

## BAB III

### METODOLOGI PENELITIAN

#### 3.1 Gambaran Umum Objek Penelitian

Objek penelitian yang digunakan dalam penelitian ini adalah teks yang mengandung kekerasan seksual verbal di media sosial yang menggunakan keyword berupa lonte” [19], “digilir”, “janda”, “cabe-cabean”, dan jablay” [20], yang dimana akan dipakai untuk melakukan klasifikasi teks di media sosial. Klasifikasi teks yang dilakukan adalah klasifikasi terhadap teks yang diduga mengandung konten berupa kekerasan seksual. Sumber data untuk permodelan akan menggunakan data yang dikumpulkan dengan menggunakan sebuah platform media sosial yang bernama *Twitter*. Peneliti akan menggunakan *Twitter API* dengan tujuan untuk mengumpulkan data tweet yang dianggap mengandung konten berupa kekerasan seksual. Selanjutnya, model akan dibuat dengan menggunakan algoritma *AWD-LSTM* yang akan dijadikan referensi untuk memprediksi data baru yang didapatkan dari berbagai *platform* media sosial, yaitu berupa data teks yang diduga mengandung kekerasan seksual verbal. Data baru tersebut akan dibandingkan dengan data yang telah didapatkan dari model yang telah dibuat. Perbandingan tersebut bertujuan bertujuan untuk mengetahui apakah data baru tersebut memiliki konten yang berkaitan dengan kekerasan seksual.

#### 3.2 Metode Penelitian

Dalam metode penelitian, terdapat 3 arsitektur untuk *Data Mining* yang bernama KDD, CRISP-DM, dan SEMMA, yang dimana dapat dijadikan referensi untuk dijadikan sebagai alur penelitian [21]. Berikut merupakan tabel perbandingan dari ketiga arsitektur tersebut.

Tabel 3.1. Perbandingan Arsitektur Data Mining [21]

Model dari Proses Data Mining	KDD	CRISP-DM	SEMMA
Jumlah Langkah dari Model	9	6	3

Model dari Proses Data Mining	KDD	CRISP-DM	SEMMA
Nama Langkah dari Model	Developing and Understanding of the Application	Business Understanding	-----
	Creating a Target Data Set	Data Understanding	Sample
	Data Cleaning and Pre-processing		Explore
	Data Transformation	Data Preparation	Modify
	Choosing the suitable Data Mining Task	Modeling	Model
	Choosing the suitable Data Mining Algorithm		
	Employing Data Mining Algorithm		
	Interpreting Mined Patterns	Evaluation	Assessment
	Using Discovered Knowledge	Deployment	-----
Orientasi Pemakaian	Research Oriented	Company Oriented	Company Oriented

Berdasarkan tabel 3.1, penelitian ini akan memakai arsitektur *Knowledge Discovery in Databases* sebagai alur penelitian. Arsitektur tersebut dipilih sebagai alur penelitian dikarenakan metode tersebut memiliki alur proses yang lebih lengkap dan akurat dari sisi modelling dibandingkan dengan arsitektur lainnya. Dapat dilihat pada tabel 3.1, metode KDD memiliki 3 langkah yang lebih spesifik dalam melakukan permodelan dibandingkan dengan metode lainnya dan metode KDD lebih berorientasi terhadap penelitian. Oleh karena hal tersebut, arsitektur KDD lebih cocok dipakai untuk penelitian dibandingkan dengan arsitektur CRISP-DM dan SEMMA yang dimana lebih berorientasi terhadap penerapan di perusahaan [21].

### 3.2.1 Alur Penelitian

Pada penelitian ini, peneliti akan menggunakan metode *Knowledge Discovery in Databases* (KDD) yang terdiri dari :

## 1. Understanding Goal

Kasus kekerasan seksual berbasis gender online mengalami loncatan kasus mengalami pelonjakan angka yang cukup drastis. Hal tersebut dapat dilihat dari jumlah kasus sebanyak 241 kasus pada tahun 2019 dan jumlah kasus sebanyak 940 kasus yang terjadi saat tahun 2020 semenjak munculnya pandemi Covid-19 terutama kasus yang terjadi pada media sosial [6]. Dalam menjalankan penanggulangan terhadap kekerasan seksual, terdapat layanan – layanan yang disediakan pemerintah yang dimana masih kurang cukup untuk mengatasi jumlah peningkatan kekerasan seksual yang terjadi. Dengan meningkatnya kasus kekerasan seksual yang terjadi, layanan tersebut akan kewalahan sehingga tidak akan dapat menangani semua kasus kekerasan seksual yang terjadi. Oleh karena itu, penelitian ini dilakukan untuk menyelesaikan suatu masalah, yaitu minimnya sebuah sistem yang dapat dijadikan referensi untuk dapat mendeteksi kekerasan seksual verbal yang terjadi. Penelitian ini dibuat untuk dapat mendeteksi kekerasan seksual verbal yang terjadi di media sosial dan juga dapat dijadikan referensi untuk pembuatan suatu sistem yang dapat mendeteksi kekerasan seksual verbal yang terjadi di media sosial.

## 2. Selection

Sumber data yang akan digunakan pada penelitian ini bersumber dari data yang akan diambil langsung dengan menggunakan Twitter API. Data yang diperoleh merupakan data tweet yang memiliki keyword berupa keyword berupa lonte” [19], “digilir”, “janda”, “cabe-cabean”, dan jablay” [20] yang dimana diduga mengandung konten berupa kekerasan seksual. Pengambilan data sendiri dilakukan mulai dari tanggal 29 Maret 2022 sampai 5 April 2022. Setelah itu, data yang sudah dikumpulkan akan disaring oleh peneliti untuk diberikan label oleh orang lain yang dimana akan menghasilkan kelas akhir yang berupa “Yes” atau “No” yang dimana “Yes” melambangkan bahwa *tweet* tersebut mengandung teks yang terindikasi sebagai kekerasan seksual verbal, dan “No” melambangkan bahwa *tweet* tersebut tidak terindikasi sebagai kekerasan seksual verbal. Penyaringan data yang dilakukan oleh peneliti didasarkan pada kemampuan

berbahasa Indonesia yang dimana memiliki nilai minimal “B” dalam pelajaran Bahasa Indonesia di universitas peneliti. Data tersebut akan diberikan untuk diberikan label dengan aturan yang mengharuskan minimal tiga orang untuk memberikan label yang dimana orang ketiga akan menjadi penengah apabila dua orang pertama memiliki pendapat yang berbeda dan label kelas akhir akan ditentukan berdasarkan hasil modus dari hasil labelling yang dilakukan oleh ketiga orang tersebut.

### 3. Data Cleaning and Preprocessing

Pada tahapan *data cleaning and preprocessing*, akan dilakukan proses berupa pembersihan data yang terdiri dari penghapusan *noise*, penanganan terhadap data yang hilang, mengurutkan data berdasarkan waktu input, dan melihat apakah terdapat perubahan signifikan terhadap data yang akan dipakai. Tahapan ini akan melakukan beberapa proses seperti *remove duplicate data* yang dimana merupakan proses untuk menghilangkan data untuk menghindari bias terhadap suatu sentimen, *case folding* merupakan proses untuk mengubah kata menjadi huruf kecil, *tokenize* untuk pemisahan teks menjadi potongan-potongan kecil yang disebut sebagai token, *filtering* untuk mengurangi *noise* dengan cara menghapus atribut yang dapat mempengaruhi model, *stopwords removal* untuk menghapus kata yang tidak memberikan informasi penting untuk permodelan, dan *stemming and lemmatization* yang digunakan untuk proses normalisasi kata yang sama akan tetapi berbeda imbuhanannya [15].

### 4. Transformation

Pada tahapan *transformation*, akan dilakukan proses berupa pemilihan fitur yang sesuai dengan data yang akan digunakan. Selain itu, akan dilakukan reduksi terhadap dimensi dan filtrasi yang dimana akan mempersiapkan data tersebut untuk permodelan. Tahapan ini akan melakukan proses seperti *extraction and parsing* yang berfungsi untuk menyesuaikan bentuk dan format data untuk permodelan yang dibuat, *translation and mapping* yang berfungsi untuk melakukan pemetaan sehingga data yang ada lebih mudah untuk dipahami oleh

model dan terakhir adalah *indexing and ordering* yang merupakan proses pengaturan data yang diurutkan agar dapat menghasilkan model yang baik [22].

#### 5. Data Mining Task

Pada penelitian ini metode yang akan dipakai adalah metode klasifikasi. Metode klasifikasi merupakan sebuah metode yang melakukan prediksi untuk mengetahui penetapan suatu label terhadap kelas yang ada yang didasarkan pada pelatihan model untuk memprediksi data baru [23]. Metode klasifikasi akan dipakai dengan tujuan untuk mengetahui hasil label terakhir yang berupa “Yes” atau “No” yang dimana “Yes” melambangkan bahwa *tweet* tersebut mengandung teks yang terindikasi sebagai kekerasan seksual verbal, dan “No” melambangkan bahwa *tweet* tersebut tidak terindikasi sebagai kekerasan seksual verbal

#### 6. Data Mining Algorithm Selection

Dalam penelitian ini, algoritma yang akan digunakan adalah algoritma *Average Stochastic Gradient Descent Weight-Dropped LSTM* (AWD-LSTM) yang didasarkan pada penelitian terdahulu [7]. Algoritma tersebut dipilih karena menghasilkan model dengan nilai akurasi precision, recall, dan f1-score sebesar : 96%, 95%, 97%, dan 96% dan juga karena kemampuannya dalam menangkap ketergantungan kata kompleks dengan lebih baik apabila dibandingkan dengan algoritma lain yang hanya berfokus untuk mendapatkan nilai akurasi, presisi, recall, dan F1-Score tertinggi.

#### 7. Data Mining

Pada tahapan *data mining*, data akan dibagi menjadi 70% *data train* dan 30% *data test* yang didasarkan pada penelitian terdahulu [7]. Selanjutnya, model akan diberikan hyperparameter yang berbeda yang akan dipakai dalam pelatihan model (AWD-LSTM). Setelah itu, akan dilakukan proses pelatihan yang dimana akan dilakukan pemantauan terhadap nilai *loss* dan nilai akurasi dari data *training* dan validasi terhadap setiap epoch. Proses tersebut akan dijalankan dengan menggunakan *Jupyter Notebook* yang menggunakan sumber daya komputasi

dengan spesifikasi berupa: Processor Intel(R) Core (TM) i7-4710HQ, RAM sebesar 8 GB, NVIDIA GPU dengan VRAM sebesar 2 GB

## 8. Interpretation

Evaluasi model dari klasifikasi akan dilakukan terhadap data *testing* dengan menggunakan metrik berupa: akurasi, precision, recall, dan skor F1.

$$Akurasi = \frac{Label\ Benar}{Total\ Label}$$

**Rumus 3.1 Rumus Perhitungan Akurasi [24]**

$$Precision = \frac{True\ Positive}{True\ Positive + False\ Positive}$$

**Rumus 3.2 Rumus Perhitungan Precision [24]**

$$Recall = \frac{True\ Positive}{True\ Positive + False\ Negative}$$

**Rumus 3.3 Rumus Perhitungan Recall [24]**

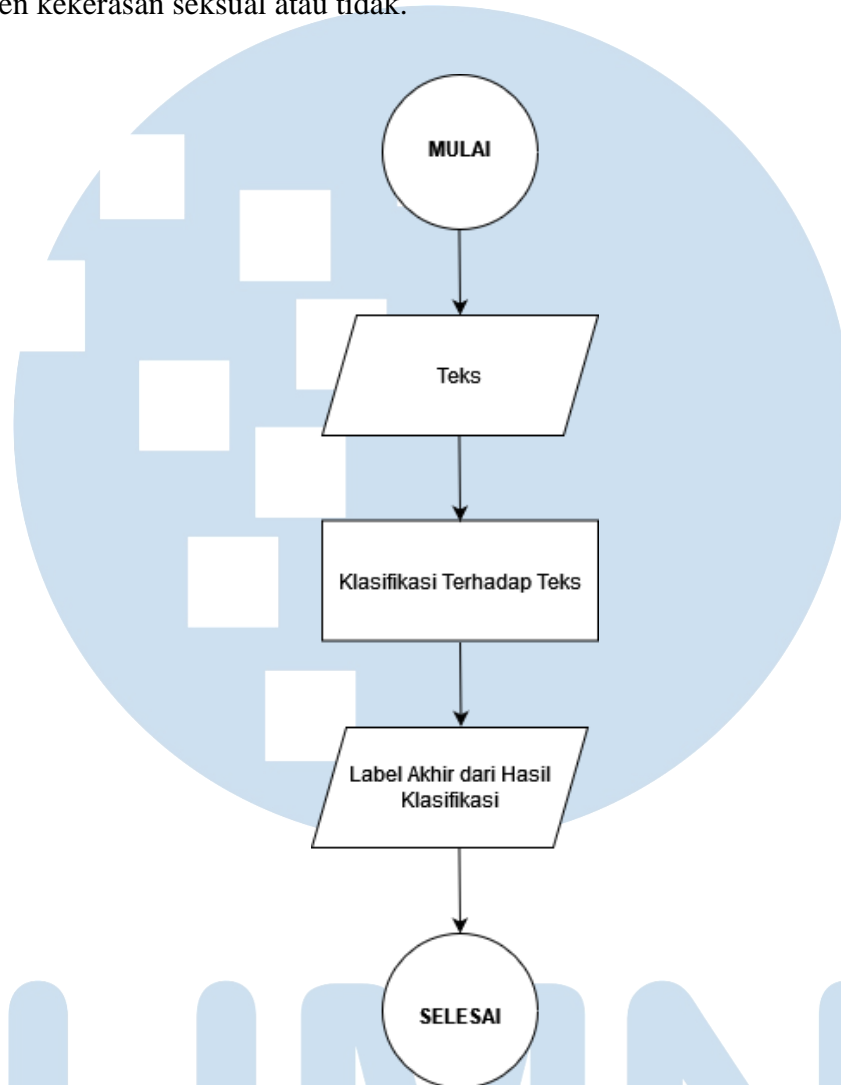
$$F1 = 2 \times \frac{Precision \times Recall}{Precision + Recall}$$

**Rumus 3.4 Rumus Perhitungan F1-Score [24]**

## 9. Consolidation

Pada tahapan *consolidation*, model yang telah dibuat akan diimplementasikan menjadi sebuah web sederhana. Web yang dibuat tersebut akan menggunakan framework yang bernama Flask yang dimana akan memungkinkan user untuk menulis teks yang dianggap mengandung konten kekerasan seksual verbal. Setelah mengunggah teks tersebut, *API* dari Flask akan melakukan ekstraksi terhadap pesan yang diunggah sehingga akan didapatkan teks sebagai input dari model klasifikasi yang sudah dibangun. Output dari model

tersebut akan menentukan apakah teks yang diunggah tersebut mengandung konten kekerasan seksual atau tidak.



**Gambar 3.1. Flowchart Implementasi Cara Kerja dari Model Web**

### 3.3 Variabel Penelitian

#### 3.3.1 Variabel Independen

Variabel Independen adalah variabel yang berpengaruh terhadap terjadinya variabel dependen [25]. Dalam penelitian ini, variabel independen yang digunakan adalah *tweet* yang memiliki keyword yang diduga mengandung konten berupa kekerasan seksual berada di media sosial Twitter.



### 3.3.2 Variabel Dependen

Variabel Dependen adalah variabel yang dipengaruhi oleh variabel Independen [25]. Dalam penelitian ini, variabel dependen yang digunakan adalah klasifikasi kelas yang dipakai untuk mengetahui apakah tweet mengandung kekerasan seksual secara verbal. Kelas yang akan dipakai dalam penelitian ini adalah tweet yang mengandung kekerasan seksual dan tweet yang tidak mengandung kekerasan seksual

### 3.4 Teknik Pengumpulan Data

Data yang dipakai merupakan data primer yang diambil langsung dengan menggunakan Twitter *API* untuk mengumpulkan data tweet dari media sosial Twitter. Periode pengumpulan data tweet akan dilakukan selama 7 hari, dimulai dari tanggal 29 Maret 2022 sampai 5 April 2022. Pengumpulan data tersebut akan didasarkan pada keyword yang didasarkan pada penelitian terdahulu mengenai kekerasan seksual verbal. Keyword tersebut terdiri dari : “lonte” [19], “digilir”, “janda”, “cabe-cabe”, dan jablay” [20]. Untuk jumlah data yang dikumpulkan akan menggunakan referensi dari penelitian tentang analisis sentimen tentang pilkada Jawa Barat dengan menggunakan data yang dikumpulkan dari Twitter. Penelitian tersebut menggunakan sebanyak 300 data tweet yang akan dibagi menjadi 2, yaitu 200 *data train*, dan 100 *data test* untuk pelatihan model yang dimana menghasilkan akurasi yang cukup baik, yaitu sebesar 84% [26].

### 3.5 Teknik Pengambilan Sampel

Teknik pengambilan sampel yang akan dipakai dalam penelitian ini adalah teknik *purposive sampling*. Teknik tersebut akan menetapkan pertimbangan atau kriteria tertentu sehingga sampel data yang diambil sesuai dengan kriteria dari penelitian yang dibuat [25]. Dalam penelitian ini, sampel data yang terpilih adalah data tweet yang memiliki keyword yang mengandung konten berupa kekerasan seksual, seperti : komentar seksual terhadap tubuh seseorang dan lelucon seksual



untuk merendahkan fisik orang lain [19]. Keyword yang akan dijadikan referensi terdiri dari “lonte” [19], “digilir”, “janda”, “cabe-cabean”, dan jablay” [20].

### 3.6 Teknik Analisis Data

Dalam analisis data, terdapat 2 arsitektur untuk *algoritma* yang bernama *Average Stochastic Gradient Descent Weight-Dropped LSTMs* (AWD-LSTMs), *AWD-Quasi-Recurrent Neural Networks* (QRNNs) yang dimana akan dibandingkan untuk dijadikan referensi algoritma dalam pembuatan model penelitian. Berikut merupakan tabel perbandingan dari kedua arsitektur tersebut.

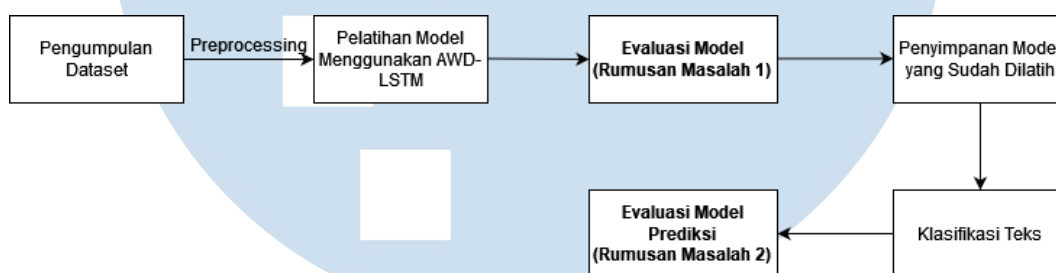
**Tabel 3.2. Perbandingan Arsitektur Algoritma [27]**

<b>DNN Architecture</b>	<b>Epochs</b>	<b>Accuracy</b>	<b>Train Loss</b>	<b>Validation Loss</b>	<b>Pre-Trained?</b>
AWD-LSTM word	31	0.494893	2.308937	2.341698	Yes
AWD-LSTM unigram	31	0.557226	1.639875	1.826841	Yes
AWD-LSTM BPE	31	0.580373	1.561393	1.703536	Yes
AWD-LSTM char	31	0.779633	0.757956	0.742808	Yes
AWD-QRNN word	31	0.515747	1.972508	2.144126	No
AWD-QRNN unigram	31	0.539951	1.790150	1.894901	No
AWD-QRNN BPE	31	0.538290	1.824709	1.896698	No
AWD-QRNN char	31	0.736358	0.944526	0.897850	No

Berdasarkan tabel 3.2, penelitian ini akan memakai *Average Stochastic Gradient Descent Weight-Dropped LSTM* (AWD-LSTM) sebagai model untuk penelitian. Arsitektur tersebut dipilih sebagai model penelitian dikarenakan memiliki akurasi rata – rata yang lebih baik dalam melakukan tokenisasi terhadap tokenizer yang berupa *word, unigram, Byte Pair Encoding* (BPE), dan *character* (Char) [27].

### 3.7 Kerangka Teori

Kerangka kerja pada penelitian dapat dilihat dari gambar 3.2. di bawah ini.



**Gambar 3.1. Kerangka Teori**

#### 3.7.1 Rumusan Masalah 1

Dataset yang sudah dikumpulkan akan dilanjutkan ke tahapan *preprocessing* yang kemudian dijadikan referensi untuk membuat permodelan dengan menggunakan algoritma *AWD-LSTM*. Data yang dilatih dengan algoritma *AWD-LSTM* memiliki tugas yang bertujuan untuk melakukan klasifikasi terhadap teks yang memiliki 2 kelas yang terdiri dari teks yang dianggap memiliki konten berupa kekerasan seksual dan teks yang dianggap tidak memiliki konten kekerasan seksual. Model yang sudah dilatih yang dimana menghasilkan performa terbaik akan dievaluasi kembali dan akan disimpan untuk dijadikan acuan untuk tahapan deployment

### 3.7.2 Rumusan Masalah 2

Penelitian ini akan mengimplementasikan model yang telah dibuat ke dalam suatu sistem yang memiliki basis web. Pembuatan model berbasis web memiliki tujuan untuk menguji model dalam memprediksi kata yang akan dijadikan data baru yang dimana dapat diperoleh dari inputan user yang memakai sistem tersebut.

Model yang dibuat akan diuji dengan menggunakan teks yang memiliki konten berupa kata yang diduga terkait dengan kekerasan seksual dan kata yang diduga tidak terkait dengan kekerasan seksual. Model yang berbasis web tersebut akan diuji untuk dapat mendeteksi apakah model sudah dapat melakukan labelling data dengan baik.

A large, light blue watermark logo of Universitas Multimedia Nusantara (UMMN) is centered on the page. It features a stylized 'U' and 'M' inside a circle, with 'N' to the right.

UMMN

UNIVERSITAS  
MULTIMEDIA  
NUSANTARA