2. *Tokenizing*

Pada proses *tokenizing*, teks akan dipecah menjadi bagian-bagian yang disebut token. Proses ini merupakan dasar bagi banyak metode dalam analisis teks.

#### 3. *Stopwords Removal*

*Stopwords* merupakan kata yang tidak memiliki arti. Pada proses *stopwords removal*, kata-kata yang tidak memiliki arti akan dihapus. Kata-kata yang kurang penting akan disaring sehingga meninggalkan dataset dengan data yang bersih.

#### 4. *Stemming*

Proses *stemming* merupakan proses untuk mengubah kata ke kata dasar. Tujuan dari *stemming* adalah untuk menyamakan kata-kata yang memiliki makna serupa dalam bentuk berbeda.

### 2.5 Feature Extraction

*Feature extraction* adalah proses mengubah teks mentah menjadi representasi numerik yang dapat digunakan oleh model pembelajaran mesin. Pada penelitian ini, digunakan tiga teknik ekstraksi fitur:

1. **TF-IDF (Term Frequency-Inverse Document Frequency)** digunakan untuk mengetahui seberapa penting suatu kata dalam dokumen. *TF* digunakan dalam mengukur kata muncul dalam dokumen, semakin sering kata muncul maka nilai *TF* akan semakin tinggi. *IDF* merupakan pengukur kata yang dipengaruhi oleh seberapa umum atau langka kata dalam dokumen. Kombinasi *TF* dan *IDF* akan menghasilkan vektor numerik yang menjadi representasi kata dalam dokumen.

2. **Word2Vec** merupakan *feature extraction* yang bekerja dengan menggunakan *Continuous Bag of Words (CBOW)* dan *Skip-gram*. *CBOW* berfungsi dalam memprediksi kata target dengan kata konteks di sekitar. *Skip-gram* berfungsi dalam memprediksi kata konteks dengan kata target.
3. **Doc2Vec** adalah *feature extraction* yang bekerja dengan menggunakan *Distributed Memory (DM)* dan *Distributed Bag of Words (DBOW)*. *DM* berfungsi dalam mempertimbangkan konteks dari kata dan ID dokumen sebagai input prediksi kata selanjutnya. *DBOW* berfungsi untuk memprediksi kata dari dokumen dengan menggunakan vektor dokumen sebagai konteks.

## 2.6 K-Fold Cross Validation

*K-Fold Cross Validation* merupakan metode validasi dengan membagi *dataset* menjadi bagian (*fold*) dengan ukuran yang sama. Proses pelatihan model dilakukan sebanyak jumlah *fold* yang dipilih, di mana dalam setiap iterasi, satu *fold* digunakan sebagai data uji (*testing set*) dan sisanya digunakan sebagai data latih (*training set*). Hasil dari evaluasi setiap *fold* akan dirata-rata untuk mendapatkan estimasi performa model yang lebih stabil dan generalisasi yang baik [23]. Kelebihan dari metode ini adalah dalam memaksimalkan data untuk pelatihan dan pengujian, serta mengurangi kemungkinan hasil evaluasi bias terhadap data.

## 2.7 Voting Classifier

*Voting Classifier* merupakan metode *ensemble learning* yang menggabungkan prediksi dari beberapa model untuk menghasilkan klasifikasi akhir yang lebih akurat dan stabil. *Voting classifier* berfungsi untuk meningkatkan tingkat akurasi dengan menggabungkan hasil dari model yang digabungkan [24]. Terdapat dua jenis utama *voting classifier*.

1. *Hard Voting* merupakan metode mengambil mayoritas hasil prediksi (klasifikasi) dari seluruh model.
2. *Soft Voting* merupakan metode mengambil rata-rata dari probabilitas prediksi masing-masing model, dan memilih kelas dengan nilai probabilitas tertinggi [25].

Metode ini sangat efektif apabila model-model yang digabungkan memiliki karakteristik dan kekuatan yang berbeda, sehingga dapat saling melengkapi [19].

## 2.8 LSTM

*LSTM* merupakan salah satu arsitektur dari *Recurrent Neural Network (RNN)* yang dirancang untuk mengatasi masalah dalam pemrosesan data sekuensial. *LSTM* memiliki kemampuan untuk mengatasi *long-term dependency* dimana model perlu untuk mengingat informasi dengan langkah jauh dalam sebuah urutan [17].

## 2.9 Extra Trees Classifier

*ETC* atau *Extremely Randomized Trees* adalah metode *ensemble* dengan basis *decision tree* yang memiliki tujuan untuk meningkatkan akurasi klasifikasi dengan cara membuat *decision tree* dan menggabungkan prediksi dari *decision tree* tersebut. *ETC* sendiri merupakan variasi dari algoritma *Random Forest* namun memiliki metode pembentukan *tree* dan *split tree* yang berbeda sehingga menjadi lebih cepat dan lebih acak [16].

## 2.10 Metrik Evaluasi

*Metrik evaluasi* merupakan metrik yang dapat digunakan sebagai indikator performa dari analisis yang dijalankan. *Metrik evaluasi* merupakan hasil dari *confusion matrix*. *Confusion matrix* merupakan matriks perbandingan nilai aktual dan nilai prediksi dari model analisis yang dibuat.

*Confusion matrix*, memiliki nilai *true positive (TP)*, *true negative (TN)*, *false positive (FP)*, *false negative (FN)*. Dari nilai *confusion matrix*, maka hasil *accuracy*, *precision*, *recall*, *f1-score* dari modelan dapat diketahui.

### 1. Accuracy

$$\text{Accuracy} = \frac{TP + TN}{TP + TN + FP + FN} \quad (2.1)$$

### 2. Precision

$$\text{Precision} = \frac{TP}{TP + FP} \quad (2.2)$$

3. *Recall*

$$\text{Recall} = \frac{TP}{TP + FN} \quad (2.3)$$

4. *F1-Score*

$$F1 = \frac{2 \cdot \text{Precision} \cdot \text{Recall}}{\text{Precision} + \text{Recall}} \quad (2.4)$$

