

## BAB II TINJAUAN PUSTAKA

### 2.1 Kalimat *Toxic*

Kalimat *toxic* dari segi komunikasi dapat dikatakan sebagai pola komunikasi yang merugikan dan tidak sehat antara individu. Teori ini menekankan pentingnya komunikasi yang efektif, menghormati, dan mendukung dalam hubungan antarmanusia (Arwan, 2018). Dimana Kalimat *toxic* cenderung memiliki efek emosional negatif pada penerima pesan atau orang yang dituju. Paparan berulang terhadap kalimat *toxic* dapat menyebabkan penurunan harga diri, depresi, kecemasan, dan stres kronis. Ini dapat mengganggu hubungan interpersonal, kinerja kerja, dan kesejahteraan secara keseluruhan (Thahir, 2014). Kalimat yang merendahkan dapat menciptakan persepsi negatif tentang kemampuan diri dan kepercayaan diri (Ulfah & Winata, 2021). Beberapa contoh kalimat atau kata *toxic* yang dipaparkan oleh (Alika dkk., 2022) dalam penelitian yang telah dilakukannya tentang “Urgensi Penggunaan Tata Bahasa yang Baik dalam Berkomentar di Aplikasi Media Sosial *Tiktok* Terhadap Kesehatan Mental dan Pembentukan Karakter pada Siswa SMP dan SMA” dan penelitian yang dilakukan oleh (Robbani dkk., 2023) yang melakukan penelitian tentang “Analisis Wacana Digital Penggunaan *Verbal Abuse* dalam Konten *Gameplay Stumble Guys* Luthfi Halimawan” dapat dilihat pada tabel 2.1 berikut:

**Tabel 2. 1 Contoh Kalimat *Toxic***

<b>Contoh Kalimat <i>Toxic</i></b>	
anak anjing	tolol
anjing	lonte
bangsat	kayak lonte

<b>Contoh Kalimat Toxic</b>	
anjir	alay
jancok	kampungan prik/bocah prik
ngentot	bego
babi	bacot
pantek	gak jelas
brengsek	lebay
kontol	lo ngapain? Keren kaga jelek iya
konthol	gaje lu
koncol	gendut
memek	kaya tante tante
bodoh	jelek
goblok	kurus banget
Oh, ini yang Winnery-Winnery nih ngentot nih. Eh, gua mau menang lagi lawan bot anjing. Jangan dipeluk! Goblok! Jangan peluk gua, anjing! Winnery ngentot gua tandain lu, goblok! Anjing, Winnery. Jangan peluk gua. Jangan peluk, gua lagi lawan bot anjing! Winnery goblok gue pengen menang sekali seumur hidup anjing. ... Aaa! Memek! Awas! Awas ngentot!	Tai ya, Winnery ya, menang lagi dia, bangsat anjing. Gua doain lu dipanggil ibu lu anjing. Awas goblok! Winnery, tolol! Anjing Winnery ni ngentot. Gue doain lu gapunya temen di sekolah, bangsat. Gausah peluk-peluk, anjing! Goblok, tengkorak merah, tolol! Diem dulu, diem

<b>Contoh Kalimat Toxic</b>	
<p>Ga usah sliding-sliding memek! Awas anjing! Goblok, bot kontol! Anjing ni orang sumpah... tai anjing! Goblok! Game Stumble, game beta memek! Anying</p>	<p>dulu, ini, ini memang ngentot yak. Ini isinya tu anak kon... ehh. Anjing, kesel gua. ... Nih Winnery aja kek anjing hidupnya ni. Winnery-Winnery kentot, anjing. ... Eh mati kan lu Winnery, goblok! Aaa! Winnery Goblok!!!</p>
<p>Liat, kan. Game beta, goblok. Benerin game lu, goblok! Memek, anjing, kesel banget gua, ngentot! Sss.. Aaa!!! Gamau, keluar, anjing, memek, tolol! Blok! Anjing. Kesel bet gua, anjiing! Tau ni, nulis review lagi ini ... kontol, ya. Gue nulis (review game) apa ya, kemarin ya? 'Memek Goldplay Button' anjing, gue nulis.</p>	<p>Gausah nempel-nempel ya, anjing! Gausah nempel-nempel! Goblok! Semua aja meluk gua tolol! Kuntoo!!! Ayam ngentot! Weh, udah dong. Anjing, kasih sekali gua menang apa, kontol! Weh, kasih guese kali menang apa, jangan kaya anjingg! ... Lu liat nih, spasinya lepas, anjing! Lu Tempel di keyboard nih, tempel ni keyboard sendiri! Nih! Lu tempel ni keyboard, goblok!</p>
<p>What the fuck, dude. Anjing, udah gitu dia (Istri Luthfi) ga main lagi, ngentot! Sumpah ya, ini yang main peluk-peluk, dua doain pas BO (menyewa pekerja seks komersial)dapet memeknya HIV anjing! Pas lu BO, memek BO-an lu berlalat kontol!</p>	<p>.. A[L] Cupuk mem... wah ngentot, anjing! Fuck, anjing! A[L] Cupuk! Sumpah jangan kaya kontol! Sumpah jangan kaya kontol! Please ini jangan kaya kontol, A[L] Cupuk, anjing! Sumpah jangan... A[L] Cupuk ... kaya kontol! Huw! Let's go! Let's fucking go! ...Let's go, sir!</p>

<b>Contoh Kalimat Toxic</b>	
	<p>Didzolimi ngentot, gua ... Mampus A[L] Cupuk! Didzolimi, diperkosa, dikentot rame-rame tapi aku tetap bangkit, anying</p>
<p>angaaan. Undang aja. Eh, gue tuh, orangnya, pengennya bermain bersama kalian guys. Karena, karena menurut gua bermain tanpa kalian tuh, ga seru. ... Gua tu orangnya, solid, bro, sama kalian bro. Ayo buat guys, ayo, masuk-masuk. Masuk, temen-temen. ... Nah ayo guys, 96264 (kode custom room), ayo temen-temen, masuk semuanya! Masuk guys! Gue orangnya solidaritas banget bro, udah gausah... anjing. Udah, bodo amat, ngentot. Lima orang aja yang boleh masuk</p>	<p>Ah, gagal anjing. Kontol! Anjing. Boro-boro crown... Udah-udah keluar, keluar, keluar. ... Uda berapa jam si streaming? Anjing udah 2 jam, cok. Main ginian doang ngentot</p>
<p>Kalau gua lebih ke, gue nih (sebagai) anak- anak juga nih. Lu, anak-anak tongkrongan juga, gitu. Jadi, kaya gua ngomong, 'Anjing lu ya, bangsat. Lu nipu gua ya, tai!' gitu kan, biasanya gue gitu kan ke anak- anak. Maksudnya, anak-anak pada ketawa juga kan. Jadi kaya, apa ya, jokes-jokes-nya itu kaya jokes-jokes anak tongkrongan gitu lho. Jadi kaya mereka curhat ke gua, gua kata-katain balik. Terus gue lagi main, gue diisengin, katain balik. Gitu lho</p>	<p>Maju, maju, maju! Ini siapa ngentot! Ahh! Itu siapa namanya, Affiliator (nama pemain lain), anjing. Goblok! Hiu, anjing! Maju-maju, bangsaat! Winnery jangan kaya kontrol lu ya! Ngentot  ni ngentot ya! Istri Luthfi (nama pemain lain)! Kontool! Memek! Aah kontrol! Gue ga main-main anjing! Memek!"</p>

## 2.2 Instagram

*Instagram* adalah *platform* media sosial yang populer dengan lebih dari 1 miliar pengguna aktif bulanan di seluruh dunia. Diluncurkan pada tahun 2010, *Instagram* fokus pada berbagi foto dan video, dan telah berkembang menjadi salah satu *platform* terkemuka dalam industri ini. Pengguna *Instagram* berasal dari berbagai latar belakang dan rentang usia, meskipun mayoritas pengguna berada dalam rentang usia 18-34 tahun. Platform ini menawarkan berbagai fitur yang memungkinkan pengguna untuk mengunggah foto dan video ke akun mereka, serta menambahkan teks, stiker, filter, dan efek lainnya untuk meningkatkan kreativitas. Fitur utama *Instagram* meliputi *Feed*, di mana pengguna dapat melihat *postingan* terbaru dari akun yang diikuti. Selain itu, terdapat *Stories*, yang memungkinkan pengguna membagikan foto dan video yang akan hilang setelah 24 jam. *IGTV* adalah fitur yang memungkinkan pengguna mengunggah video yang lebih panjang, sementara *Direct Messaging* memungkinkan pengguna untuk berkomunikasi secara pribadi dengan pengguna lain. Interaksi dan *engagement* merupakan bagian penting dari *Instagram*. Pengguna dapat menyukai, mengomentari, dan berbagi *postingan* dari pengguna lain. Fitur "*Explore*" memungkinkan pengguna menemukan konten baru yang sesuai dengan minat mereka. *Instagram* juga menjadi *platform* yang populer untuk periklanan. Bisnis dapat menggunakan berbagai opsi iklan yang ditawarkan, termasuk iklan dalam *feed*, *Stories*, dan *Explore*. Platform ini juga menyediakan profil bisnis dengan informasi kontak dan fitur analitik untuk melacak performa *postingan* dan iklan. Dengan fokus pada konten berkualitas tinggi, *Instagram* memungkinkan pengguna menggunakan *hashtag* untuk meningkatkan visibilitas dan mempermudah temuan oleh pengguna dengan minat yang sama. Secara keseluruhan, *Instagram* menyediakan ekosistem yang luas bagi individu, kreator konten, dan bisnis untuk berinteraksi dan berbagi konten secara kreatif.

## **2.3 Artificial Intelligence**

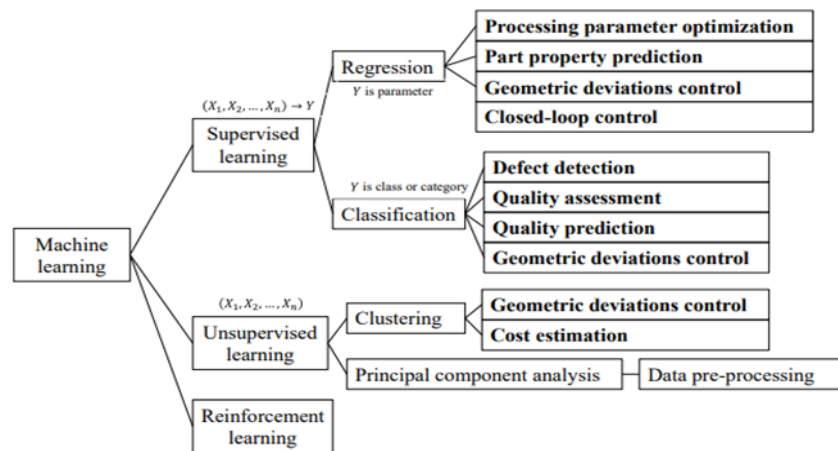
*Artificial Intelligence (AI)* atau dalam bahasa Indonesia dikenal sebagai Kecerdasan Buatan, adalah suatu bidang ilmu komputer yang bertujuan untuk membuat mesin atau sistem komputer dapat melakukan tugas yang biasanya memerlukan kecerdasan manusia (Novamizanti & Siadari, 2022). *AI* mencakup pengembangan algoritma dan model komputasional yang dirancang untuk memungkinkan mesin untuk "belajar" dari *data*, mengenali pola, mengambil keputusan, dan melakukan tugas-tugas yang cerdas (Markiewicz & Zheng, 2020).

### **2.3.1 Machine Learning**

*Machine Learning* adalah bagian dari Kecerdasan Buatan (*AI*) yang melibatkan penggunaan algoritme dan model statistik untuk memungkinkan komputer atau mesin belajar dari dan membuat prediksi atau keputusan berdasarkan *data* tanpa diprogram secara eksplisit (Carleo dkk., 2019). Dengan kata lain, ini adalah metode atau pendekatan di mana mesin dapat secara otomatis belajar dari pengalaman, meningkat seiring waktu, dan melakukan tugas tanpa diprogram secara eksplisit untuk setiap tugas tertentu.

Pembelajaran Mesin melibatkan penggunaan sejumlah besar *data*, yang dimasukkan ke dalam algoritme dan model untuk melatih mesin. Selama proses pelatihan, mesin mengidentifikasi pola, tren, dan hubungan dalam *data* dan menyesuaikan parameter modelnya untuk mengoptimalkan kinerjanya. Setelah mesin dilatih, ia dapat menggunakan pengetahuan yang dipelajarinya untuk membuat prediksi atau keputusan pada *data* baru yang tidak terlihat (Zhou dkk., 2017).

Menurut (Meng dkk., 2020) *Machine Learning* dapat dikategorikan menjadi tiga jenis seperti pada gambar 2.1 berikut:



**Gambar 2. 1 Machine Learning Method**

### 1. *Supervised Learning*

Dalam jenis Pembelajaran Mesin ini, model dilatih pada *data* berlabel, di mana *data* masukan disertai dengan label keluaran yang sesuai. *Model* belajar dari *data* berlabel untuk membuat prediksi atau mengklasifikasikan *data* baru yang tidak terlihat ke dalam kategori yang telah ditentukan sebelumnya. Beberapa contoh dari jenis pembelajaran mesin *supervised learning* ini adalah deteksi penyakit berdasarkan gejala pada pasien, prediksi harga rumah berdasarkan fitur-fitur seperti ukuran, lokasi, dan jumlah kamar tidur, klasifikasi jenis hewan berdasarkan ciri-ciri fisik, identifikasi jenis tumbuhan berdasarkan karakteristik daun dll.

### 2. *Unsupervised Learning*

Dalam jenis Pembelajaran Mesin ini, model dilatih pada *data* yang tidak berlabel, di mana *data* masukan tidak memiliki label keluaran yang menyertainya. *Model* belajar mengidentifikasi pola atau struktur dalam *data*, seperti pengelompokan atau pengurangan dimensi. Beberapa contoh dari jenis pembelajaran mesin *unsupervised learning* ini adalah pengelompokan berita berdasarkan topik yang serupa, segmentasi konsumen

berdasarkan perilaku belanja, reduksi dimensi pada *data* gambar untuk analisis visual dan analisis sentimen pada ulasan produk tanpa label positif/negatif.

### 3. *Reinforcement Learning*

Dalam Pembelajaran Mesin jenis ini, model belajar dari interaksi dengan lingkungan dan menerima umpan balik dalam bentuk penghargaan atau hukuman. *Model* belajar mengambil tindakan di lingkungan untuk memaksimalkan imbalan kumulatif dari waktu ke waktu. Beberapa contoh dari jenis pembelajaran mesin reinforcement *learning* ini adalah pengembangan agen permainan video yang belajar bermain game secara optimal, pengaturan pengiriman inventaris oleh robot dalam gudang, pengembangan sistem perdagangan saham yang memutuskan kapan harus membeli atau menjual dan pengembangan agen cerdas untuk bermain permainan catur atau *Go*.

#### **2.3.2 *Natural Language Processing***

*NLP* adalah singkatan dari Natural Language Processing, yang merupakan subbidang Kecerdasan Buatan (*AI*) yang berfokus pada memungkinkan komputer atau mesin untuk memahami, menafsirkan, dan berinteraksi dengan bahasa manusia dengan cara yang bermakna dan bermanfaat. *NLP* melibatkan pengembangan algoritma dan model yang dapat memproses dan menganalisis bahasa manusia, termasuk teks dan ucapan, untuk mengekstrak makna, mengidentifikasi pola, dan menghasilkan respons (Markiewicz & Zheng, 2020).

#### **2.4 *Decision Tree***

Algoritma *Decision Tree* (*DT*), yang termasuk dalam kelas algoritma pembelajaran yang diawasi sebagian besar lebih disukai untuk memecahkan masalah klasifikasi tetapi bagaimanapun juga, itu dapat digunakan dalam



mengklasifikasikan serta dalam kasus regresi. Ini terdiri dari simpul dalam yang mewakili struktur cabang, kumpulan *data*, yang mewakili keputusan yang diberikan oleh algoritma, dan setiap simpul daun mewakili hasil. Ada dua simpul, pertama adalah simpul keputusan, yang digunakan untuk mengambil keputusan dan memiliki berbagai cabang; dan yang kedua adalah *leaf node*, yang merupakan output dari decision *node* dan tidak memiliki cabang lagi. Itu berutang namanya pada pohon karena kesamaan bentuknya. Simpul akar adalah titik awal yang selanjutnya berkembang ke berbagai cabang menjadikannya struktur seperti pohon. Pohon keputusan hanya membagi pohon menjadi sub-pohon berdasarkan jawaban atas pertanyaan, yaitu apakah ya atau tidak (Bansal dkk., 2022). Berikut merupakan persamaan yang digunakan pada algoritma decision tree

$$Gain(S,A) = Entropy(S) - \sum_{i=1}^n \frac{|S_i|}{|S|} Entropy(S_i) \quad 2.1$$

$$Entropy(S) = - \sum_{i=1}^n p_i \log_2 p_i \quad 2.2$$

Berdasarkan persamaan diatas maka dapat di jelaskan sebagai berikut:

1. Menyiapkan *data training*.
2. Menentukan akar dari pohon.
3. Hitung *Gain* Untuk memilih atribut sebagai akar, didasarkan pada nilai *Gain* tertinggi dari atribut-atribut yang ada. Untuk menghitung *Gain* dapat menggunakan persamaan 2.1 diatas.
4. Ulangi langkah kedua hingga setiap cabang terpenuhi. Sementara itu, untuk penghitungan nilai *Entropy* dapat menggunakan persamaan 2.2
5. Proses partisipasi pohon keputusan akan berhenti saat semua cabang dalam *node N* mendapat kelas yang sama.

## 2.5 Confusion Matrix

*Confusion Matrix* adalah sebuah *matrix* yang memvisualisasikan kinerja algoritma klasifikasi menggunakan *data matrix*, dengan membandingkan

klasifikasi yang di prediksi dengan klasifikasi aktual dalam bentuk *false positive*, *true positive*, *false negative* dan *true negative* (Awad & Khanna, 2015).

*Confusion Matrix* menggunakan tabel 2.2 seperti *matrix* di bawah ini:

**Tabel 2. 2 Confusion Matrix**

		<i>Predicted</i>	
		<i>Positive</i>	<i>Negative</i>
<i>Actual</i>	<i>Positive</i>	<i>TP</i>	<i>FN</i>
	<i>Negative</i>	<i>FP</i>	<i>TN</i>

Keterangan:

*TP* : Total kasus *Positive* yang dikategorikan sebagai *Positive*

*FP* : Total kasus *Negative* yang dikategorikan sebagai *Positive*

*TN* : Total kasus *Negative* yang dikategorikan sebagai *Negative*

*FN* : Total kasus *Positive* yang dikategorikan sebagai *Negative*

Berdasarkan nilai *True Negative (TN)*, *Values Positive (FP)*, *Values Negative (FN)*, dan *True Positive (TP)* dapat diperoleh nilai akurasi, presisi dan *Recall* menggunakan persamaan berikut

$$Akurasi = \frac{(TP + TN)}{(TP + TN + FP + FN)} \quad 2.3$$

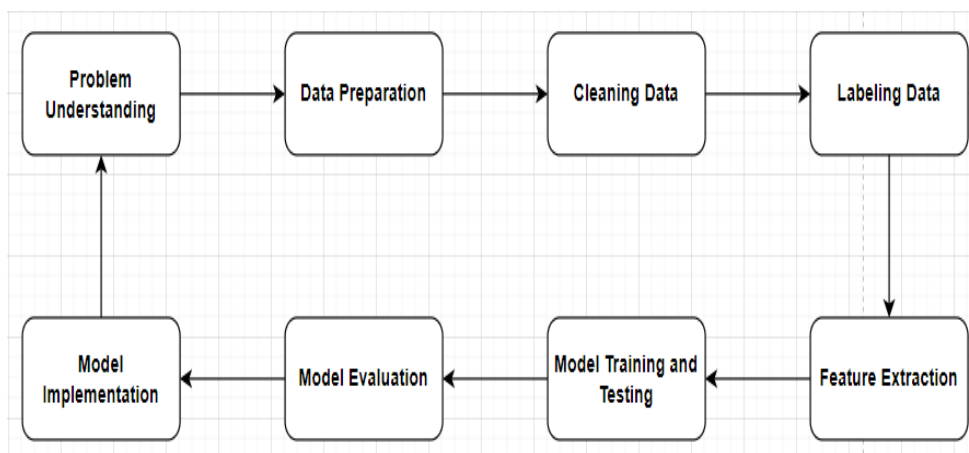
$$Precision = \frac{TP}{(TP + FP)} \quad 2.4$$

$$Recall = \frac{TP}{(TP + FN)} \quad 2.5$$

## 2.6 Metode Pengembangan *Machine Learning Model*

Berdasarkan studi literatur yang dilakukan oleh peneliti maka, dapat diberikan analisa bahwa dalam penelitian yang akan dilakukan ini terdapat kerangka pemikiran yang dijadikan untung menunjang keberlangsungannya penelitian ini yang disebut dengan *Machine Learning Life Cycle*.

*Machine Learning Life Cycle* merupakan salah satu model pengembangan model *machine learning* pada kecerdasan buatan (*AI*). Tujuan dari *machine learning life cycle* adalah untuk menemukan solusi untuk masalah yang diberikan dengan menerapkan model *Machine Learning*. *Machine Learning Life Cycle* dibagi menjadi delapan tahap utama seperti pada gambar 2.2 dibawah ini:



**Gambar 2. 2** *Machine Learning Life Cycle*

Pada tahapan pembuatan model klasifikasi *Machine Learning* menggunakan *Machine Learning Life Cycle* memiliki tahapan-tahapan sebagai berikut:

1. *Problem Understanding*

Tahap ini melibatkan pemahaman yang komprehensif tentang masalah yang akan dipecahkan, termasuk identifikasi jenis masalah, pemahaman konteks penggunaan model, serta penetapan batasan dan tujuan yang jelas.

## 2. *Data Preparation*

Pada tahap ini dilakukan persiapan *data* yang akan digunakan dalam pelatihan dan pengujian model, termasuk pengumpulan *data* yang relevan, penggabungan *data* dari sumber yang berbeda, pengaturan *data* dalam format yang sesuai, dan pembagian *data* menjadi *data* latih (*data training*) dan *data* uji (*data testing*) (Hameed & Naumann, 2020).

## 3. *Data Cleaning*

Tahap *data cleaning* ini melibatkan pembersihan *data* untuk memastikan *data* yang digunakan dalam model adalah *data* yang berkualitas dan valid, termasuk penghapusan *data* yang tidak relevan, penggantian atau penghapusan *data* yang hilang atau rusak, serta penanganan duplikasi atau *outliers* dalam *data* (Chu dkk., 2016).

## 4. *Data Labeling*

Tahap *data labeling* merupakan tahap pemberian label pada *data* yang akan digunakan sebagai *data* latih model, yang mengindikasikan kategori atau kelas yang ingin diprediksi oleh model, baik secara manual atau menggunakan teknik pemrosesan *data* atau pemrosesan bahasa alami (*NLP*) (Fredriksson dkk., 2020).

## 5. *Feature Extraction*

*Feature extraction* merupakan tahap mengubah *data* menjadi representasi numerik yang dapat digunakan sebagai fitur dalam model (Salahat & Qasaimah, 2017), termasuk perubahan *data* teks menjadi vektor kata atau vektor dokumen, penggunaan metode seperti *bag-of-words* atau *word embeddings*, serta ekstraksi fitur-fitur lainnya yang relevan untuk masalah yang dipecahkan..

## 6. *Training dan Testing Model*

Yang akan dilakukan pada tahap ini yaitu melatih model menggunakan algoritma atau teknik *machine learning* yang sesuai, dengan membagi *data* latih dan *data* uji, mengfitkan model pada

*data* latih, dan menguji performa model pada *data* uji untuk membuat prediksi atau klasifikasi *data* (Spjuth dkk., 2021).

#### 7. *Model Evaluation*

Tahap ini dilakukan sebagai pengukuran kinerja model menggunakan metrik evaluasi yang sesuai (Spjuth dkk., 2021), seperti akurasi, presisi, *recall*, atau *F1-score*, untuk mengevaluasi seberapa baik model dalam melakukan prediksi atau klasifikasi *data* serta memvalidasi model dengan menggunakan teknik *k-fold cross validation*.

#### 8. *Implementation*

Dalam tahap ini dilakukan dengan mengimplementasikan model dalam lingkungan produksi atau aplikasi yang diinginkan, termasuk integrasi model ke dalam sistem produksi, pengujian ulang model di lingkungan produksi, pengelolaan versi model, serta pemeliharaan model untuk memastikan kinerjanya tetap optimal seiring waktu (Spjuth dkk., 2021).

## 2.7 Penelitian Terkait

Tabel 2. 3 Penelitian Terkait

Judul, Penulis dan Tahun	Jumlah dan Atribut Dataset	Algoritma	<i>Preprocessing</i>	Feature Selection	Validasi	Open Source Dataset	Akurasi
Klasifikasi Ujaran Kebencian pada Cuitan dalam Bahasa Indonesia (Antariksa dkk., 2019)	20.305 <i>data</i> dengan dua atribut	Bernoulli Naïve Bayes, Multinomial Naïve Bayes, Logistic Regression dan Support	Menghapus karakter spesial, menghapus huruf acak, menghapus <i>data</i> yang tidak lengkap	<i>TF-IDF</i> , n-grams, dan word2vec	<i>Confusion matrix</i>	-	Bernoulli Naïve Bayes + <i>TF-IDF</i> 94%, Bernoulli Naïve Bayes + N-gram Char Level 80%, Multinomial Naïve Bayes + Word2Vec 98%, Logistic Regression + <i>TF-IDF</i> 98%, Logistic Regression + Ngram Char Level 98%, Logistic

Judul, Penulis dan Tahun	Jumlah dan Atribut Dataset	Algoritma	<i>Preprocessing</i>	Feature Selection	Validasi	Open Source Dataset	Akurasi
		Vector Machine					Regression + Word2Vec 98%, SVM + TF-IDF 98%, SVM + N-gram Word Level 98% dan SVM + Word2Vec 98%
Analisis Sentimen Pola Pikir Masyarakat Indonesia Terkait Virus Covid-19 Dalam Media Sosial <i>Twitter</i> Menggunakan	22.856 <i>data</i> dengan tiga atribut	Metode Rule Based Leksikon	<i>Data cleaning, lowering case, emoticon conversion, stemming, and stopwords data</i>	<i>Tokenizing</i> kata unigram (1 suku kata) dan bigram (2 suku kata)	<i>Confusion matrix</i>	-	<i>Accuracy</i> sebesar 81 %, nilai <i>precision</i> sebesar 93 %, nilai <i>recall</i> sebesar 95 %, dan nilai <i>f1score</i> sebesar 83%.

Judul, Penulis dan Tahun	Jumlah dan Atribut Dataset	Algoritma	<i>Preprocessing</i>	Feature Selection	Validasi	Open Source Dataset	Akurasi
Metode Rule Based Leksikon (Riskiyanti, 2022)							
Deteksi Cyberbullying pada Facebook Menggunakan Algoritma K-Nearest Neighbor (Detect Cyberbullying on Facebook Using K-Nearest Neighbor	3.000 <i>data</i> dengan tiga atribut	KNN	Convert <i>Emoticon</i> dan <i>Punctuation</i> , <i>Case Folding</i> , <i>Tokenizing</i> , <i>Filtering</i> , <i>Stemming</i> dan <i>Weighting</i> <i>Labelling</i>	Normalisasi	<i>Confusion Matrix</i> , <i>cross validation</i>	-	1-NN 70.80%, influence punctuation 71.40%, lowercase dan uppercase influence 71.43%, 3-NN 65.93%



<b>Judul, Penulis dan Tahun</b>	<b>Jumlah dan Atribut Dataset</b>	<b>Algoritma</b>	<i>Preprocessing</i>	<b>Feature Selection</b>	<b>Validasi</b>	<b>Open Source Dataset</b>	<b>Akurasi</b>
Algorithm) (Hasan & Wati, 2021)							
Analisa Tanggapan Terhadap Psbb Di Indonesia Dengan Algoritma <i>Decision Tree</i> Pada <i>Twitter</i> (Aditya Quantano Surbakti dkk., 2021)	2.439 <i>data</i> dengan dua atribut	<i>Decision Tree, Naïve Bayes dan K-NN</i>	<i>Cleansing, Stemming, Tokenisasi dan Stopword</i>	-	<i>Confusion Matrix, cross validation</i>	-	<i>Decision Tree 84.78%, Naive Bayes 84.73% dan K-NN 77.75%</i>
PENDETEKSIAN <i>HATE SPEECH</i> PADA SOSIAL	13.169 dengan	<i>Support Vector Machine</i>	<i>Case Folding, Remove Special Character,</i>	<i>Text Vectorization</i>	<i>Confusion Matrix</i>	-	<i>Decision Tree 79%, Support Vector Machine</i>

Judul, Penulis dan Tahun	Jumlah dan Atribut Dataset	Algoritma	Preprocessing	Feature Selection	Validasi	Open Source Dataset	Akurasi
MEDIA INDONESIA DENGAN ALGORITMA SUPPORT VECTOR MACHINE (SVM) DAN DECISION TREE (Kusuma & Pamungkas, 2023)	dua atribut	(SVM), BaggingClassifier dan Decision Tree	Remove Username, Remove Url Remove White Space Dan Remove Number				83% dan BaggingClassifier 82%.
Deteksi Komentar Cyberbullying pada Media Sosial	170 data dengan	Random Forest	Cleaning, Case Folding Tokenization &	Word Embedding FastText	Confusion Matrix	-	Random Forest sebesar 84%

Judul, Penulis dan Tahun	Jumlah dan Atribut Dataset	Algoritma	<i>Preprocessing</i>	Feature Selection	Validasi	Open Source Dataset	Akurasi
<i>Instagram Menggunakan Algoritma Random Forest (Santoso dkk., 2023)</i>	dua atribut		<i>Normalization, Stopwords Removal, Stemming</i>				
<i>Big Five Personality Detection on Twitter Users Using Gradient Boosted Decision Tree Method</i>	275 data dengan lima atribut	<i>Gradient Boosted Decision Tree</i>	<i>Data Cleaning, Case Folding, Tokenization, Stopwords Removal dan Stemming</i>	<i>TF-IDF dan SMOTE</i>	<i>Confusion Matrix</i>	-	Hasil pada skenario pertama menunjukkan akurasi tertinggi pada uji fitur sentimen sebesar 41,82%, dengan presisi dan recall yang sama pada angka 41,82%, serta nilai f1-score sebesar 41,61%.

Judul, Penulis dan Tahun	Jumlah dan Atribut Dataset	Algoritma	<i>Preprocessing</i>	Feature Selection	Validasi	Open Source Dataset	Akurasi
(Sugiono & Maharani, 2023)							<p>Pada skenario kedua dengan penerapan SMOTE. Akurasi tertinggi pada uji fitur emosi meningkat menjadi 60,36% (+27,63%), presisi meningkat menjadi 59,63% (+33,64%), recall meningkat menjadi 60,36% (+27,63%), dan f1-score meningkat menjadi 60,78% (+27,61%).</p>

Judul, Penulis dan Tahun	Jumlah dan Atribut Dataset	Algoritma	Preprocessing	Feature Selection	Validasi	Open Source Dataset	Akurasi
Studi Komparasi Metode Analisis Sentimen Naïve Bayes, SVM, dan Logistic Regression Pada Piala Dunia 2022 (Anbari & Sugiantoro, 2023)	10.680 data dengan dua atribut	Naïve Bayes, SVM, dan Logistic Regression	Remove Username, Remove Url, Remove Emoticon, Remove Special Character	TF-IDF	Confusion Matrix	-	Bernaouli Naïve Bayes menghasilkan nilai parameter presisi adalah 71%, parameter recall 99%, dan akurasi 76%. Sedangkan metode Support Vector Classifier menghasilkan nilai parameter presisi adalah 94%, parameter recall 93%, dan akurasi 92%. Adapun metode Logistic Regression menghasilkan nilai parameter presisi

Judul, Penulis dan Tahun	Jumlah dan Atribut Dataset	Algoritma	<i>Preprocessing</i>	Feature Selection	Validasi	Open Source Dataset	Akurasi
							adalah 93%, parameter recall 93%, dan akurasi 92%.
<i>Multi-Tier Sentiment Analysis of Social Media Text Using Supervised Machine Learning</i> (Rahman dkk., 2023)	156.060 data dengan lima atribut	<i>SVM, Decision tree dan Naïve Bayes</i>	<i>Remove Stopword, Remove Special Character, Lemmatization dan Tokenization</i>	<i>TF-IDF</i>	<i>Confusion Matrix</i>	-	SVM Multi-tier model memperoleh akurasi sebesar 0.52%, Decision tree Multi-tier model 0.52% dan Naïve Bayes Multi-tier model 0.56% sedangkan untuk SVM Single-tier model memperoleh akurasi sebesar 0.51%, Decision

Judul, Penulis dan Tahun	Jumlah dan Atribut Dataset	Algoritma	<i>Preprocessing</i>	Feature Selection	Validasi	Open Source Dataset	Akurasi
							tree Single-tier model 0.51% dan Naïve Bayes Single-tier model memperoleh akurasi sebesar 0.55%
PERBANDINGAN ALGORITMA <i>SUPPORT VECTOR MACHINE</i> DENGAN <i>DECISION TREE</i> PADA APLIKASI	1.500 data dengan dua atribut	<i>SVM</i> dan <i>Decision Tree</i>	<i>Case Folding, Tokenization, Stemming, Filtering</i>	<i>TF-IDF</i>	<i>Confusion Matrix</i>	-	<i>SVM</i> melakukan perhitungan probabilitas pada data dengan pembobotan kata menggunakan <i>TF-IDF</i> mendapatkan hasil nilai akurasi sebesar 84.2%, sedangkan pada algoritma

<b>Judul, Penulis dan Tahun</b>	<b>Jumlah dan Atribut Dataset</b>	<b>Algoritma</b>	<i>Preprocessing</i>	<b>Feature Selection</b>	<b>Validasi</b>	<b>Open Source Dataset</b>	<b>Akurasi</b>
RUANG GURU (Hassanah dkk., 2023)							Decision Tree mendapatkan nilai akurasi sebesar 70%.



Secara umum, penelitian-penelitian di atas menunjukkan bahwa penerapan algoritma *Machine Learning* membantu dalam mendeteksi kalimat *toxic* pada media sosial. Namun, terdapat variasi akurasi yang dihasilkan tergantung pada berbagai faktor seperti ukuran *dataset*, atribut, algoritma, preprocessing, dan *feature selection* yang digunakan. Oleh karena itu beberapa aspek yang memberikan nilai tambah atau kebaruan dalam penelitian ini. Pertama, penelitian ini memiliki fokus pada penentuan ukuran *dataset* yang sesuai, memungkinkan pemilihan sampel yang lebih representatif untuk analisis. Kedua, dengan berfokus pada algoritma *Decision Tree*, pendekatan yang diambil dapat menghasilkan solusi yang berbeda dan lebih akurat dibandingkan dengan menggunakan beberapa algoritma sekaligus seperti dalam penelitian terdahulu. Ketiga, penggunaan fitur ekstraksi *TF-IDF* memberikan kontribusi pada peningkatan kualitas *dataset*. Keempat, dengan membatasi penelitian pada *platform Instagram*, peneliti dapat menghasilkan hasil yang lebih terfokus dan khusus pada media sosial tersebut.