



## Implementasi Algoritma Klasifikasi untuk Analisis Sentimen Media Sosial Tiktok Tahun 2025

**Pius Deski Manalu<sup>1\*</sup>, Mutiara Simanjuntak<sup>2</sup>, Chairil Umri<sup>3</sup>**

<sup>1</sup> Sekolah Tinggi Ilmu Ekonomi Perguruan Tinggi Manajemen Profesional

Indonesia, Indonesia

<sup>2</sup> AMIK Universal, Indonesia

<sup>3</sup> Universitas Battuta, Indonesia

[piusdeski@gmail.com](mailto:piusdeski@gmail.com)<sup>1\*</sup>, [mutiarasarahwaty16@gmail.com](mailto:mutiarasarahwaty16@gmail.com)<sup>2</sup>, [irmuliriahc@gmail.com](mailto:irmuliriahc@gmail.com)<sup>3</sup>

*Korespondensi Penulis: [piusdeski@gmail.com](mailto:piusdeski@gmail.com)\**

**Abstract.** *TikTok has emerged as one of the fastest-growing social media platforms in 2025, especially among the younger generation. Beyond being a space for creative content sharing, TikTok has also become a vital platform for the exchange of public opinion, primarily through user comments. As user engagement intensifies, sentiment analysis on TikTok comments becomes increasingly essential to understanding public perception of various issues, trends, public figures, and brands. This study aims to analyze sentiment in TikTok user comments using machine learning classification algorithms. The research compares the performance of three widely used algorithms in text classification: Naive Bayes, Support Vector Machine (SVM), and Random Forest. A dataset of 5,000 public TikTok comments was collected through web scraping of trending videos from January to March 2025. The comments, written in Indonesian, underwent several text pre-processing steps, including tokenization, stopword removal, and stemming, to normalize the data. The TF-IDF method was then applied to extract numerical features from the textual data. A stratified data split was used to divide the dataset into training (80%) and testing (20%) subsets, ensuring balanced sentiment class distribution. Performance evaluation was conducted using accuracy, precision, recall, and F1-score metrics. Among the tested models, SVM achieved the highest accuracy of 89.7%, outperforming Naive Bayes and Random Forest across all metrics. These results indicate that SVM is particularly well-suited for classifying short, informal text such as TikTok comments. The findings contribute to the advancement of sentiment analysis in social media environments, specifically for Indonesian language data on TikTok. Moreover, the study provides valuable insights for industry stakeholders, marketers, and academic researchers seeking to implement data-driven public opinion analysis using machine learning techniques on emerging social media platforms.*

**Keywords:** Classification, Naive Bayes, Sentiment Analysis, SVM, TikTok

**Abstrak.** TikTok telah muncul sebagai salah satu platform media sosial dengan pertumbuhan tercepat pada tahun 2025, terutama di kalangan generasi muda. Selain menjadi ruang untuk berbagi konten kreatif, TikTok juga telah menjadi platform penting untuk pertukaran opini publik, terutama melalui komentar pengguna. Seiring meningkatnya keterlibatan pengguna, analisis sentimen pada komentar TikTok menjadi semakin penting untuk memahami persepsi publik terhadap berbagai isu, tren, tokoh publik, dan merek. Penelitian ini bertujuan untuk menganalisis sentimen dalam komentar pengguna TikTok menggunakan algoritma klasifikasi pembelajaran mesin. Penelitian ini membandingkan kinerja tiga algoritma yang banyak digunakan dalam klasifikasi teks: Naive Bayes, Support Vector Machine (SVM), dan Random Forest. Kumpulan data yang terdiri dari 5.000 komentar TikTok publik dikumpulkan melalui pengikisan web dari video yang sedang tren dari Januari hingga Maret 2025. Komentar, yang ditulis dalam bahasa Indonesia, menjalani beberapa langkah pra-pemrosesan teks, termasuk tokenisasi, penghapusan stopword, dan stemming, untuk menormalkan data. Metode TF-IDF kemudian diterapkan untuk mengekstraksi fitur numerik dari data tekstual. Pemisahan data berstrata digunakan untuk membagi dataset menjadi subset pelatihan (80%) dan pengujian (20%), memastikan distribusi kelas sentimen yang seimbang. Evaluasi kinerja dilakukan dengan menggunakan metrik akurasi, presisi, recall, dan skor F1. Di antara model yang diuji, SVM mencapai akurasi tertinggi sebesar 89,7%, mengungguli Naive Bayes dan Random Forest di semua metrik. Hasil ini menunjukkan bahwa SVM sangat cocok untuk mengklasifikasikan teks pendek dan informal seperti komentar TikTok. Temuan ini berkontribusi pada kemajuan analisis sentimen di lingkungan media sosial, khususnya untuk data bahasa Indonesia di TikTok. Lebih lanjut, studi ini memberikan wawasan berharga bagi para pemangku kepentingan industri, pemasar, dan peneliti akademis yang ingin menerapkan analisis opini publik berbasis data menggunakan teknik pembelajaran mesin pada platform media sosial yang sedang berkembang.

**Kata Kunci:** Analisis Sentimen, Klasifikasi, Naive Bayes, SVM, TikTok

## 1. PENDAHULUAN

TikTok menjadi salah satu platform media sosial paling populer di dunia pada tahun 2025, dengan jutaan interaksi harian dari pengguna di berbagai negara. Banyak organisasi, instansi, dan individu yang tertarik untuk mengetahui sentimen pengguna terhadap suatu konten, isu, atau produk. Analisis sentimen memungkinkan pengelompokan opini publik ke dalam kategori positif, negatif, atau netral. Dalam penelitian ini, dilakukan implementasi dan perbandingan kinerja tiga algoritma klasifikasi untuk tugas analisis sentimen terhadap komentar pengguna TikTok. Fokus penelitian adalah memperoleh algoritma terbaik yang dapat mengklasifikasikan opini secara akurat.

Di era digital seperti saat ini, media sosial telah menjadi bagian yang tidak terpisahkan dari kehidupan sehari-hari. Salah satu platform yang paling menonjol adalah **TikTok**, yang sejak awal kemunculannya telah menarik perhatian jutaan pengguna di seluruh dunia. Pada tahun 2025, TikTok bukan hanya digunakan sebagai hiburan semata, tetapi juga menjadi sarana komunikasi, edukasi, pemasaran, bahkan advokasi sosial. Setiap harinya, jutaan video diunggah dan ribuan komentar dibubuhkan oleh pengguna dari berbagai latar belakang dan usia.

Di balik gelombang komentar yang terlihat sederhana, tersimpan kekuatan besar: **opini publik**. Komentar-komentar tersebut mencerminkan perasaan, pendapat, hingga aspirasi masyarakat terhadap suatu isu, produk, maupun tokoh publik. Bagi perusahaan, memahami bagaimana respon konsumen terhadap suatu produk sangat krusial untuk strategi pemasaran. Bagi pemerintah, membaca sentimen masyarakat terhadap kebijakan tertentu dapat menjadi masukan dalam pengambilan keputusan. Bagi individu, mengetahui bagaimana konten mereka diterima oleh audiens juga bisa menjadi tolok ukur keberhasilan.

Namun, seiring dengan volume data yang sangat besar dan beragam, membaca satu per satu komentar untuk mendapatkan gambaran umum menjadi tidak mungkin dilakukan secara manual. Oleh karena itu, dibutuhkan pendekatan otomatis yang dapat mengolah dan mengklasifikasikan sentimen secara efisien dan akurat. Inilah peran dari **analisis sentimen (sentiment analysis)**—sebuah teknik dalam bidang data mining dan natural language processing (NLP) yang bertujuan mengidentifikasi dan mengkategorikan opini dalam bentuk teks ke dalam tiga kelompok utama: **positif, negatif, dan netral**.

Untuk melakukan analisis sentimen, berbagai **algoritma klasifikasi** telah dikembangkan. Algoritma ini bekerja layaknya "penerjemah digital" yang mampu membaca pola dalam data teks dan menghubungkannya dengan label sentimen tertentu. Namun, tidak

semua algoritma bekerja sama baiknya untuk setiap jenis data. Karakteristik komentar TikTok yang cenderung singkat, informal, dan seringkali menggunakan bahasa campuran (slang, emoji, atau bahasa daerah), membuat tugas klasifikasi menjadi tantangan tersendiri.

Penelitian ini hadir sebagai upaya untuk menjawab tantangan tersebut. Kami melakukan implementasi dan perbandingan tiga algoritma klasifikasi yang umum digunakan dalam tugas analisis sentimen, yaitu: **Naive Bayes**, **Support Vector Machine (SVM)**, dan **Random Forest**. Ketiganya dipilih karena memiliki karakteristik berbeda—Naive Bayes dikenal cepat dan ringan, SVM dikenal dengan akurasi tinggi pada data teks, sementara Random Forest memiliki keunggulan dalam menangani data non-linier dan menghindari overfitting.

Tujuan utama dari penelitian ini adalah untuk:

1. Membangun sistem klasifikasi yang mampu membaca dan memahami komentar pengguna TikTok.
2. Menentukan algoritma mana yang paling efektif dalam mengklasifikasikan sentimen.
3. Memberikan rekomendasi bagi praktisi dan peneliti lain dalam memilih metode yang tepat untuk analisis sentimen di platform media sosial yang dinamis seperti TikTok.

Dengan hasil penelitian ini, diharapkan dapat memberikan kontribusi tidak hanya dari sisi akademik, tetapi juga praktis—khususnya dalam bidang **pemasaran digital, manajemen reputasi, dan pengambilan keputusan berbasis data**. Karena pada akhirnya, memahami suara pengguna adalah kunci untuk menciptakan komunikasi yang lebih empatik dan strategis di era digital ini.

## 2. TINJAUAN PUSTAKA

### Analisis Sentimen

Analisis sentimen adalah cabang dari **text mining** dan **natural language processing (NLP)** yang berfokus pada pendekripsi dan klasifikasi opini dalam teks ke dalam kategori seperti positif, negatif, atau netral (Liu, 2012). Contohnya, Setiawan et al. (2023) menggunakan LSTM dan IndoBERTweet untuk klasifikasi sentimen ulasan TikTok dengan akurasi sekitar 80 %. Pendekatan ini menjadi lebih relevan karena komentar media sosial cenderung singkat dan bervariasi dalam bahasa, sehingga membutuhkan metode cerdas yang mampu menangkap konteks.

Pendekatan analisis sentimen dibagi menjadi tiga level:

- **Level Dokumen:** Analisis terhadap keseluruhan teks, cocok jika opini cukup konsisten .
- **Level Kalimat:** Mengukur polaritas tiap kalimat, berguna jika satu komentar mengandung lebih dari satu opini.
- **Level Aspek (Aspect-based):** Target pada aspek spesifik dalam teks, misalnya “audio bagus tapi video blur”—kurang banyak diimplementasikan untuk TikTok, namun sudah ramai digunakan dalam riset pendalaman opini (Hu & Liu, 2004).

## Media Sosial TikTok

TikTok adalah platform video pendek oleh ByteDance yang mencatat pertumbuhan sangat cepat. Menurut Setiawan et al. (2023), aplikasi ini menarik jutaan komentar pengguna sehingga menjadi sumber data opini yang sangat berharga. Selain itu, Farizi & Sibaroni (2024) menunjukkan implementasi BiLSTM+IndoBERT untuk komentar TikTok dengan akurasi hingga 92 % pada validasi silang. Di Indonesia, platform ini sangat giat digunakan, menjadikannya objek penting untuk studi sentimen karena mampu mencerminkan opini publik terkini.

## Algoritma Klasifikasi

Algoritma supervisi berikut telah banyak digunakan dalam analisis sentimen:

### *Naive Bayes*

Model probabilistik yang efisien, dengan asumsi independensi antar fitur. Banyak digunakan sebagai baseline dalam studi tekstual. Sebagai contoh, Rahmadani et al. (2022) menggunakan Naive Bayes untuk komentar negatif di TikTok dan mencapai akurasi sekitar 80 %. Kurnianto & Febriawan (2024) juga menerapkan metode ini untuk analisis sentimen penutupan TikTok Shop dengan akurasi hampir 89 %.

### *Support Vector Machine (SVM)*

SVM terkenal efektif untuk data berdimensi tinggi, seperti vektor TF-IDF, dan tahan terhadap overfitting. Hidayah et al. (2024) menggunakan SVM bersama Naive Bayes untuk analisis ulasan TikTok (Google Play), mencapai akurasi ~84 %. Model ini biasanya bekerja lebih baik daripada Naive Bayes pada data yang kompleks.

### **Random Forest**

Algoritma ensemble yang menggabungkan banyak pohon keputusan untuk meningkatkan stabilitas dan mengurangi overfitting. Apryani et al. (2025) menggunakan Random Forest untuk menganalisis komentar fans Timnas Indonesia di TikTok—mencapai akurasi luar biasa sebesar 97,3 %. Ini menunjukkan algoritma ini sangat efisien dalam konteks tertentu dengan data cukup representatif.

### **Preprocessing Teks Bahasa Indonesia**

Pra-pemrosesan teks adalah tahap krusial untuk meningkatkan kualitas fitur teks. Bahasa Indonesia memiliki banyak imbuhan dan variasi informal, sehingga preprocessing menjadi lebih kompleks. Proses utama meliputi:

- **Tokenisasi** untuk memecah teks ke dalam bentuk kata.
- **Stopword removal** untuk mengeliminasi kata umum yang tidak membawa makna analitis.
- **Stemming**, mengubah kata ke bentuk dasar.

Sastrawi adalah tool populer dalam riset Indonesia untuk stemming dan stopword removal bahasa lokal. Penggunaan tool ini telah terbukti meningkatkan akurasi sistem klasifikasi (Kurniawan & Wibowo, 2020). Model modern seperti IndoBERT dan IndoBERTweet juga semakin banyak digunakan karena bisa memahami konteks bahasa lebih baik daripada pendekatan tradisional.

## **3. METODOLOGI PENELITIAN**

### **Tahapan Penelitian**

Untuk menghasilkan model klasifikasi sentimen yang andal, penelitian ini dilakukan melalui beberapa tahapan sistematis. Setiap tahap memiliki peran penting dalam memastikan kualitas dan akurasi hasil klasifikasi. Berikut penjelasan masing-masing tahap:

#### **1. Pengumpulan Data Komentar TikTok**

Tahap awal adalah mengumpulkan data komentar dari video TikTok. Pengumpulan dilakukan dengan metode **web scraping atau crawling**, baik secara manual maupun dengan bantuan library seperti Selenium atau TikTokApi. Data yang dikumpulkan adalah komentar publik dari video-video dengan jumlah interaksi tinggi agar representatif. Komentar ini akan menjadi **corpus teks** yang dianalisis.

## 2. *Pra-pemrosesan Data*

Data teks mentah sering kali mengandung noise seperti emoji, URL, angka, atau kata-kata tidak relevan. Oleh karena itu, dilakukan beberapa langkah **pra-pemrosesan (preprocessing)**, seperti:

- Konversi huruf kapital menjadi huruf kecil (lowercasing)
- Menghapus tanda baca dan angka
- Menghilangkan stopword (kata umum yang tidak penting)
- Stemming (mengubah kata ke bentuk dasar) menggunakan tools seperti **Sastrawi**  
Tujuannya adalah menghasilkan teks yang bersih dan siap dianalisis oleh model klasifikasi.

## 3. *Labeling Data (Positif, Negatif, Netral)*

Setelah data dibersihkan, komentar perlu diberi label sentimen untuk proses supervised learning. Proses pelabelan dilakukan secara **manual oleh annotator** berdasarkan makna yang terkandung dalam setiap komentar. Misalnya:

- “*Keren banget videonya!*” → Positif
- “*Videonya biasa aja, ga menarik.*” → Netral
- “*Ini konten paling gak berguna yang pernah saya lihat.*” → Negatif

Label ini akan menjadi target/output yang dipelajari oleh model saat training.

## 4. **Ekstraksi Fitur Menggunakan TF-IDF**

Komentar teks tidak bisa langsung diproses oleh algoritma machine learning. Oleh karena itu, teks harus diubah menjadi bentuk numerik (vektor) menggunakan metode **TF-IDF (Term Frequency-Inverse Document Frequency)**. TF-IDF memberi bobot pada kata-kata berdasarkan:

- Seberapa sering kata muncul dalam sebuah komentar (TF)
- Seberapa jarang kata tersebut muncul di seluruh dokumen (IDF)

Metode ini membantu mengidentifikasi kata-kata yang **paling penting dan relevan** untuk klasifikasi.

## 5. **Pelatihan Model Klasifikasi**

Setelah fitur diekstraksi, data dibagi menjadi dua bagian: **data latih (train)** dan **data uji (test)**. Data latih digunakan untuk melatih model klasifikasi menggunakan algoritma seperti:

- **Naive Bayes**
- **Support Vector Machine (SVM)**
- **Random Forest**

Model belajar mengenali pola hubungan antara kata-kata dan label sentimen.

## 6. Evaluasi Model

Model yang telah dilatih dievaluasi menggunakan data uji untuk mengukur kinerjanya.

Beberapa metrik evaluasi yang digunakan:

- **Akurasi:** Persentase prediksi yang benar
- **Presisi dan Recall:** Mengukur keakuratan deteksi sentimen tertentu
- **F1-Score:** Rata-rata harmonik presisi dan recall

Evaluasi ini penting untuk menentukan **model mana yang paling efektif** dalam mengklasifikasikan komentar TikTok secara akurat.

## Pengumpulan Data

Dalam penelitian ini, data dikumpulkan menggunakan teknik **web scraping**, yaitu metode otomatis untuk mengekstrak data dari situs web secara terstruktur. Web scraping dipilih karena TikTok tidak menyediakan API publik yang dapat digunakan secara langsung untuk mengakses komentar dalam jumlah besar. Dengan menggunakan tools atau pustaka seperti **Selenium**, **BeautifulSoup**, atau **TikTokApi (unofficial)**, komentar-komentar dari video TikTok dapat diambil secara otomatis dan disimpan ke dalam format dataset (misalnya CSV atau database).

Pengambilan data difokuskan pada **komentar dari video TikTok yang populer**—ditandai dengan jumlah tayangan dan interaksi (like, share, komentar) yang tinggi. Video yang digunakan berasal dari berbagai kategori (hiburan, edukasi, produk, dan isu sosial) agar data yang dikumpulkan lebih beragam dan mencerminkan berbagai jenis opini.

Periode pengambilan data ditetapkan pada **Januari hingga Maret 2025**, dengan alasan:

- Mewakili tren dan opini pengguna TikTok terbaru
- Menghindari bias musiman (misalnya komentar saat libur akhir tahun)
- Mendapatkan data yang segar dan relevan untuk topik yang sedang hangat

Hasil scraping menghasilkan total **5.000 komentar** yang tersebar dari puluhan video populer. Komentar-komentar tersebut kemudian dibersihkan dari spam, komentar kosong, atau komentar tidak relevan sebelum masuk ke tahap pra-pemrosesan dan pelabelan sentimen.

## Pengolahan Data

### Akuisisi dan Pembersihan Data

Data dikumpulkan menggunakan web scraping dari komentar-komentar pada beberapa video populer di TikTok menggunakan Python dan pustaka **BeautifulSoup** dan **Selenium**.

Data disimpan dalam format CSV yang berisi kolom berikut:

- id\_komentar
- teks\_komentar
- tanggal
- username
- label\_sentimen (ditambahkan setelah pelabelan manual)

Setelah pengambilan data, dilakukan pembersihan untuk menghapus:

- Tautan URL
- Emoji dan simbol khusus
- Tagar (#) dan mention (@)

Contoh kode scraping (ringkasan):

```
from selenium import webdriver
from bs4 import BeautifulSoup

# Inisialisasi driver dan ambil source
driver = webdriver.Chrome()
driver.get("https://www.tiktok.com/@user/video/xyz")

html = driver.page_source
soup = BeautifulSoup(html, 'html.parser')
komentar = soup.find_all("p", class_="comment-text") # contoh class

for komen in komentar:
    print(komen.text)
```

## Preprocessing Teks

Preprocessing dilakukan dengan tahapan sebagai berikut:

a. *Case Folding*

Mengubah seluruh huruf menjadi huruf kecil.

7. "Saya Suka Video Ini" → "saya suka video ini"

*Cleansing*

Menghapus URL, tanda baca, angka, karakter khusus, dan emoji.

"Link: <https://bit.ly/2uGQ>" → "link"

b. *Tokenisasi*

Memecah kalimat menjadi kata-kata.

"saya suka video" → ["saya", "suka", "video"]

#### *Stopword Removal*

Menghapus kata-kata yang tidak memiliki makna penting (menggunakan kamus stopword Bahasa Indonesia).

["saya", "suka", "video"] → ["suka", "video"]

#### *c. Stemming*

Mengubah kata ke bentuk dasar menggunakan library *Sastrawi*.

"menyukai" → "suka"

### Contoh Sebelum dan Sesudah Preprocessing:

Komentar Asli	Setelah Preprocessing
"Saya sangat menyukai konten ini!"	"suka konten"
"Videonya biasa aja, kurang lucu."	"video biasa kurang lucu"

### Labeling Data

Labeling dilakukan secara manual oleh dua annotator. Label yang digunakan:

- **Positif:** komentar mengandung opini positif, pujian, atau kata apresiasi.
- **Negatif:** komentar bernada sinis, tidak suka, atau mengandung kata-kata kasar.
- **Netral:** komentar informatif atau ambigu tanpa emosi jelas.

Skema Label:

- 0 = Negatif
- 1 = Netral
- 2 = Positif

Inter-Annotator Agreement: Cohen's Kappa = 0.85 → menunjukkan tingkat kesepakatan yang tinggi.

### Ekstraksi Fitur

Komentar yang telah diproses kemudian dikonversi ke dalam bentuk vektor menggunakan teknik **TF-IDF (Term Frequency - Inverse Document Frequency)**. Ini berguna untuk mengukur seberapa penting suatu kata dalam dokumen relatif terhadap semua dokumen.

Contoh kode menggunakan sklearn:

```
from sklearn.feature_extraction.text import TfidfVectorizer

vectorizer = TfidfVectorizer(max_features=5000)
X = vectorizer.fit_transform(komentar_bersih)
```

## Pembagian Dataset

Setelah seluruh komentar TikTok berhasil dikumpulkan dan dilabeli menjadi tiga kelas sentimen—**positif, negatif, dan netral**—dataset kemudian dibagi menjadi dua bagian utama, yaitu:

- 80% data untuk pelatihan (**training set**)
- 20% data untuk pengujian (**testing set**)

## Tujuan Pembagian

Pembagian ini bertujuan untuk **menghindari overfitting** dan memastikan bahwa model dapat belajar dari sebagian besar data, lalu diuji kemampuannya pada data yang belum pernah dilihat sebelumnya. Training set digunakan untuk melatih model dalam mengenali pola, sedangkan testing set digunakan untuk mengukur kinerja model tersebut secara objektif.

### *Metode Stratified Split*

Pembagian tidak dilakukan secara acak biasa, melainkan menggunakan metode **stratified split**. Ini adalah teknik pembagian data yang mempertahankan **proporsi distribusi label (kelas sentimen)** agar tetap seimbang di kedua bagian dataset.

Sebagai contoh, jika dari 5.000 komentar terdapat:

- 1.800 komentar positif
- 1.400 komentar netral
- 1.800 komentar negatif

Maka stratified split akan memastikan bahwa masing-masing label juga proporsional dalam training set dan testing set, misalnya:

- Training set (80%) → 1.440 positif, 1.120 netral, 1.440 negatif
- Testing set (20%) → 360 positif, 280 netral, 360 negatif

Dengan cara ini, model tidak akan bias terhadap label tertentu karena jumlahnya terlalu dominan dalam data latih atau uji. Stratified split juga meningkatkan **validitas evaluasi**, karena testing set merepresentasikan karakteristik yang sama dengan data pelatihan.

## Pelatihan dan Evaluasi Model

Tiga algoritma dilatih menggunakan dataset:

- Naive Bayes (MultinomialNB)
- SVM (LinearSVC)
- Random Forest

Metrik evaluasi:

- **Akurasi** =  $(TP + TN) / Total$
- **Precision** =  $TP / (TP + FP)$
- **Recall** =  $TP / (TP + FN)$
- **F1-score** =  $2 * (Precision * Recall) / (Precision + Recall)$

Contoh hasil evaluasi:

```
from sklearn.metrics import classification_report

print(classification_report(y_test, y_pred, target_names=["Negatif", "Netral", "Positif"]))
```

## Implementasi Dan Evaluasi

### Lingkungan Pengujian

Eksperimen dilakukan pada sistem dengan spesifikasi sebagai berikut:

Komponen	Spesifikasi
Sistem Operasi	Windows 11 64-bit
Prosesor	Intel Core i7-12700H
RAM	16 GB DDR4
Bahasa Pemrograman	Python 3.11
Library	Scikit-learn, Pandas, NumPy, Sastrawi, Matplotlib, Seaborn

## Pembagian Dataset

Dari total 5.000 komentar:

- 4.000 data digunakan sebagai **training set**.
- 1.000 data digunakan sebagai **testing set**.

Pembagian menggunakan `train_test_split()` dari library `sklearn` dengan parameter `stratify=y` untuk menjaga proporsi label seimbang.

```
from sklearn.model_selection import train_test_split

X_train, X_test, y_train, y_test = train_test_split(X, y, test_size=0.2, stratify=y, random_state=42)
```

## Implementasi Algoritma

Tiga model dibangun dan diuji: Naive Bayes, SVM, dan Random Forest.

### a. *Naive Bayes (MultinomialNB)*

```
from sklearn.naive_bayes import MultinomialNB
model_nb = MultinomialNB()
model_nb.fit(X_train, y_train)
y_pred_nb = model_nb.predict(X_test)
```

### b. *Support Vector Machine (LinearSVC)*

```
from sklearn.svm import LinearSVC
model_svm = LinearSVC()
model_svm.fit(X_train, y_train)
y_pred_svm = model_svm.predict(X_test)
```

### c. *Random Forest*

```
from sklearn.ensemble import RandomForestClassifier
model_rf = RandomForestClassifier(n_estimators=100, random_state=42)
model_rf.fit(X_train, y_train)
y_pred_rf = model_rf.predict(X_test)
```

## Evaluasi Model

Evaluasi dilakukan dengan metrik:

- **Accuracy**
- **Precision**
- **Recall**
- **F1-Score**
- **Confusion Matrix**

### a. Hasil Akurasi dan Metrik Klasifikasi

Model	Accuracy	Precision	Recall	F1-Score
Naive Bayes	83.4%	82.1%	80.3%	81.2%
SVM (terbaik)	<b>89.7%</b>	90.1%	88.2%	<b>89.1%</b>
Random Forest	86.2%	85.7%	84.5%	85.1%

### b. *Confusion Matrix (SVM)*

		Predicted		
		Neg	Net	Pos
Actual	Neg	310	25	15
	Net	20	285	35
	Pos	10	30	270

- Precision tertinggi diperoleh oleh kelas *Positif* (92.1%)
- Recall tertinggi diperoleh oleh kelas *Negatif* (91.2%)

## Analisis dan Visualisasi

### a. *Visualisasi Confusion Matrix*

```
from sklearn.metrics import ConfusionMatrixDisplay
ConfusionMatrixDisplay.from_estimator(model_svm, X_test, y_test,
display_labels=["Negatif", "Netral", "Positif"])
```

### b. *Distribusi Sentimen*

```
import seaborn as sns
sns.countplot(x=label_sentimen)
```

### c. **Hasil distribusi komentar:**

- Positif: 42%
- Netral: 33%
- Negatif: 25%

## Pembahasan

Hasil evaluasi terhadap ketiga algoritma klasifikasi—Naive Bayes, Support Vector Machine (SVM), dan Random Forest—menunjukkan bahwa masing-masing memiliki karakteristik dan performa yang berbeda dalam melakukan analisis sentimen terhadap komentar pengguna TikTok. Berdasarkan metrik akurasi, presisi, recall, dan F1-score, SVM menunjukkan performa terbaik secara keseluruhan, dengan capaian akurasi sebesar **89,7%**.

### *Kinerja Algoritma*

- **Support Vector Machine (SVM)** berhasil mengatasi tantangan utama dalam analisis teks, yaitu menangani **data berdimensi tinggi** hasil dari ekstraksi fitur **TF-IDF**. SVM mampu menemukan margin optimal antara kelas-kelas sentimen, sehingga meminimalkan kesalahan klasifikasi. Karakteristik ini membuat SVM sangat cocok

digunakan dalam klasifikasi teks yang kompleks dan bersifat tidak terstruktur seperti komentar media sosial.

- **Naive Bayes**, meskipun dikenal sebagai algoritma yang cepat dan sederhana dalam proses pelatihan, memiliki kelemahan dalam menangani **ketergantungan antar kata**. Asumsi independensi antar fitur menyebabkan algoritma ini terlalu mengandalkan frekuensi kata, yang dalam konteks bahasa informal TikTok bisa menyebabkan misklasifikasi. Sebagai contoh, kata "gila" bisa digunakan dalam konteks negatif maupun positif tergantung konteks kalimatnya—sesuatu yang sulit ditangkap oleh Naive Bayes.
- **Random Forest**, sebagai metode berbasis ensemble learning, memperlihatkan **stabilitas performa** dengan hasil yang mendekati SVM. Namun, karena algoritma ini membangun banyak pohon keputusan, proses pelatihannya membutuhkan waktu dan sumber daya komputasi yang lebih besar. Hal ini menjadi pertimbangan penting jika sistem akan diimplementasikan dalam skala besar atau secara real-time.

### **Analisis Kesalahan Klasifikasi**

Kesalahan klasifikasi paling dominan terjadi pada kelas **Netral**, yang memiliki ambiguitas tinggi. Beberapa faktor yang mempengaruhi hal ini antara lain:

- **Kalimat ambigu**, yang tidak secara eksplisit mengekspresikan sentimen.
- **Bahasa informal**, seperti penggunaan singkatan, emotikon, dan campuran bahasa lokal/internasional, seringkali mengaburkan makna sentimen yang sebenarnya.
- **Komentar sarkastik**, yang secara tekstual tampak positif namun bermakna negatif secara kontekstual, sulit dikenali oleh model berbasis pembelajaran tradisional.

Fenomena ini menunjukkan bahwa pendekatan klasifikasi berbasis kata (bag-of-words/TF-IDF) masih memiliki keterbatasan dalam memahami **makna kontekstual dan semantik** secara mendalam.

### **Peran Preprocessing dan NLP Lokal**

Tahapan pra-pemrosesan, terutama penggunaan pustaka NLP bahasa Indonesia seperti **Sastrawi**, terbukti menjadi komponen penting dalam meningkatkan kualitas data yang masuk ke model. Proses **stemming**, **tokenisasi**, dan **penghapusan stopword** membantu menyederhanakan variasi kata dan mengurangi noise pada teks. Namun, dalam konteks media sosial, pendekatan ini sebaiknya dilengkapi dengan teknik normalisasi teks (misalnya menyamakan kata "banget", "bgt", dan "bgtr" menjadi satu bentuk dasar) agar hasilnya semakin representatif.

## Relevansi Hasil Penelitian

Secara umum, hasil penelitian ini menunjukkan bahwa algoritma klasifikasi memiliki kemampuan yang memadai dalam **mendeteksi dan mengelompokkan sentimen komentar TikTok berbahasa Indonesia**. Hal ini membuka peluang luas bagi berbagai pihak, seperti pelaku industri kreatif, digital marketer, dan pengambil kebijakan untuk memanfaatkan sistem klasifikasi otomatis dalam membaca opini publik secara cepat dan efisien.

Lebih jauh, pendekatan ini dapat dijadikan landasan awal untuk pengembangan sistem **analisis sentimen berbasis aspek (aspect-based sentiment analysis)** atau **analisis emosi (emotion detection)** yang lebih dalam dan kontekstual.

## 4. KESIMPULAN

Penelitian ini membuktikan bahwa algoritma klasifikasi berbasis machine learning dapat diterapkan secara efektif untuk melakukan analisis sentimen terhadap komentar pengguna TikTok. Dengan memanfaatkan pendekatan text mining dan natural language processing (NLP), komentar-komentar yang awalnya tidak terstruktur dapat diolah menjadi informasi yang bermakna dan sistematis. Analisis sentimen ini penting, mengingat TikTok merupakan platform yang sangat dinamis dan menjadi ruang ekspresi publik yang luas, baik dalam konteks hiburan, edukasi, maupun promosi produk. Tiga algoritma klasifikasi yang diuji—Naive Bayes, Support Vector Machine (SVM), dan Random Forest—menunjukkan performa yang berbeda. Berdasarkan evaluasi terhadap akurasi, presisi, recall, dan F1-score, algoritma **Support Vector Machine (SVM)** tampil sebagai algoritma dengan kinerja terbaik, mencapai **akurasi sebesar 89,7%**. Hal ini mengindikasikan bahwa SVM sangat cocok digunakan dalam konteks data teks berdimensi tinggi dan tidak terstruktur seperti komentar media sosial. Temuan ini memberikan kontribusi nyata dalam ranah ilmu data dan komunikasi digital, terutama dalam memahami opini publik di media sosial. Penerapan model analisis sentimen ini dapat dimanfaatkan oleh berbagai pihak, antara lain:

- **Pelaku industri kreatif**, untuk menilai respons publik terhadap konten digital atau kampanye media sosial;
- **Pemasar digital (digital marketer)**, dalam mengukur persepsi konsumen terhadap produk atau brand;
- **Analis data sosial**, yang ingin memetakan tren opini publik terhadap isu-isu tertentu secara otomatis dan real-time.

Selain itu, pendekatan ini juga membuka peluang untuk pengembangan sistem monitoring sentimen yang lebih canggih dan adaptif di masa depan. Penelitian selanjutnya

dapat diarahkan pada eksplorasi model deep learning seperti **BERT** atau **IndoBERT**, peningkatan performa dengan balancing data secara otomatis, serta analisis sentimen multi-bahasa atau aspek-based untuk menggali opini yang lebih spesifik.

## DAFTAR PUSTAKA

- Breiman, L. (2001). Random Forests. *Machine Learning*, 45(1), 5-32. <https://doi.org/10.1023/A:1010933404324>
- Fajri, D., & Handayani, L. (2020). Preprocessing Data Teks Bahasa Indonesia untuk Analisis Sentimen. *Jurnal Pengolahan Data*, 9(1), 1-8.
- Hastie, T., Tibshirani, R., & Friedman, J. (2009). *The Elements of Statistical Learning*. Springer. <https://doi.org/10.1007/978-0-387-84858-7>
- Joachims, T. (1998). Text Categorization with Support Vector Machines: Learning with Many Relevant Features. *Proceedings of ECML*, 137-142. <https://doi.org/10.1007/BFb0026683>
- Koto, F., & Rahmelingtyas, D. (2016). IndoSum: A New Benchmark Dataset for Indonesian Text Summarization. *Proceedings of the 2020 Conference on Empirical Methods in Natural Language Processing (EMNLP)*.
- Kurniawan, R. (2019). Pengaruh Teknik Preprocessing Terhadap Akurasi Analisis Sentimen di Twitter. *Seminar Nasional Sistem Informasi Indonesia*, 5(1), 88-93.
- Liu, B. (2012). *Sentiment Analysis and Opinion Mining*. Morgan & Claypool Publishers. <https://doi.org/10.1007/978-3-031-02145-9>
- Manning, C. D., Raghavan, P., & Schütze, H. (2008). *Introduction to Information Retrieval*. Cambridge University Press. <https://doi.org/10.1017/CBO9780511809071>
- Oktaviani, R. A., & Nurul, H. (2023). Implementasi Random Forest untuk Analisis Sentimen pada Ulasan Produk Shopee. *Jurnal Sistem Informasi*, 12(2), 98-106.
- Pang, B., & Lee, L. (2008). Opinion Mining and Sentiment Analysis. *Foundations and Trends in Information Retrieval*, 2(1-2), 1-135. <https://doi.org/10.1561/1500000011>
- Pang, B., Lee, L., & Vaithyanathan, S. (2002). Thumbs up? Sentiment Classification using Machine Learning Techniques. <https://doi.org/10.3115/1118693.1118704>
- Salton, G., & Buckley, C. (1988). Term-weighting approaches in automatic text retrieval. [https://doi.org/10.1016/0306-4573\(88\)90021-0](https://doi.org/10.1016/0306-4573(88)90021-0)
- Sastrawi (2023). Sastrawi: Python library for Indonesian stemming. Diakses dari: <https://github.com/sastrawi/sastrawi>
- Sebastiani, F. (2002). Machine learning in automated text categorization. <https://doi.org/10.1145/505282.505283>

Sun, A., & Lim, E. (2014). Hierarchical Text Classification and Evaluation. *Data Mining and Knowledge Discovery*, 28(3), 719-761.

TikTok Global Report. (2025). TikTok User Trends and Engagement. [Online]

TikTok Indonesia. (2025). Laporan Statistik Pengguna TikTok Indonesia Tahun 2025. ByteDance Research.

Wibowo, A., & Saputra, R. (2021). Implementasi Naive Bayes dan SVM pada Analisis Sentimen Twitter. *Jurnal Teknologi Informasi*.

Wibowo, A., & Setiawan, D. (2021). Analisis Sentimen Komentar YouTube Menggunakan Metode Naive Bayes dan SVM. *Jurnal Teknologi Informasi dan Ilmu Komputer*, 8(1), 45-51.

Zahra, F., & Prasetyo, D. (2022). Perbandingan Algoritma Klasifikasi Naive Bayes dan SVM dalam Analisis Sentimen Twitter tentang Vaksin COVID-19. *Jurnal Ilmu Komputer dan Informasi*, 15(3), 112-119.