

ELT(EXTRACT, LOAD, TRANSFORM) PROCESS

“Analisis Sentimen Komentar YouTube
Warga Indonesia terhadap Danantara”

By :

Al Fitra Nur Ramadhani

202210370311264

Data, Information, and Knowledge B

BACKGROUND

- Media sosial kini telah menjadi platform utama bagi masyarakat untuk mengekspresikan opini, emosi, dan tren yang berkembang. Salah satu topik yang sedang menarik perhatian publik Indonesia adalah Danantara Indonesia, Badan Pengelola Investasi Daya Anagata Nusantara, yang memiliki peran strategis dalam mengonsolidasikan dan mengoptimalkan investasi pemerintah untuk mendukung pertumbuhan ekonomi nasional.
- Dalam proyek ini, saya melakukan analisis sentimen terhadap komentar-komentar yang ditinggalkan oleh warga Indonesia pada video YouTube Raymond Chin yang membahas Danantara. Data yang digunakan dikumpulkan melalui YouTube API v3, yang memungkinkan ekstraksi komentar secara langsung dari platform YouTube.

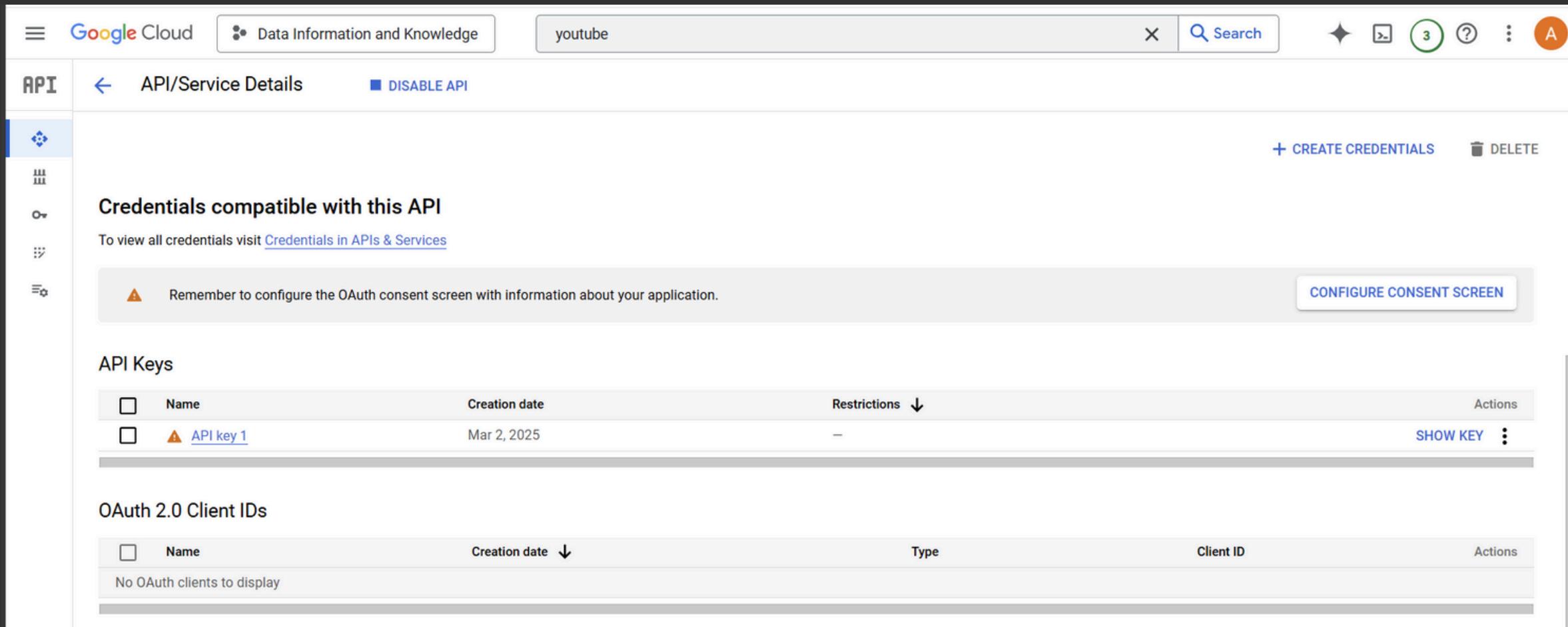
TOOLS

- VS Code
- Python
- Mongo DB
- Youtube API V3
- Youtube Platform



EXTRACT

LETS GO EXECUTE (EXTRACT)



- Untuk dapatkan API Key nya, pergi ke Google Cloud dan cari Youtube Data API V3
- Click Create Credentials untuk dapatkan API Key mu

LETS GO EXECUTE (EXTRACT)

```
● ● ●  
1 dev = "Your API Key"
```

```
● ● ●  
1 import googleapiclient.discovery  
2 import pandas as pd  
3  
4 api_service_name = "youtube"  
5 api_version = "v3"  
6 DEVELOPER_KEY = dev  
7  
8 youtube = googleapiclient.discovery.build(  
9     api_service_name, api_version, developerKey=DEVELOPER_KEY)  
10  
11 request = youtube.commentThreads().list(  
12     part="snippet",  
13     videoId="Tl6fDFA4xNI",  
14     maxResults=100  
15 )  
16  
17 comments = []  
18  
19 # Execute the request.  
20 response = request.execute()  
21  
22 # Get the comments from the response.  
23 for item in response['items']:  
24     comment = item['snippet']['topLevelComment']['snippet']  
25     public = item['snippet']['isPublic']  
26     comments.append([  
27         comment['authorDisplayName'],  
28         comment['publishedAt'],  
29         comment['likeCount'],  
30         comment['textOriginal'],  
31         public  
32     ])
```

- Copy Your API Key, dan Youtube ID pada Program , untuk mengambil comments dari Youtube

```
● ● ●  
1 while (1 == 1):  
2     try:  
3         nextPageToken = response['nextPageToken']  
4     except KeyError:  
5         break  
6     nextPageToken = response['nextPageToken']  
7     # Create a new request object with the next page token.  
8     nextRequest = youtube.commentThreads().list(part="snippet", videoId="Tl6fDFA4xNI", maxResults=100, pageToken=nextPageToken)  
9     # Execute the next request.  
10    response = nextRequest.execute()  
11    # Get the comments from the next response.  
12    for item in response['items']:  
13        comment = item['snippet']['topLevelComment']['snippet']  
14        public = item['snippet']['isPublic']  
15        comments.append([  
16            comment['authorDisplayName'],  
17            comment['publishedAt'],  
18            comment['likeCount'],  
19            comment['textOriginal'],  
20            public  
21        ])  
22  
23 df = pd.DataFrame(comments, columns=['author', 'updated_at', 'like_count', 'text','public'])  
24 df.info()
```

LETS GO EXECUTE (EXTRACT)

```
● ● ●  
1 df.info()
```

```
● ● ●  
1 df.head(10)
```

- Cek Data nya apakah sudah dimuat, dan berapa entri, data type, dll
- Jika ingin lebih yakin, cek data nya dalam bentuk tabel dengan df.head()

	author	updated_at	# like_count	text	public
0	@ahmadsyaifuddin3553	2025-03-02T02:50:13Z	0	Ayo taruh an ga bakal an di korupsi	True
1	@lenyselvianti4031	2025-03-02T02:49:21Z	0	Semoga ada transparansi	True
2	@andrihermansyah7	2025-03-02T02:47:54Z	0	Itulah yang ditakutkan. Eksekusi amkt	True

```
● ● ●  
1 df.to_csv('Data_Scrapped_Youtube_RaymondChin.csv', index=False)  
2
```

- Langkah terakhir, download datanya dalam bentuk CSV

LOAD

LETS GO EXECUTE (LOAD)

The screenshot shows the MongoDB Atlas interface. On the left, a sidebar lists various services: Overview, DATABASE (Clusters selected), SERVICES (Atlas Search, Stream Processing, Triggers, Migration, Data Federation), SECURITY (Quickstart, Backup, Database Access, Network Access, Advanced), and New On Atlas (6). The main area displays Cluster0 details: Version 8.0.4, Region GCP / Jakarta (asia-southeast2), Type Replica Set - 3 nodes, Backups Inactive, Linked App Services None linked, and links to ATLAS SQL and ATLAS SEARCH. A prominent callout box suggests loading sample datasets. Below this, a monitoring section shows R/W metrics (0.0B/s, 0.0B/s), Connections (3.0), and Network traffic (In 5.85 B/s, Out 224.45 B/s). A green 'Upgrade' button is visible.

The screenshot shows a step-by-step guide for connecting to Cluster0. Step 1: Set up connection security (done). Step 2: Choose a connection method (done). Step 3: Connect. The 'Connecting with MongoDB Driver' section includes: 1. Select your driver and version (Driver: Python, Version: 4.7 or later). It recommends installing the latest driver version. 2. Install your driver (Run the following on the command line: `python -m pip install "pymongo[srv]"`). It notes to use appropriate Python 3 executable. 3. Add your connection string into your application code (Use this connection string in your application: `mongodb+srv://alfitranurr:<db_password>@cluster0.vfngh.mongodb.net/?retryWrites=true&w=majority&appName=Cluster0`). It advises replacing `<db_password>` with the password for the `alfitranurr` database user. Ensure any option params.

- Pergi ke MongoDB untuk membuat Database, kemudian Click Connect dan pilih Driver, kemudian sesuaikan dengan seperti pada gambar, dan copy connection string nya

LETS GO EXECUTE (LOAD)



```
1 import pandas as pd
2 import numpy as np
3 import matplotlib.pyplot as plt
4 import seaborn as sns
5 from wordcloud import WordCloud, STOPWORDS
6 from pymongo import MongoClient
7 import emoji
8 import csv
9 import json
```



```
1 # 1. Koneksi ke MongoDB Atlas
2 # Ganti <db_password> dengan kata sandi MongoDB Anda
3 mongo_uri = "mongodb+srv://alfitranurr:YourCode@cluster0.vfngh.mongodb.net/?retryWrites=true&w=majority&appName=Cluster0"
4 client = MongoClient(mongo_uri)
5
6 # Pilih database dan collection
7 db = client['danantara_db'] # Nama database, bisa diganti sesuai keinginan
8 collection = db['youtube_comments'] # Nama collection, bisa diganti sesuai keinginan
```

- Masuk ke dalam Program, dan berikut adalah beberapa library yang harus di impor

- Pastekan URL yang sudah di copy dari MongoDb dan ubah sesuai PW mu, dan sesuaikan untuk nama database dan tabel nya untuk ditampilkan di Mongo DB

LETS GO EXECUTE (LOAD)



```
1 # 2. Fungsi untuk memuat data CSV ke MongoDB
2 def load_csv_to_mongodb(csv_file_path):
3     # Baca file CSV menggunakan pandas
4     df = pd.read_csv(csv_file_path)
5
6     # Konversi DataFrame ke list of dictionaries
7     data = df.to_dict(orient='records')
8
9     # Hapus data lama di collection (opsional, uncomment jika ingin reset)
10    # collection.delete_many({})
11
12    # Insert data ke MongoDB
13    collection.insert_many(data)
14    print(f"Berhasil memuat {len(data)} dokumen ke MongoDB!")
```



```
1 # 3. Fungsi untuk membaca data dari MongoDB ke DataFrame
2 def read_from_mongodb():
3     # Ambil semua data dari collection
4     data = list(collection.find())
5
6     # Konversi ke DataFrame
7     df = pd.DataFrame(data)
8
9     # Hapus kolom '_id' yang otomatis dibuat oleh MongoDB (opsional)
10    if '_id' in df.columns:
11        df = df.drop('_id', axis=1)
12
13    return df
```

- Berikut adalah program untuk memuat data CSV ke dalam mongo DB, dengan mengubah awalnya berbentuk dataframe menjadi list of dictionary

- Kemudian read file yang ada di mongoDB

LETS GO EXECUTE (LOAD)



```
1 # 4. Main execution
2 if __name__ == "__main__":
3     # Ganti dengan path file CSV Anda
4     csv_file_path = "Data_Scrapped_Youtube_RaymondChin.csv"
5
6     # Load data CSV ke MongoDB
7     load_csv_to_mongodb(csv_file_path)
8
9     # Baca data dari MongoDB ke DataFrame
10    df = read_from_mongodb()
11
12    # Tampilkan informasi dasar DataFrame
13    print(df.info())
14    print(df.head())
```

- Execute Program didalam main, jika sudah akan menampilkan seperti pada gambar

```
Berhasil memuat 1758 dokumen ke MongoDB!
<class 'pandas.core.frame.DataFrame'>
RangeIndex: 3515 entries, 0 to 3514
Data columns (total 5 columns):
 #   Column      Non-Null Count  Dtype  
--- 
 0   author       3514 non-null   object 
 1   updated_at   3515 non-null   object 
 2   like_count   3515 non-null   int64  
 3   text         3515 non-null   object 
 4   public       3515 non-null   bool   
dtypes: bool(1), int64(1), object(3)
memory usage: 113.4+ KB
None
```

	author	updated_at	like_count	\
0	@ahmadsyaifuddin3553	2025-03-02T02:50:13Z	0	
1	@lenyelvianti4031	2025-03-02T02:49:21Z	0	
2	@andrihermansyah7	2025-03-02T02:47:54Z	0	
3	@agussupriadi6083	2025-03-02T02:43:54Z	0	
4	@Bonsai-q4f	2025-03-02T02:41:49Z	0	

```
text  public
0 Ayo taruh an ga bakal an di korupsi kaya jiwas...  True
1                               Semoga ada transparansi😊  True
2 ttulah yang ditakutkan Eksokusi amburadul an  True
```



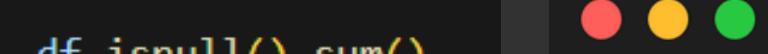
TRANSFORM

LETS GO EXECUTE (TRANSFORM)



```
1 # Tentukan kolom yang berisi tipe list
2 list_columns = ['text'] # Ganti dengan nama kolom yang sesuai
3
4 # Mengonversi kolom dengan tipe list menjadi string
5 for col in list_columns:
6     df[col] = df[col].apply(lambda x: str(x) if isinstance(x, list) else x)
7
8 # Cek duplikat setelah konversi
9 duplikat = df.duplicated()
10 print(f"Jumlah duplikat: {duplikat.sum()}")
11 print("Baris yang duplikat:")
12 print(df[duplikat])
13
```

- Cek duplikat data nya



```
df.isnull().sum()
```

```
✓ 0.0s
```

```
author      1
updated_at   0
like_count   0
text         0
public       0
dtype: int64
```



```
1 # Menghapus baris dengan nilai null di kolom 'text'
2 df = df.dropna(subset=['author'])
```



```
1 # Periksa kembali jumlah nilai null
2 print(df.isnull().sum())
```

- Cek Missing Value
- Ternyata ada 1 di kolom “Author”, kemudian bersihkan dan tampilkan lagi untuk memastikan apakah benar-benar sudah dibersihkan

LETS GO EXECUTE (TRANSFORM)

```
df.info()  
[2]: ✓ 0.0s  
  
[  
    <class 'pandas.core.frame.DataFrame'>  
    Index: 1757 entries, 0 to 1756  
    Data columns (total 5 columns):  
     #   Column      Non-Null Count  Dtype     
     ---  --          -----          ---  
     0   author       1757 non-null   object    
     1   updated_at   1757 non-null   object    
     2   like_count   1757 non-null   int64     
     3   text         1757 non-null   object    
     4   public        1757 non-null   bool      
    dtypes: bool(1), int64(1), object(3)  
    memory usage: 70.3+ KB
```

- Cek lagi datanya apakah sudah benar
- Dari yang semula 1758, kini menjadi 1757 setelah PreProcessing

```
● ● ●  
1 # Menyimpan df_cleaned sebagai file CSV  
2 df.to_csv('data_cleaned_Youtube_RaymondChin.csv', index=False)  
3  
  
● Simpan datanya sebagai csv  
  
● ● ●  
1 # Menghapus semua data lama  
2 collection.delete_many({})  
3  
4 # Mengonversi DataFrame menjadi list of dictionaries  
5 data_baru = df.to_dict(orient='records')  
6  
7 # Memuat data baru ke MongoDB  
8 collection.insert_many(data_baru)  
9 print(f"Data berhasil dimuat kembali dengan {len(data_baru)} dokumen.")  
10
```

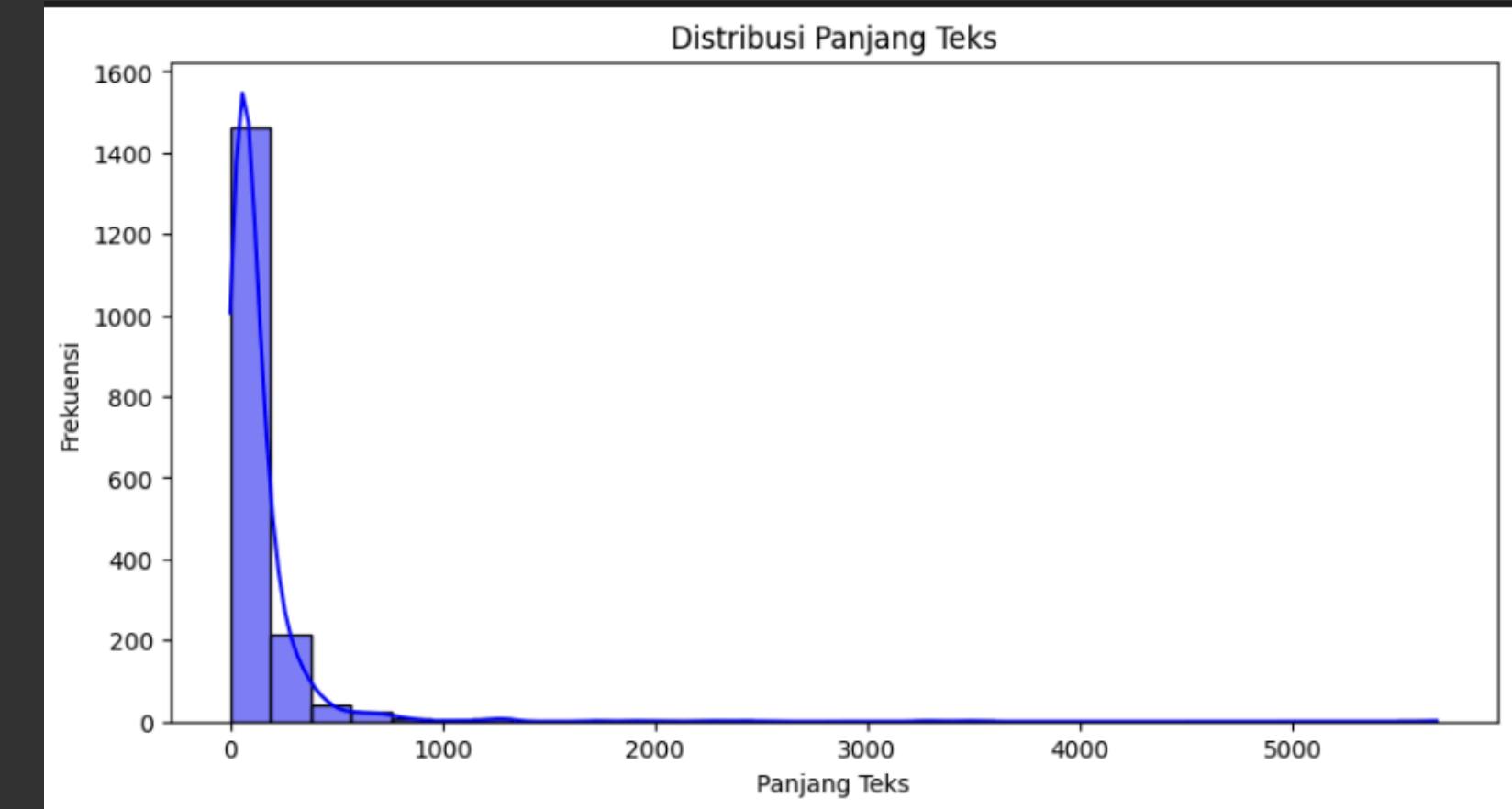
- Load lagi datanya ke mongo DB untuk perubahan data terbaru

ANALYZE & VISUALIZATION

LETS GO EXECUTE (ANALYZE & VISUALIZATION)



```
1 # Analisis panjang teks
2 df['text_length'] = df['text'].apply(lambda x: len(str(x)))
3
4
5 plt.figure(figsize=(10,5))
6 sns.histplot(df['text_length'], bins=30, kde=True, color='blue')
7 plt.xlabel("Panjang Teks")
8 plt.ylabel("Frekuensi")
9 plt.title("Distribusi Panjang Teks")
10 plt.show()
11
```



- Grafik menunjukkan bahwa sebagian besar teks dalam dataset kita memiliki panjang yang pendek, sekitar 0-500 karakter, dengan jumlah terbanyak mencapai lebih dari 1.400, sementara teks panjang di atas 1.000 karakter sangat jarang ditemukan. Hal ini menunjukkan bahwa konten kita didominasi oleh teks pendek

LETS GO EXECUTE (ANALYZE & VISUALIZATION)

```
1 # WordCloud untuk melihat kata yang sering muncul
2 text_corpus = " ".join(df['text'].dropna().astype(str))
3 # Tema warna (misalnya: biru ke ungu)
4
5
6 stopwords = set(STOPWORDS)
7 stopwords.update(["yg", "dg", "rt", "dgn", "ny", "d", 'klo',
8                  'kalo', 'amp', 'biar', 'bikin', 'bilang',
9                  'gak', 'ga', 'krn', 'nya', 'nih', 'sih',
10                 'si', 'tau', 'tdk', 'tuh', 'utk', 'ya',
11                 'jd', 'jgn', 'sdh', 'aja', 'n', 't',
12                 'nyg', 'hehe', 'pen', 'u', 'nan', 'loh', 'rt',
13                 'wkwk', 'wkwkwk', 'nggak', 'gakk', 'gk', 'tp', 'lg',
14                 'bgt', 'dr', 'gw', 'gue', 'gpp', 'aja', 'deh', 'kok',
15                 'doi', 'banget', 'bangett', 'cuma', 'kali', 'yaa',
16                 'udah', 'k', 'sy', 'm', 'bang', 'kak', 'kakak', 'min',
17                 'mas', 'mba', 'mbak', 'anjay', 'anjir', 'anj', 'lah',
18                 'lho', 'loh', 'lu', 'loe', 'lo', 'kamu', 'aku', 'gue',
19                 'gw', 'gua', 'gue', 'kita', 'mereka', 'kalian', 'aku', 'saya',
20                 'dia', 'mereka', 'kamu', 'kalian', 'kau', 'engkau', 'anda',
21                 'beliau', 'ia', 'mereka', 'kita', 'ini', 'itu', 'sana', 'sini',
22                 'mana', 'kapan', 'kenapa', 'bagaimana', 'berapa', 'siapa', 'dimana',
23                 'apa', 'siapa', 'mana', 'kapan', 'mengapa', 'bagaimana', 'berapa',
24                 'kabur', 'aja', 'dulu', 'indonesia', 'luar', 'negeri'])
25
26 wordcloud = WordCloud(width=800,height=400,background_color="black",colormap="cool",stopwords=STOPWORDS,
27                         max_words=200,contour_width=3,contour_color="white").generate(text_corpus)
28
29
30 plt.figure(figsize=(10,5))
31 plt.imshow(wordcloud, interpolation='bilinear')
32 plt.axis("off")
33 plt.title("Word Cloud Post Tiktok")
34 plt.show()
```



- berikut adalah wordcloud sebelum adanya preprocessing

LETS GO EXECUTE (ANALYZE & VISUALIZATION)



```
1 # Check the result of demojize
2
3 df['text'].apply(lambda x: emoji.demojize(x) if isinstance(x, str) else x).head()
```

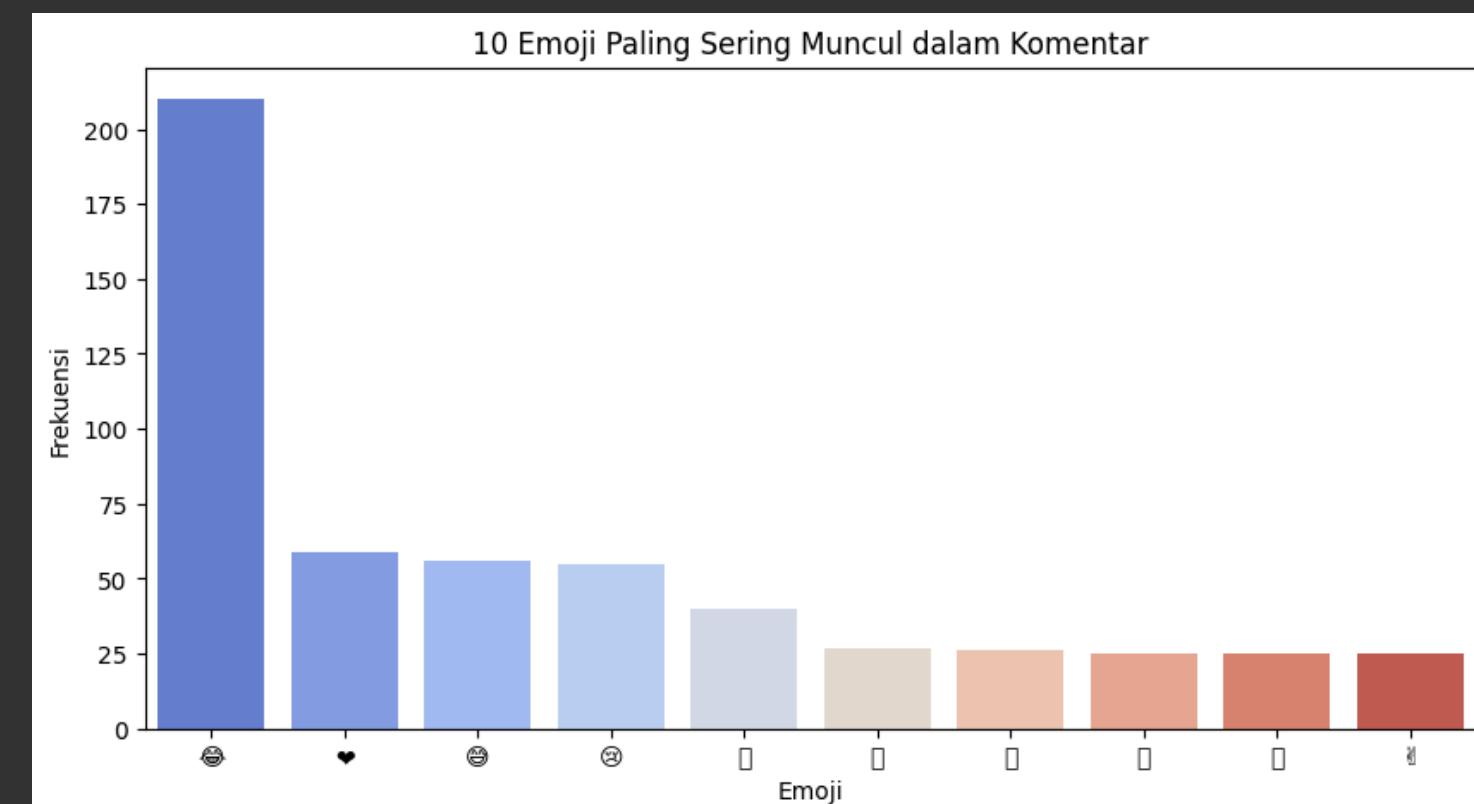
```
0 Ayo taruh an ga bakal an di korupsi kaya jiwas...
1 Semoga ada transparansi:crying_face:
2 Itulah yang ditakutkan. Eksekusi amburadul, ap...
3 DANANTARA ITU LEMBAGA NEGARA BUKAN KOH....TERU...
4 Jk berkaca dari SDM penguasa negeri ini sih sg...
```

Name: text, dtype: object

- Mengubah emoji menjadi sebuah kata, dan memvisualisasikan dalam top 10 Emoji paling sering muncul



```
1 def extract_emojis(text):
2     return [char for char in text if emoji.is_emoji(char)]
3
4 all_emojis = [emj for text in df['text'].dropna().astype(str) for emj in extract_emojis(text)]
5 emoji_freq = Counter(all_emojis)
6 top_emojis = emoji_freq.most_common(10)
7 emoji_chars, emoji_counts = zip(*top_emojis)
8
9 plt.figure(figsize=(10, 5))
10 sns.barplot(x=list(emoji_chars), y=list(emoji_counts), palette="coolwarm")
11 plt.xlabel("Emoji")
12 plt.ylabel("Frekuensi")
13 plt.title("10 Emoji Paling Sering Muncul dalam Komentar")
14 plt.show()
```



LETS GO EXECUTE (ANALYZE & VISUALIZATION)



```
1 import pandas as pd
2 import re
3 import string
4 import nltk
5 from Sastrawi.Stemmer.StemmerFactory import StemmerFactory
6 from nltk.corpus import stopwords
7 from nltk.tokenize import regexp_tokenize
```



```
1 # Ensure the NLTK stopwords list is downloaded
2 nltk.download('stopwords')
```



```
1 # Membuat stemmer
2 factory = StemmerFactory()
3 stemmer = factory.create_stemmer()
```



```
1 # Fungsi untuk menghapus URL dan HTML
2 def remove_urls_and_html(text):
3     # Menghapus URL
4     text = re.sub(r'http\S+', '', text)
5     # Menghapus tag HTML
6     text = re.sub(r'<.*?>', '', text)
7     return text
```



```
1 # Fungsi untuk menghapus emoji
2 def remove_emoji(text):
3     return emoji.replace_emoji(text, replace='')
```



```
1 # Fungsi cleansing teks
2 def cleansing(df):
3     df_clean = df.astype(str) # Mengubah semua entri menjadi string
4     df_clean = df_clean.str.lower() # Mengubah menjadi huruf kecil
5     df_clean = [re.sub(r'\d+', "", i) for i in df_clean] # Menghapus angka
6     df_clean = [re.sub(r'[^w]', ' ', i) for i in df_clean] # Menghapus karakter non-alfanumerik
7     df_clean = [re.sub(r'\s+', ' ', i) for i in df_clean] # Mengganti spasi ganda menjadi satu
8     df_clean = [i.strip() for i in df_clean] # Menghapus spasi di awal dan akhir
9     return df_clean
```



```
1 # Daftar stopwords Bahasa Indonesia dari NLTK
2 stopwords_list = stopwords.words('indonesian')
3 stopwords_list.extend([
4     "yg", "dg", "rt", "dgn", "ny", "d", "klo", "kalo", "amp", "biar", "bikin", "bilang",
5     "gak", "ga", "krn", "nya", "nih", "sih", "si", "tau", "tdk", "tuh", "utk", "ya", "jd",
6     "jgn", "sdh", "aja", "n", "t", "nyg", "hehe", "pen", "u", "nan", "loh", "rt", "wkwk",
7     "wkwk", "nggak", "gakk", "gk", "tp", "lg", "bgt", "dr", "gw", "gue", "gpp", "aja", "deh",
8     "kok", "doi", "banget", "bangett", "cuma", "kali", "yaa", "udah", "k", "sy", "m", "bang",
9     "kak", "kakak", "min", "mas", "mba", "mbak", "anjay", "anjir", "anj", "lah", "lho", "loh",
10    "lu", "loe", "lo", "kamu", "aku", "gue", "gw", "gua", "gue", "kita", "mereka", "kalian",
11    "aku", "saya", "dia", "mereka", "kamu", "kalian", "kau", "engkau", "anda", "beliau", "ia",
12    "mereka", "kita", "ini", "itu", "sana", "sini", "mana", "kapan", "kenapa", "bagaimana",
13    "berapa", "siapa", "dimana", "apa", "siapa", "mana", "kapan", "mengapa", "bagaimana",
14    "berapa", "kabur", "aja", "dulu", "indonesia", "luar", "negeri", "dan"
15  ])
```



```
1 # Fungsi untuk menghapus stopwords
2 def remove_stopwords(tokens):
3     return [word for word in tokens if word not in stopwords_list]
```

- Preprocessing lebih lanjut untuk sentiment analysis nya
- Preprocessing mencakup Stopword, Stemming dengan Sastrawi, Tokenize, Lower case, menghapus angka, URL maupun HTML, menghapus emoji, dll.

LETS GO EXECUTE (ANALYZE & VISUALIZATION)

```
● ● ●  
1 # Fungsi utama untuk memproses DataFrame  
2 def process_dataframe(df):  
3     df_clean = cleansing(df['text'])  
4  
5     # Menambahkan kolom 'clean_comment' untuk teks yang sudah dibersihkan  
6     df['clean_comment'] = df_clean  
7  
8     # Tokenisasi  
9     df['tokens'] = df_clean # Simpan hasil cleansing langsung ke kolom tokens  
10  
11    # Tokenisasi teks  
12    df['tokens'] = df['tokens'].apply(tokenize_text)  
13  
14    # Menghapus stopwords  
15    df['tokens_nostopwords'] = df['tokens'].apply(remove_stopwords)  
16  
17    return df
```

```
● ● ●  
1 # Proses DataFrame  
2 df_processed = process_dataframe(df)
```

```
● ● ●  
1 # Check the result  
2 df['nonstop_comment'] = df['tokens'].apply(lambda x: [word for word in x if word not in stopwords_list])  
3 df.head()
```

- Proses dataframe dengan membuat kolom baru bernama clean_comment, tokens, tokens_nostopwords, nonstop_comment

#	word_count	clean_comment	tokens	tokens_nostopwords	nonstop_comment
19	ayo taruh an ga bakal an di korupsi k	['ayo', 'taruh', 'an', 'ga', 'bakal', 'an', 'di', 'korupsi', 'k']	['ayo', 'taruh', 'an', 'an', 'korupsi', 'kay']	['ayo', 'taruh', 'an', 'an', 'korupsi', 'kay']	['ayo', 'taruh', 'an', 'an', 'korupsi', 'kay']
3	semoga ada transparansi	['semoga', 'ada', 'transparansi']	['semoga', 'transparansi']	['semoga', 'transparansi']	['semoga', 'transparansi']
55	itulah yang ditakutkan eksekusi amb	['itulah', 'yang', 'ditakutkan', 'eksekus']	['ditakutkan', 'eksekusi', 'amburadul']	['ditakutkan', 'eksekusi', 'amburadul']	['ditakutkan', 'eksekusi', 'amburadul']
10	danantara itu lembaga negara bukar	['danantara', 'itu', 'lembaga', 'negara']	['danantara', 'lembaga', 'negara', 'kor']	['danantara', 'lembaga', 'negara', 'kor']	['danantara', 'lembaga', 'negara', 'kor']
10	jk berkaca dari sdm penguasa negeri	['jk', 'berkaca', 'dari', 'sdm', 'penguasa']	['jk', 'berkaca', 'sdm', 'penguasa', 'sgt']	['jk', 'berkaca', 'sdm', 'penguasa', 'sgt']	['jk', 'berkaca', 'sdm', 'penguasa', 'sgt']

LETS GO EXECUTE (ANALYZE & VISUALIZATION)

```
● ● ●  
1 # Fungsi utama untuk memproses DataFrame  
2 def process_dataframe(df):  
3     df_clean = cleansing(df['text'])  
4  
5     # Menambahkan kolom 'clean_comment' untuk teks yang sudah dibersihkan  
6     df['clean_comment'] = df_clean  
7  
8     # Tokenisasi  
9     df['tokens'] = df_clean # Simpan hasil cleansing langsung ke kolom tokens  
10  
11    # Tokenisasi teks  
12    df['tokens'] = df['tokens'].apply(tokenize_text)  
13  
14    # Menghapus stopwords  
15    df['tokens_nostopwords'] = df['tokens'].apply(remove_stopwords)  
16  
17    return df
```

```
● ● ●  
1 # Proses DataFrame  
2 df_processed = process_dataframe(df)
```

```
● ● ●  
1 # Check the result  
2 df['nonstop_comment'] = df['tokens'].apply(lambda x: [word for word in x if word not in stopwords_list])  
3 df.head()
```

- Proses dataframe dengan membuat kolom baru bernama clean_comment, tokens, tokens_nostopwords, nonstop_comment

#	word_count	clean_comment	tokens	tokens_nostopwords	nonstop_comment
19	ayo taruh an ga bakal an di korupsi k	['ayo', 'taruh', 'an', 'ga', 'bakal', 'an', 'di', 'korupsi', 'k']	['ayo', 'taruh', 'an', 'an', 'korupsi', 'kay']	['ayo', 'taruh', 'an', 'an', 'korupsi', 'kay']	['ayo', 'taruh', 'an', 'an', 'korupsi', 'kay']
3	semoga ada transparansi	['semoga', 'ada', 'transparansi']	['semoga', 'transparansi']	['semoga', 'transparansi']	['semoga', 'transparansi']
55	itulah yang ditakutkan eksekusi amb	['itulah', 'yang', 'ditakutkan', 'eksekus']	['ditakutkan', 'eksekusi', 'amburadul']	['ditakutkan', 'eksekusi', 'amburadul']	['ditakutkan', 'eksekusi', 'amburadul']
10	danantara itu lembaga negara bukar	['danantara', 'itu', 'lembaga', 'negara']	['danantara', 'lembaga', 'negara', 'kor']	['danantara', 'lembaga', 'negara', 'kor']	['danantara', 'lembaga', 'negara', 'kor']
10	jk berkaca dari sdm penguasa negeri	['jk', 'berkaca', 'dari', 'sdm', 'penguasa']	['jk', 'berkaca', 'sdm', 'penguasa', 'sgt']	['jk', 'berkaca', 'sdm', 'penguasa', 'sgt']	['jk', 'berkaca', 'sdm', 'penguasa', 'sgt']

LETS GO EXECUTE (ANALYZE & VISUALIZATION)

1

```
1 import matplotlib.pyplot as plt
2 from wordcloud import WordCloud
3 from nltk.corpus import stopwords
4
5 # Make sure stopwords are downloaded
6 nltk.download('stopwords')
7
8 # Get the stopwords from NLTK and convert it to a set
9 stopwords_list = set(stopwords.words('indonesian')) # Convert to a set to avoid issues with iteration
10
11 # Combine all the tokens (without stopwords)
12 text_corpus_preprocessed = " ".join([" ".join(tokens) for tokens in df['tokens_nostopwords']])
13
14 # Create the word cloud
15 wordcloud_preprocessed = WordCloud(
16     width=800, height=400, background_color="black", colormap="cool",
17     stopwords=stopwords_list, max_words=200, contour_width=3, contour_color="white"
18 ).generate(text_corpus_preprocessed)
19
20 # Plot the wordcloud
21 plt.figure(figsize=(10,5))
22 plt.imshow(wordcloud_preprocessed, interpolation='bilinear')
23 plt.axis("off")
24 plt.title("Word Cloud Comments Youtube (Preprocessed)")
25 plt.show()
```

Word Cloud Comments Youtube (Preprocessed)



- Wordcloud setelah adanya preprocess lanjutan

LETS GO EXECUTE (ANALYZE & VISUALIZATION)



```
1 # Kamus Leksikon positif dan negatif
2 positive_url = "https://raw.githubusercontent.com/fajri91/InSet/master/positive.tsv"
3 negative_url = "https://raw.githubusercontent.com/fajri91/InSet/master/negative.tsv"
4
5 pos_lex = set(pd.read_csv(positive_url, sep='\t', header=None)[0].str.lower())
6 neg_lex = set(pd.read_csv(negative_url, sep='\t', header=None)[0].str.lower())
7
8 # Fungsi untuk menentukan sentimen
9 def determine_sentiment(text):
10     if isinstance(text, str):
11         positive_count = sum(1 for word in text.split() if word in pos_lex)
12         negative_count = sum(1 for word in text.split() if word in neg_lex)
13         sentiment_score = positive_count - negative_count
14
15         if sentiment_score > 0:
16             sentiment = "Positif"
17         elif sentiment_score < 0:
18             sentiment = "Negatif"
19         else:
20             sentiment = "Netral"
21
22         return sentiment_score, sentiment
23
24     return 0, "Netral"
25
26 # Tentukan sentimen dan skor untuk setiap ulasan
27 df[['Score', 'Sentiment']] = df['clean_comment'].apply(lambda x: pd.Series(determine_sentiment(x)))
28
29 # Tampilkan hasilnya
30 print(df.head(5))
31
```

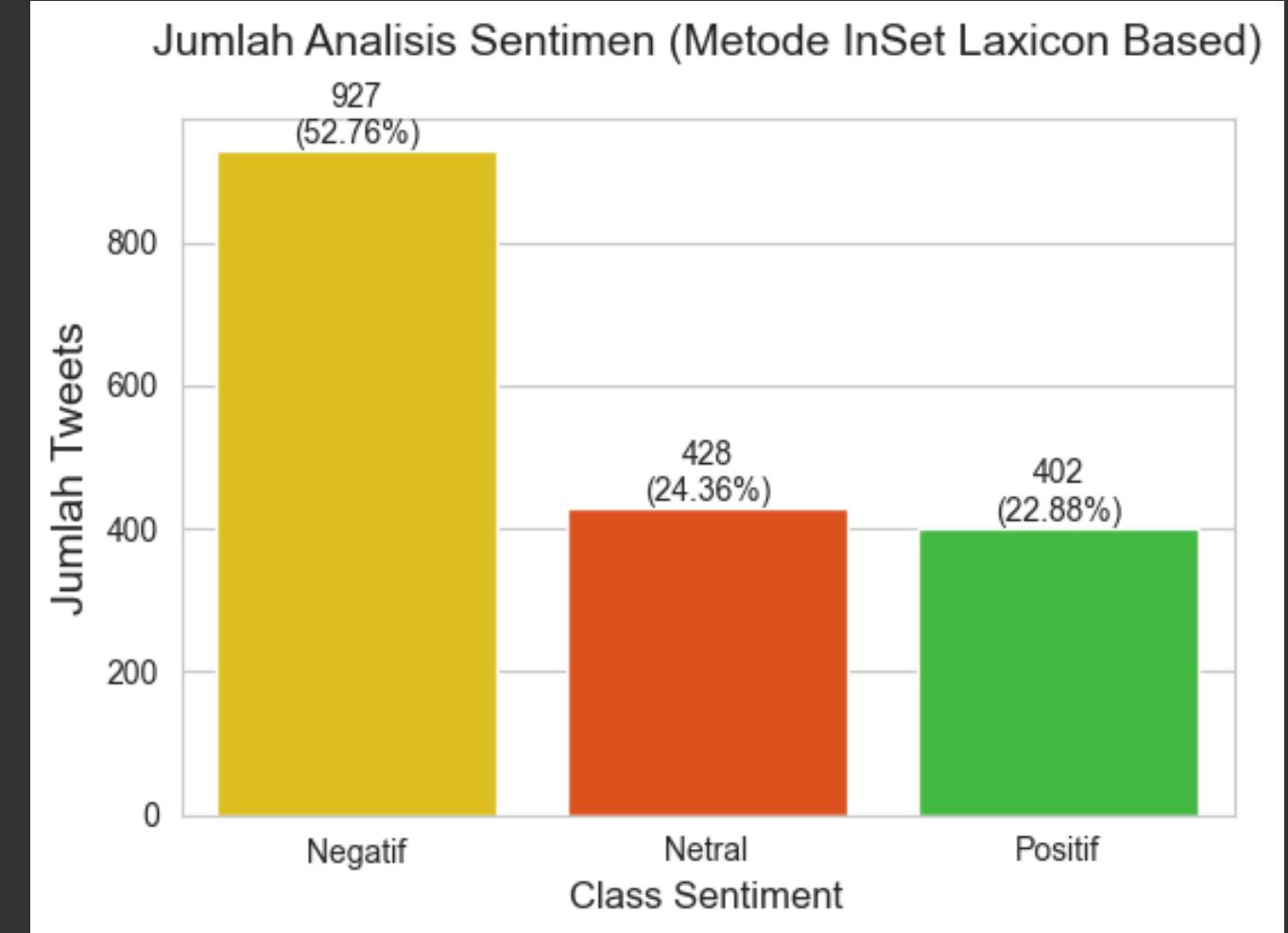
```
...
1 [semoga, transparansi] 1 Positif
2 [ditakutkan, eksekusi, amburadul, pemerintah, ... -16 Negatif
3 [danantara, lembaga, negara, koh, pengelolanya... -3 Negatif
4 [jk, berkaca, sdm, penguasa, sgt, pesimis] 0 Netral
```

- Ambil kamus Inset Lexicon Based untuk Positif dan Negatif dari Github

LETS GO EXECUTE (ANALYZE & VISUALIZATION)



```
1 sentiment_count = df['Sentiment'].value_counts()
2 sns.set_style('whitegrid')
3
4
5 custom_colors = ["#FFD700", "#FF4500", "#32CD32"]
6 fig, ax = plt.subplots(figsize=(6, 4))
7 ax = sns.barplot(x=sentiment_count.index, y=sentiment_count.values, palette=custom_colors)
8
9
10 plt.title("Jumlah Analisis Sentimen (Metode InSet Laxicon Based)", fontsize=14, pad=20)
11 plt.xlabel("Class Sentiment", fontsize=12)
12 plt.ylabel("Jumlah Tweets", fontsize=14)
13
14
15 total = len(df['Sentiment'])
16 for i, count in enumerate(sentiment_count.values):
17     percentage = f"{100 * count / total:.2f}%"
18     ax.text(i, count + 0.10, f"{count}\n({percentage})", ha='center', va='bottom')
19
20 plt.show()
```



- Visualisasi Hasilnya, dan seperti yang tertera pada diagram bahwa Negatif berada pada posisi paling tinggi.
- Menandakan bahwa Danantara pada Video Youtube Raymond Chin memiliki konotasi sentimen negatif di kalangan masyarakat.

MACHINE LEARNING MODEL

LETS GO EXECUTE (MACHINE LEARNING MODEL)



```
1 from sklearn.model_selection import train_test_split
2 # Pembagian data
3 x_train, x_test, y_train, y_test = train_test_split(df['clean_comment'], df['Sentiment'], test_size=0.2, random_state=42)
4
```



1 x_train



1 y_train



1 x_test



1 y_test

- Membuat model machine learning dengan Train, Test Split terlebih dahulu

LETS GO EXECUTE (MACHINE LEARNING MODEL)

Teks Representation 1 : Bag_of Words

```
from sklearn.feature_extraction.text import CountVectorizer
vectorizer = CountVectorizer()
x = vectorizer.fit_transform(x_train)
train_bow = pd.DataFrame(x.toarray(), columns=vectorizer.get_feature_names_out())

# Menampilkan hasil
print("Representasi Bag of Words:")
print(train_bow)
✓ 0.0s ━ Open 'train_bow' in Data Wrangler
```

Representasi Bag of Words:

	_jujur	aada	aamiiin	aamiin	aamin	aarrrrrggghhhhhhhhhh	abadi	\
0	0	0	0	0	0	0	0	
1	0	0	0	0	0	0	0	
2	0	0	0	0	0	0	0	
3	0	0	0	0	0	0	0	
4	0	0	0	0	0	0	0	
...	
1400	0	0	0	0	0	0	0	
1401	0	0	0	0	0	0	0	
1402	0	0	0	0	0	0	0	
1403	0	0	0	0	0	0	0	
1404	0	0	0	0	0	0	0	

```
train_bow.columns
✓ 0.0s Python
Index(['_jujur', 'aada', 'aamiiin', 'aamiin', 'aamin',
       'aarrrrrggghhhhhhhhhh', 'abadi', 'abah', 'abang', 'abangnya',
       ...
       'bedah', 'ideologi', 'kita', 'mas', 'tolong', 'AXLbikin', 'AXL', 'AXL',
       'DEWadOrA', 'DEWaDOяA'],
      dtype='object', length=5877)

x_test_bow = vectorizer.transform(x_test)
test_bow = pd.DataFrame(x_test_bow.toarray(), columns=vectorizer.get_feature_names_out())
✓ 0.0s Python
```

- Text Representation Bag of Words dengan 3 Model Machine Learning (Random Forest, SVM, dan Naive Bayes)

LETS GO EXECUTE (MACHINE LEARNING MODEL)

Teks Representation 1 : Bag_of_Words

The screenshot shows a Jupyter Notebook interface with three code cells labeled 1. Random Forest, 2. SVM, and 3. Naive Bayes.

1. Random Forest:

```
from sklearn.ensemble import RandomForestClassifier
rf_class = RandomForestClassifier(n_estimators=100, random_state=42)
rf_class.fit(train_bow, y_train)

test_rf_class = rf_class.predict(test_bow)

from sklearn.metrics import classification_report
print('\nRandom Forest Classification Report\n')
print(classification_report(y_test, test_rf_class, target_names=['Negatif','Netral','Positif']))
```

Execution time: 2.0s

2. SVM:

```
from sklearn import svm
svm_class1 = svm.LinearSVC( random_state=42)
svm_class1.fit(train_bow, y_train)

test_svm_class1=svm_class1.predict(test_bow)

print('\nClassification Report\n')
print(classification_report(y_test, test_svm_class1, target_names=['Negatif','Netral','Positif']))
```

Execution time: 0.0s

3. Naive Bayes:

```
from sklearn.naive_bayes import MultinomialNB
nb_class = MultinomialNB()
nb_class.fit(train_bow, y_train)
test_nb_class = nb_class.predict(test_bow)

print('\nNaive Bayes Classification Report\n')
print(classification_report(y_test, test_nb_class, target_names=['Negatif','Netral','Positif']))
```

Execution time: 0.1s

Classification Reports:

Random Forest Classification Report		Classification Report		Naive Bayes Classification Report	
		precision	recall	f1-score	support
Negatif	0.73	0.90	0.80	0.77	184
Netral	0.61	0.60	0.60	0.25	92
Positif	0.83	0.39	0.54	0.54	76
accuracy			0.71	0.71	352
macro avg	0.72	0.63	0.65	0.64	352
weighted avg	0.72	0.71	0.69	0.64	352

	precision	recall	f1-score	support	
Negatif	0.83	0.81	0.82	184	
Netral	0.53	0.60	0.56	92	
Positif	0.71	0.63	0.67	76	
accuracy			0.72	0.72	352
macro avg	0.69	0.68	0.68	0.52	352
weighted avg	0.72	0.72	0.72	0.64	352

LETS GO EXECUTE (MACHINE LEARNING MODEL)

Teks Representation 2 : TF-IDF

```
from sklearn.feature_extraction.text import TfidfVectorizer
tfidf_vectorizer = TfidfVectorizer()
X = tfidf_vectorizer.fit_transform(x_train)
train_tfidf=pd.DataFrame(X.toarray(),columns=vectorizer.get_feature_names_out())
train_tfidf.head()

✓ 0.4s └ Open 'train_tfidf' in Data Wrangler
```

	#_jujur	#_aada	#_aamiin	#_aamiin
0	0.0	0.0	0.0	0.0
1	0.0	0.0	0.0	0.0
2	0.0	0.0	0.0	0.0
3	0.0	0.0	0.0	0.0
4	0.0	0.0	0.0	0.0

5 rows x 5,877 cols 10 ▾ per page << < Page 1 of 1 > >


```
test_tfidf = tfidf_vectorizer.transform(x_test)
TFIDF_test=pd.DataFrame(test_tfidf.toarray(),columns=vectorizer.get_feature_names_out())
TFIDF_test.head()

✓ 0.5s └ Open 'TFIDF_test' in Data Wrangler
```

	#_jujur	#_aada	#_aamiin	#_aamiin
0	0.0	0.0	0.0	0.0
1	0.0	0.0	0.0	0.0
2	0.0	0.0	0.0	0.0
3	0.0	0.0	0.0	0.0
4	0.0	0.0	0.0	0.324

- Text Reprentation TF - IDF dengan 3 Model Machine Learning (Random Forest, SVM, dan Naive Bayes)

LETS GO EXECUTE (MACHINE LEARNING MODEL)

Teks Representation 2 : TF-IDF

1. Random Forest

```
from sklearn.ensemble import RandomForestClassifier
rf_class2 = RandomForestClassifier(n_estimators=100, random_state=42)
rf_class2.fit(train_tfidf, y_train)

test_rf_class2 = rf_class2.predict(test_tfidf)
print('\nRandom Forest Classification TFIDF Report\n')
print(classification_report(y_test, test_rf_class2, target_names=['Negatif','Netral','Positif']))

```

	precision	recall	f1-score	support
Negatif	0.70	0.91	0.79	184
Netral	0.59	0.52	0.55	92
Positif	0.81	0.34	0.48	76

2. SVM

```
svm_class2 = svm.LinearSVC( random_state=42)
svm_class2.fit(train_tfidf, y_train)

test_svm_class2 = svm_class2.predict(test_tfidf)
print('\nClassification Report SVM TFIDF\n')
print(classification_report(y_test, test_svm_class2, target_names=['Negatif','Netral','Positif']))

Classification Report SVM TFIDF
precision    recall    f1-score   support
Negatif      0.77     0.88     0.82      184
Netral       0.59     0.46     0.52      92
Positif      0.75     0.70     0.72      76
```

3. Naive Bayes

```
from sklearn.naive_bayes import MultinomialNB
nb_class2 = MultinomialNB()
nb_class2.fit(train_tfidf, y_train)
test_nb_class2 = nb_class2.predict(test_tfidf)

print('\nNaive Bayes Classification Report\n')
print(classification_report(y_test, test_nb_class2, target_names=['Negatif','Netral','Positif']))

Naive Bayes Classification Report
precision    recall    f1-score   support
Negatif      0.55     1.00     0.71      184
Netral       0.75     0.03     0.06       92
Positif      1.00     0.18     0.31       76
accuracy          0.57
macro avg      0.77     0.41     0.36      352
weighted avg   0.70     0.57     0.45      352
```

LETS GO EXECUTE (MACHINE LEARNING MODEL)

Best Model Accuracy

```
● ● ●  
1 def find_best_model(y_test, models_predictions):  
2     best_model = None  
3     best_accuracy = -1  
4  
5     for model_name, predictions in models_predictions.items():  
6         report = classification_report(y_test, predictions, target_names=['Negatif', 'Netral', 'Positif'], output_dict=True)  
7         accuracy = report['accuracy']  
8  
9         if accuracy > best_accuracy:  
10             best_accuracy = accuracy  
11             best_model = model_name  
12  
13     return best_model, best_accuracy  
14  
15  
16 models_predictions = {  
17     'RF_BOW': test_rf_class,  
18     'SVM_BOW': test_svm_class1,  
19     'NB_BOW': test_nb_class,  
20     'RF_TFIDF': test_rf_class2,  
21     'SVM_TFIDF': test_svm_class2,  
22     'NB_TFIDF': test_nb_class2,  
23 }  
24  
25 best_model, best_accuracy = find_best_model(y_test, models_predictions)  
26 print(f"The best performing model is: {best_model} with accuracy: {best_accuracy:.2%}")
```

```
The best performing model is: SVM_TFIDF with accuracy: 73.01%
```

- Best model Accuracy ada di model TF - IDF (SVM) dengan score : 73,01%
- Dengan kata lain, perlu adanya peningkatan dan revisi pada program

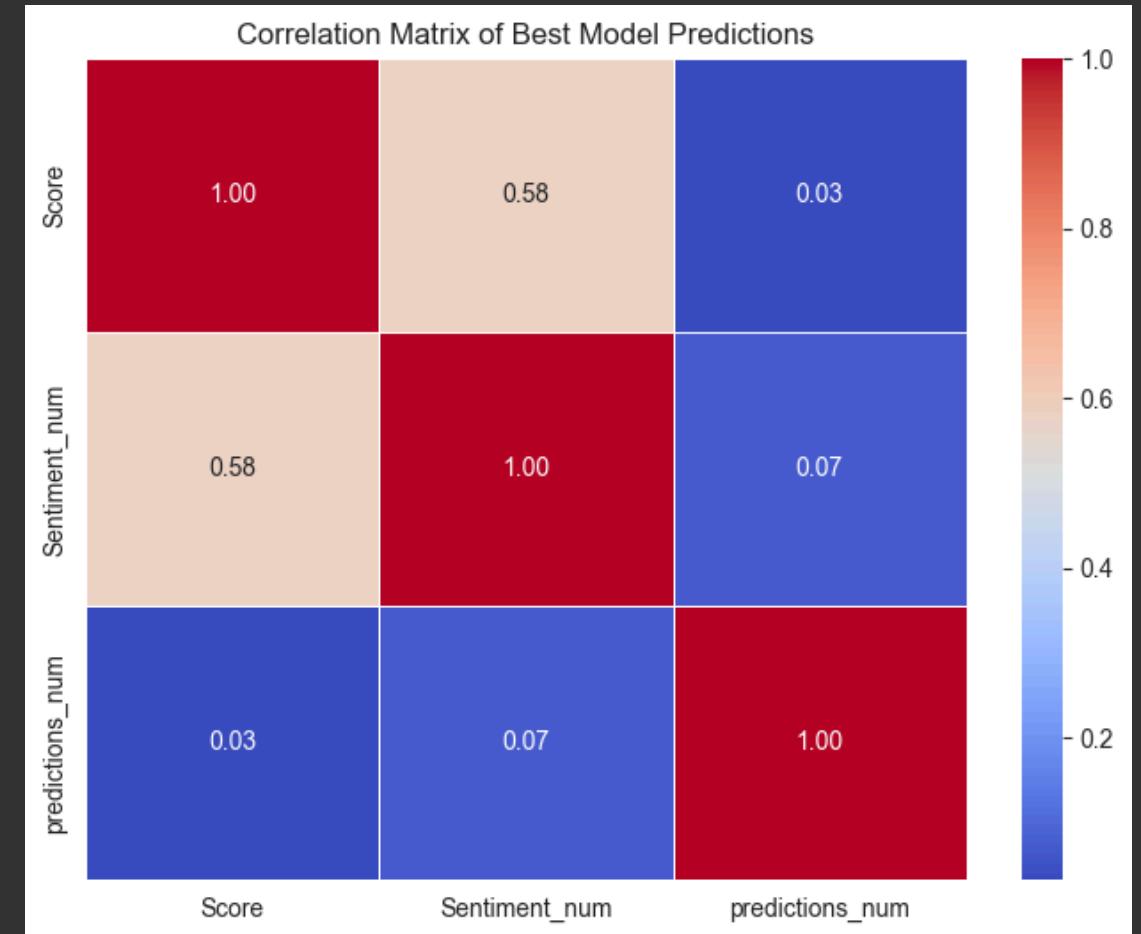


INSIGHT

INSIGHT

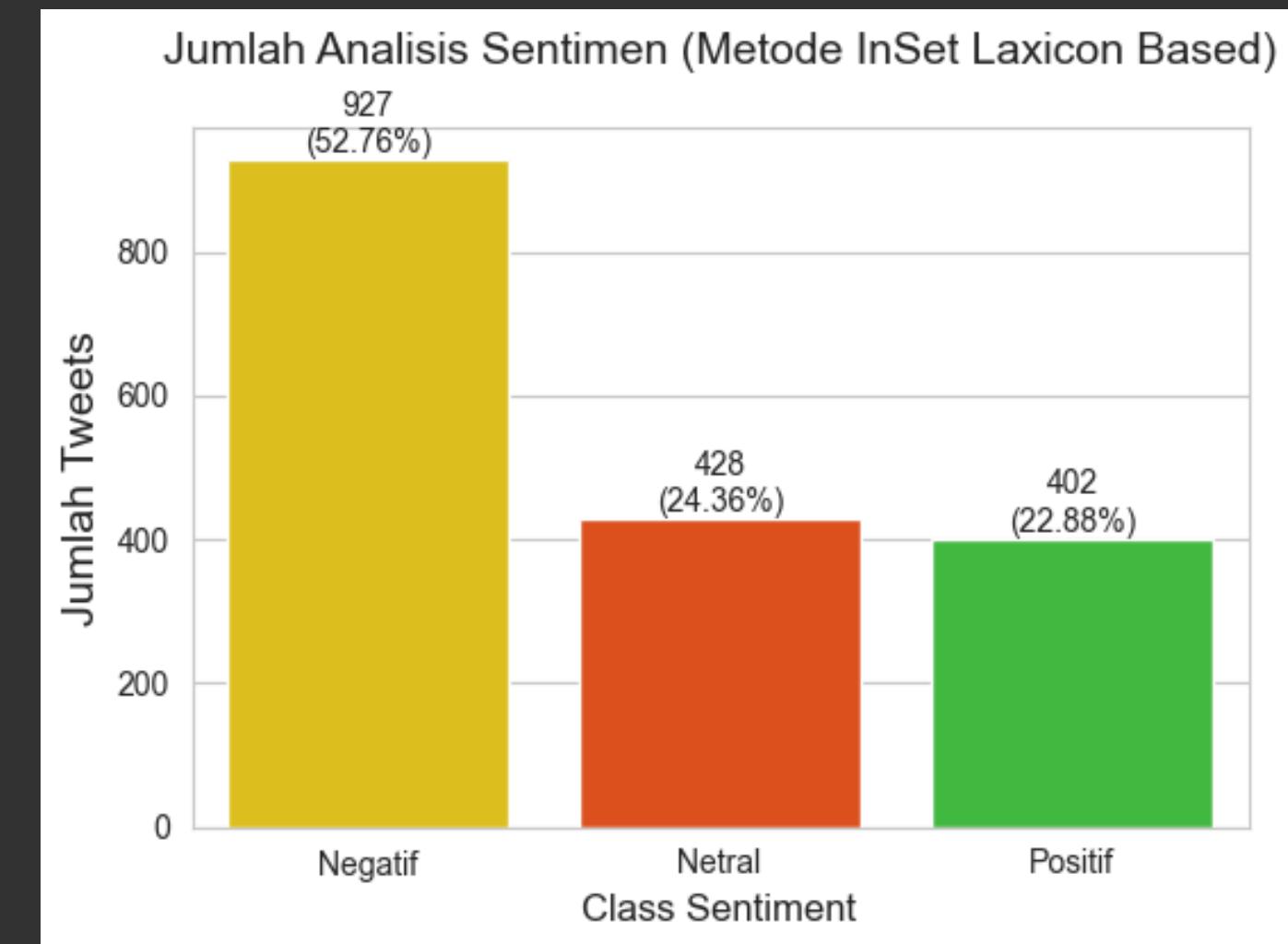
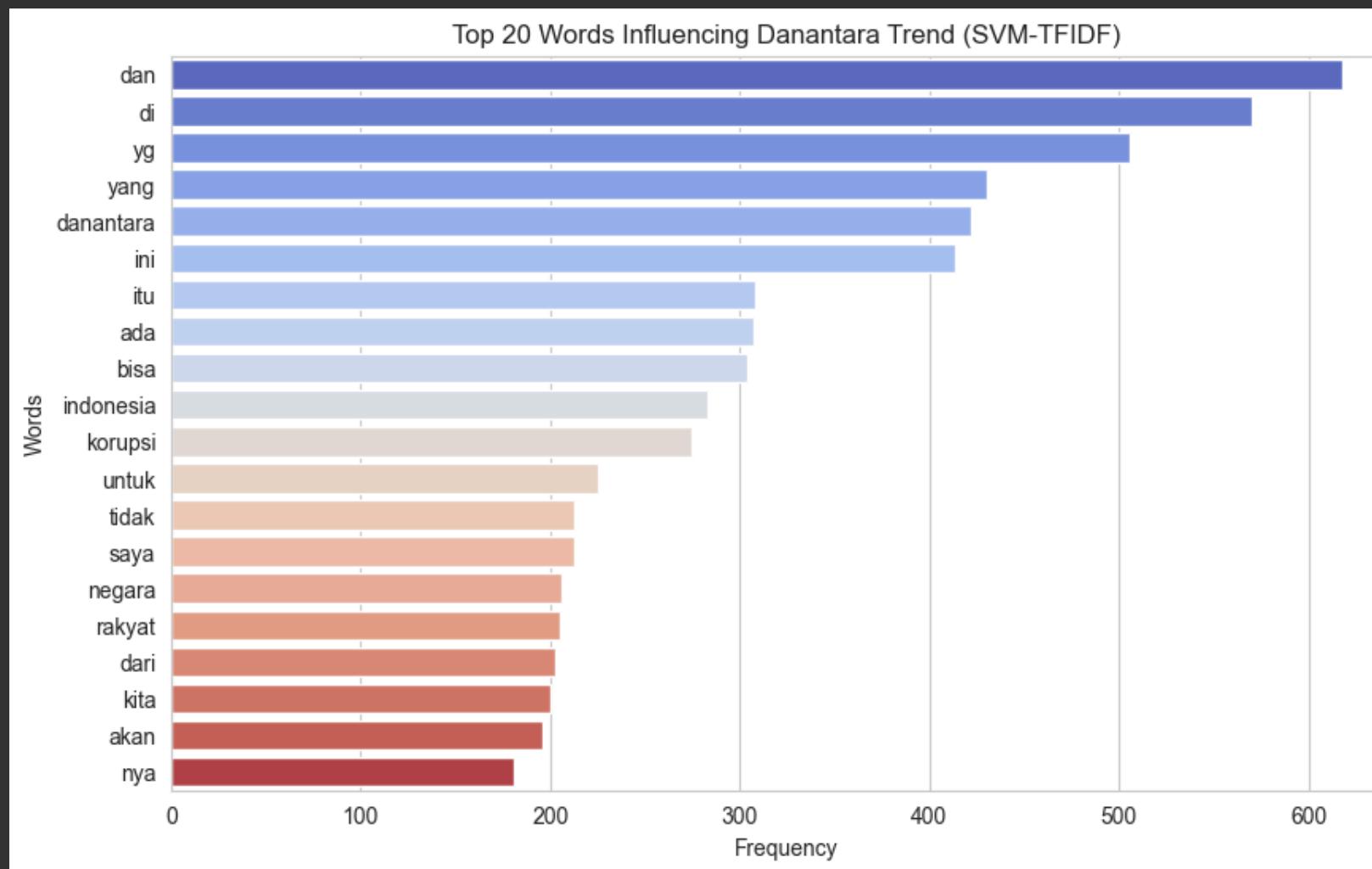


```
1 best_model_predictions = pd.DataFrame({'predictions': test_svm_class2})
2 combined_df = pd.concat([df, best_model_predictions], axis=1)
3
4 sentiment_mapping = {'Negatif': 0, 'Netral': 1, 'Positif': 2}
5 combined_df['Sentiment_num'] = combined_df['Sentiment'].map(sentiment_mapping)
6 combined_df['predictions_num'] = combined_df['predictions'].map(sentiment_mapping)
7
8 correlation_matrix = combined_df[['Score', 'Sentiment_num', 'predictions_num']].corr()
9
10 plt.figure(figsize=(8, 6))
11 sns.heatmap(correlation_matrix, annot=True, cmap='coolwarm', fmt=".2f", linewidths=.5)
12 plt.title('Correlation Matrix of Best Model Predictions')
13 plt.show()
```



- Berdasarkan Correlation Matrix tersebut, terdapat korelasi yang cukup kuat antara Score dan Sentiment_num sebesar 0.58, menunjukkan bahwa nilai skor memiliki hubungan positif dengan sentimen asli. Namun, korelasi antara predictions_num dan Sentiment_num hanya 0.07, menandakan bahwa **model masih memiliki akurasi rendah dalam memprediksi sentimen sesuai data asli**.
- Insight: **Model perlu ditingkatkan**, terutama dalam mendeteksi sentimen positif dan netral, dengan optimasi fitur atau pemilihan model yang lebih kompleks.

INSIGHT



Analisis top 20 words tentang Danantara Indonesia menunjukkan bahwa masyarakat memandang Danantara sebagai agen perubahan dalam mendorong investasi strategis dan pemulihan ekonomi nasional. Namun, hasil analisis sentimen dengan metode InSet Lexicon Based menunjukkan mayoritas sentimen yang muncul adalah negatif (52.76%), menandakan adanya kekhawatiran masyarakat terhadap efektivitas dan transparansi pengelolaan investasi. Kata-kata seperti korupsi, tidak, dan masalah memperkuat adanya ketidakpercayaan terhadap tata kelola yang bersih. Meskipun demikian, kehadiran kata bisa, manfaat, dan kita menunjukkan harapan bahwa Danantara mampu membawa dampak positif bagi kesejahteraan masyarakat jika dikelola dengan baik dan transparan. Insight ini mencerminkan perlunya peningkatan komunikasi dan transparansi dalam menjawab kekhawatiran publik serta memperkuat peran Danantara sebagai agen perubahan yang inklusif dan berkelanjutan.

THANK YOU