

# **Analisis Sentimen Komen Instagram**

**"Kenapa Elektabilitas Anies Tak Kunjung Naik?"**

## **KELOMPOK 3**

- ANDREW LAUWIRA - 2440035904**
- BELVA FITHRIYAH - 2440084231**
- FARRAH A MAHARANI - 2440109732**
- SHAREN IVANA - 2440086155**



# I.D.E

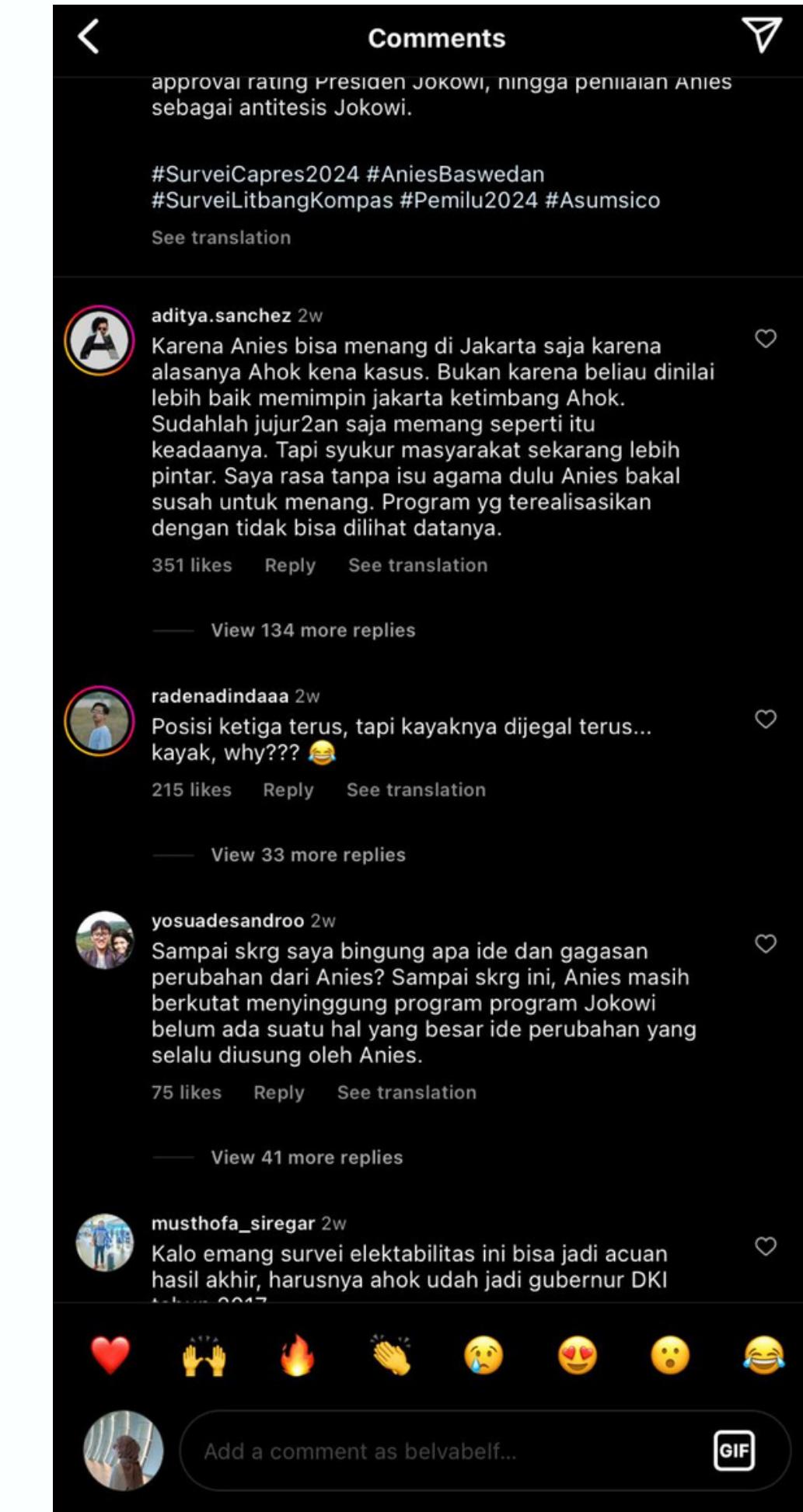
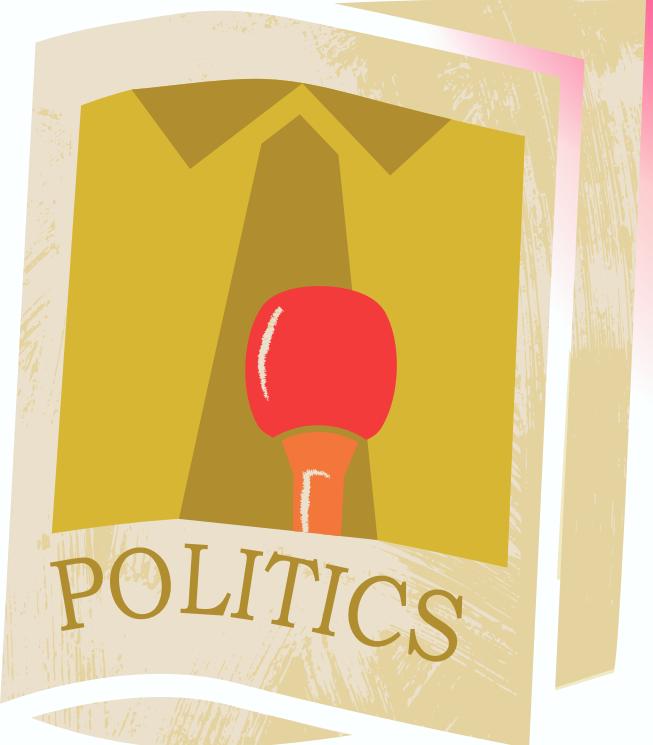
A thin vertical white line runs along the right edge of the slide.

# POLITIK

DI DUNIA POLITIK BEGITU BANYAK PRO KONTRA AKAN SUATU HAL, TERMASUK SUATU PERNYATAAN. PENTINGNYA UNTUK KITA MENCoba UNTUK MENGANALISA TENTANG PERNYATAAN-PERNYATAAN TERSEBUT UNTUK DAPAT MENGETAHUI SISI POSITIF DAN NEGATIFNYA. OLEH KARENA ITU, KITA MELAKUKAN ANALISIS SENTIMEN UNTUK DAPAT MEMBEDAKAN PERNYATAAN POSITIF DAN NEGATIF DARI KOMEN DI INSTAGRAM.



# INSTAGRAM @ASUMSICO





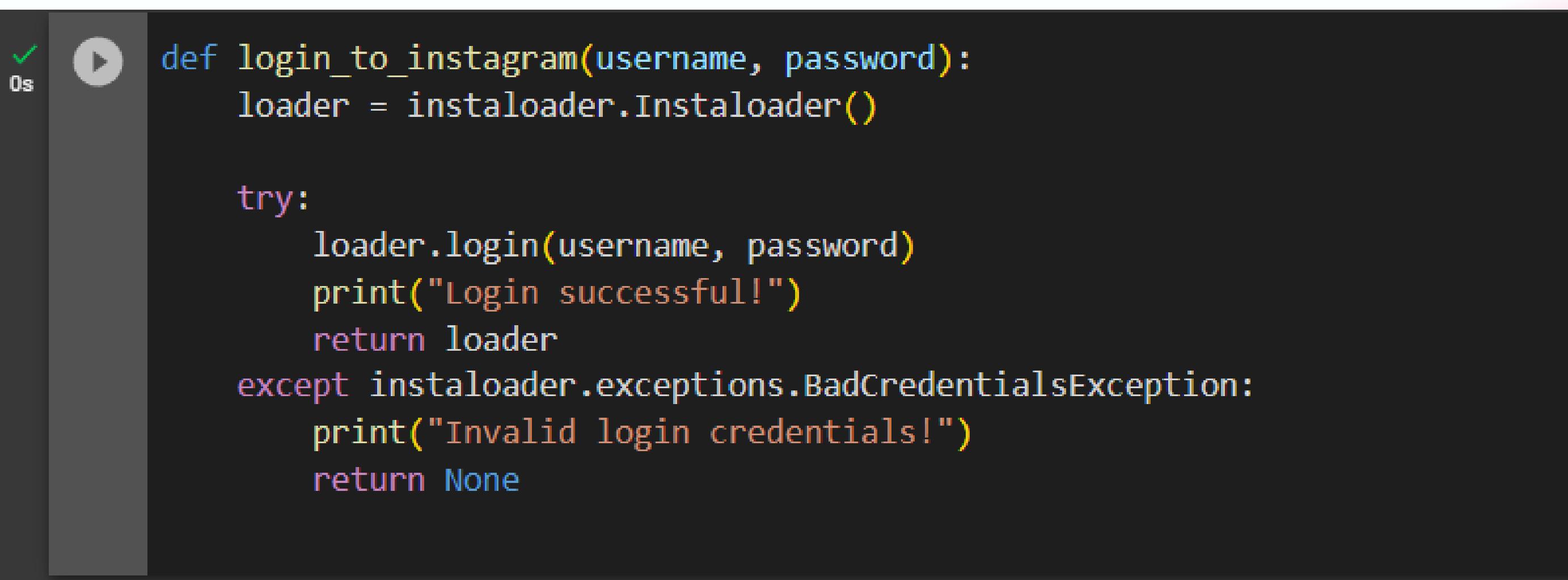
# SCRAPING KOMEN INSTAGRAM KE CSV

```
!pip install instaloader
import instaloader
import random
import time
import csv

[+] Looking in indexes: https://pypi.org/simple, https://us-python.pkg.dev/colab-wheels/public/simple/
Requirement already satisfied: instaloader in /usr/local/lib/python3.10/dist-packages (4.9.6)
Requirement already satisfied: requests>=2.4 in /usr/local/lib/python3.10/dist-packages (from instaloader) (2.27.1)
Requirement already satisfied: urllib3<1.27,>=1.21.1 in /usr/local/lib/python3.10/dist-packages (from requests>=2.4->instaloader) (1.26.15)
Requirement already satisfied: certifi>=2017.4.17 in /usr/local/lib/python3.10/dist-packages (from requests>=2.4->instaloader) (2022.12.7)
Requirement already satisfied: charset-normalizer~=2.0.0 in /usr/local/lib/python3.10/dist-packages (from requests>=2.4->instaloader) (2.0.12)
Requirement already satisfied: idna<4,>=2.5 in /usr/local/lib/python3.10/dist-packages (from requests>=2.4->instaloader) (3.4)
```

## IMPORTING NECESSARY LIBRARIES

- **Instaloader:** **instaloader adalah library Python yang digunakan untuk mendownload dan berinteraksi dengan data-data Instagram.**
- **Random:** **Library random adalah library bawaan Python yang menyediakan fungsi untuk mengenerate angka acak.**
- **Time:** **Modul time adalah modul bawaan Python yang menyediakan fungsi untuk bekerja dengan waktu.**
- **CSV:** **Modul csv adalah modul bawaan Python yang digunakan untuk membaca dan menulis file CSV (Comma-Separated Values).**



The screenshot shows a code editor window with a dark theme. In the top-left corner, there is a green checkmark icon and the text '0s'. A play button icon is also visible. The main area contains the following Python code:

```
def login_to_instagram(username, password):
    loader = instaloader.Instaloader()

    try:
        loader.login(username, password)
        print("Login successful!")
        return loader
    except instaloader.exceptions.BadCredentialsException:
        print("Invalid login credentials!")
        return None
```

## DEFINING FUNCTION

Fungsi `login_to_instagram` menerima dua parameter, yaitu `username` dan `password`, dan bertugas untuk melakukan proses login ke Instagram menggunakan library `instaloader`.

```
✓ 0s  def scrape_instagram_post(loader, shortcode):
    if loader is None:
        loader.login(username, password)
        print("Login successful!")
        return loader

    with open('SEresult.csv', 'a+', newline='', encoding='utf-8') as file:
        writer = csv.writer(file)
        try:
            post = instaloader.Post.from_shortcode(loader.context, shortcode)

            # Access the post data
            print("Post ID:", post.mediaid)
            time.sleep(random.uniform(2, 3))
            postid = post.mediaid
            print("Caption:", post.caption)
            time.sleep(random.uniform(2, 3))
            caption = post.caption
            print("Likes:", post.likes)
            time.sleep(random.uniform(2, 3))
            likes = post.likes
            writer.writerow(["post.ID", "Caption", "Likes"])
            writer.writerow([postid, caption, likes])

            # Scrape and print the comments
            print("Comments:")
            writer.writerow(["Username", "Comment"])
            for comment in post.get_comments():
                print("Username:", comment.owner.username)
                time.sleep(random.uniform(2, 3))
                username=comment.owner.username
                print("Comment:", comment.text)
                time.sleep(random.uniform(2, 3))
                comment_user=comment.text
                print("---")
                writer.writerow([username, comment_user])
        except instaloader.exceptions.InstaloaderException as e:
            print("Error occurred while scraping the post:", str(e))
```

# DEFINING FUNCTION

## Fungsi `scrape_instagram_post`

**menerima dua parameter, yaitu `loader` dan `shortcode`, dan bertanggung jawab untuk melakukan pengambilan data dari sebuah postingan Instagram menggunakan `instaloader`.**

# DEFINING FUNCTION

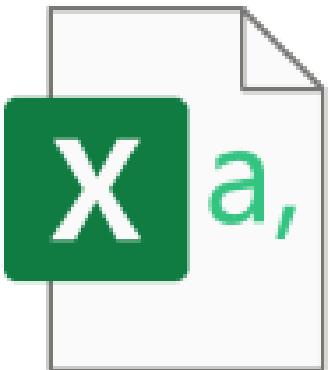
```
shortcode = 'CsqkhP3vXz6'

username = 'sharenbaubabi'
password = 'ranicantiksharenbabi'

loader = login_to_instagram(username, password)
time.sleep(random.uniform(2, 3)) # Sleep for a random duration between 2 and 3 seconds
scrape_instagram_post(loader, shortcode)
```

- **shortcode = 'CsqkhP3vXz6'**: Baris ini mengassign shortcode sebenarnya dari postingan Instagram ke dalam variabel shortcode.
- **username & password for login credentials**: Username dan pass akun ig
- **time.sleep(random.uniform(2, 3))**: Baris ini menggunakan fungsi time.sleep() dari modul time untuk memasukkan jeda acak antara 2 hingga 3 detik sebelum menjalankan instruksi berikutnya. Hal ini berguna untuk menambahkan penundaan acak agar program terlihat lebih realistik dan menghindari terlalu cepatnya permintaan ke server Instagram.

# HASIL SCRAPING:



## SEresult.csv



# Analisis Sentimen

POSITIF / NEGATIF

# DATA PREPROCESSING

## 1. Import Library

```
1 import nltk  
2 import pandas as pd  
3 import numpy as np  
4 import string  
5 import re
```

## 2. Case Folding

```
[5] 1 df = df.iloc[5:]  
2 df = df.drop(['post.ID','Likes'], axis =1)  
3 df.head(10)
```

```
▶ 1 print('Case Folding Result : \n')  
2 print(df['Caption'].head(5))  
3 print('\n\n\n\n')
```

↪ Case Folding Result :

```
5 Karena Anies bisa menang di Jakarta saja karen...  
6 Posisi ketiga terus, tapi kayaknya dijegal ter...  
7 Hasil survei litbang kompas sebelum pilkada dk...  
8 Heh,wowkwok waktu pilgub jakarta juga peringk...  
9 Menang survei kok bangga.. liat aja dilapangan...  
Name: Caption, dtype: object
```

### 3. Tokenizing

```
# import word_tokenize & FreqDist from NLTK
from nltk.tokenize import word_tokenize
from nltk.probability import FreqDist

# ----- Tokenizing -----

def remove_special(text):
    text = str(text)
    # remove tab, new line, ans back slice
    text = text.replace('\t', " ").replace('\n', " ").replace('\u', " ").replace('\\', "")
    # remove non ASCII (emoticon, chinese word, .etc)
    text = text.encode('ascii', 'replace').decode('ascii')
    text = text.lower()
    # remove mention, link, hashtag
    text = ' '.join(re.sub("(@#[\w_-]+)|(\w+:\/\/[\w_-]+)", " ", text).split())
    # remove incomplete URL
    return text.replace("http://", " ").replace("https://", " ")

df['Caption'] = df['Caption'].apply(remove_special)
```

# NLTK Tokenizing

```
# NLTK word rokenize
def word_tokenize_wrapper(text):
    return word_tokenize(text)

df['tokenized'] = df['Caption'].apply(word_tokenize_wrapper)

print('Tokenizing Result : \n')
print(df['tokenized'].head())
print('\n\n\n')

Tokenizing Result :

[karena, anies, bisa, menang, di, jakarta, saj...
[posisi, ketiga, terus, tapi, kayaknya, dijega...
[hasil, survey, litbang, kompas, sebelum, pilk...
[hehwowkwok, waktu, pilgub, jakarta, juga, pe...
[menang, survei, kok, bangga, liat, aja, dilap...
ne: tokenized, dtype: object
```

## 4. Filtering (removing stopwords)

```
1 from nltk.corpus import stopwords  
2  
3 #stopword indonesia  
4 list_stopwords = stopwords.words('indonesian')  
5  
6  
7 # menggabungkan additional stopword  
8 list_stopwords.extend(["yg", "dg", "rt", "dgn", "ny", "d", 'klo',  
9                 'kalo', 'amp', 'biar', 'bikin', 'bilang',  
10                'gak', 'ga', 'krn', 'nya', 'nih', 'sih',  
11                'si', 'tau', 'tdk', 'tuh', 'utk', 'ya',  
12                'jd', 'jgn', 'sdh', 'aja', 'n', 't',  
13                'nyg', 'hehe', 'pen', 'u', 'nan', 'loh', 'rt',  
14                '&amp', 'yah', 'awokawok', 'b aja', 'yang', 'jg'])  
15  
16
```

```
17 # read txt stopword using pandas
18 txt_stopword = pd.read_csv("stopwords.txt", names= ["stopwords"], header = None)
19
20 # convert stopword string to list & append additional stopword
21 list_stopwords.extend(txt_stopword["stopwords"][0].split(' '))
22
23
24 # memasukan list stopwords ke dalam corpus
25 list_stopwords = set(list_stopwords)
26
27
28 #remove stopword pada list token
29 def stopwords_removal(words):
30     return [word for word in words if word not in list_stopwords]
31
32 df['tokens_WSW'] = df['tokenized'].apply(stopwords_removal)
33
34
35 print(df['tokens_WSW'].head())
```

```
5 [anies, menang, jakarta, alasanya, ahok, kena,...
6 [posisi, ketiga, kayaknya, dijegal, kayak, why]
7 [hasil, survey, litbang, kompas, pilkada, dki,...
8 [hehwowkwowk, pilgub, jakarta, peringkat, mulu]
9 [menang, survei, bangga, liat, dilapangan, mas...
Name: tokens_WSW, dtype: object
```

## 5. Normalization

```
1 normalized_word = pd.read_excel("normalisasi.xlsx")
2
3 normalized_word_dict = {}
4
5 for index, row in normalized_word.iterrows():
6     if row[0] not in normalized_word_dict:
7         normalized_word_dict[row[0]] = row[1]
8
9 def normalized_term(document):
10    return [normalized_word_dict[term] if term in normalized_word_dict else term for term in document]
11
12 df['normalized'] = df['tokens_WSW'].apply(normalized_term)
13
14 df['normalized'].head(10)
```

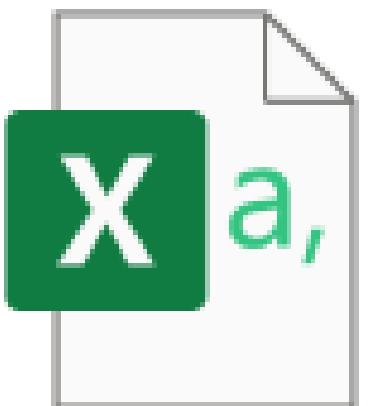
# 6. Stemming

```
1 from Sastrawi.Stemmer.StemmerFactory import StemmerFactory
2 import swifter
3
4
5 # create stemmer
6 factory = StemmerFactory()
7 stemmer = factory.create_stemmer()
8
9 # stemmed
10 def stemmed_wrapper(term):
11     return stemmer.stem(term)
12
13 term_dict = {}
14
15 for document in df['normalized']:
16     for term in document:
17         if term not in term_dict:
18             term_dict[term] = ''
```

```
20 print(len(term_dict))
21
22 for term in term_dict:
23     term_dict[term] = stemmed_wrapper(term)
24     print(term, ":" ,term_dict[term])
25
26 print(term_dict)
27 print("===")
28
29
30 # apply stemmed term to dataframe
31 def get_stemmed_term(document):
32     return [term_dict[term] for term in document]
33
34 df['tokens_stemmed'] = df['normalized'].swifter.apply(get_stemmed_term)
35 print(df['tokens_stemmed'])
```

```
↳ program : program
terealisakan : realisasi
datanya : data
posisi : posisi
ketiga : tiga
kayaknya : kayak
diiepal : iepal
```

# HASIL SETELAH DILAKUKAN PREPROCESSING



## process

# MEMBUAT MODEL BUAT TRAIN TEST

[https://github.com/Devildances/TwitterSentimentAnalysis\\_Final\\_Project/blob/master/data/indonesia\\_tweet\\_clean\\_tweets.csv](https://github.com/Devildances/TwitterSentimentAnalysis_Final_Project/blob/master/data/indonesia_tweet_clean_tweets.csv)



`clean_tweets`

# METODE BERT

## **Bidirectional Encoder Representations from Transformers**

**BERT adalah model bahasa yang dikembangkan oleh Google pada tahun 2018. Model ini menggunakan arsitektur Transformer yang merupakan jaringan saraf rekurensi (RNN) yang kuat dalam pemrosesan data berurutan.**

**BERT secara khusus dirancang untuk tugas pemahaman bahasa alami (natural language understanding, NLU) seperti pemrosesan teks, analisis sentimen, pemisahan kalimat, dan tugas-tugas lainnya. Keunikan utama BERT terletak pada kemampuannya untuk memahami konteks kata dalam kalimat dengan mempertimbangkan kata-kata yang ada sebelum dan sesudahnya, yang dikenal sebagai pemrosesan arah ganda (bidirectional processing).**

# SENTIMENT ANALYSIS MENGGUNAKAN METODE BERT

## a. import library

### ▼ Import necessary library

```
✓ [1] !pip install transformers
     !pip install Sastrawi

     import pandas as pd
     import numpy as np
     import torch
     from transformers import BertTokenizer, BertForSequenceClassification, AdamW, get_linear_schedule_with_warmup
     from sklearn.model_selection import train_test_split
     from torch.utils.data import TensorDataset, DataLoader, RandomSampler, SequentialSampler
     from sklearn.metrics import f1_score
     from Sastrawi.StopWordRemover.StopWordRemoverFactory import StopWordRemoverFactory
```

## b. Tokenize dengan BERT

### ▼ Tokenize data using BERT tokenizer

```
✓ [3] tokenizer = BertTokenizer.from_pretrained('bert-base-multilingual-cased', do_lower_case=True)

# Encode train and test data using BERT tokenizer
train_encodings = tokenizer(train_text.tolist(), truncation=True, padding=True)
test_encodings = tokenizer(test_text.tolist(), truncation=True, padding=True)

# Convert encoded inputs to PyTorch tensors
train_inputs = torch.tensor(train_encodings['input_ids'])
train_masks = torch.tensor(train_encodings['attention_mask'])
train_labels = torch.tensor(train_labels)
test_inputs = torch.tensor(test_encodings['input_ids'])
test_masks = torch.tensor(test_encodings['attention_mask'])
test_labels = torch.tensor(test_labels)

# Create data loaders for efficient batching
batch_size = 16
train_data = TensorDataset(train_inputs, train_masks, train_labels)
train_sampler = RandomSampler(train_data)
train_dataloader = DataLoader(train_data, sampler=train_sampler, batch_size=batch_size)
test_data = TensorDataset(test_inputs, test_masks, test_labels)
test_sampler = SequentialSampler(test_data)
test_dataloader = DataLoader(test_data, sampler=test_sampler, batch_size=batch_size)
```

▷ Downloading (...)solve/main/vocab.txt: 100%  996k/996k [00:00<00:00, 8.46MB/s]

▷ Downloading (...)okenizer\_config.json: 100%  29.0/29.0 [00:00<00:00, 1.72kB/s]

▷ Downloading (...)lve/main/config.json: 100%  625/625 [00:00<00:00, 15.1kB/s]

## c. Train the model

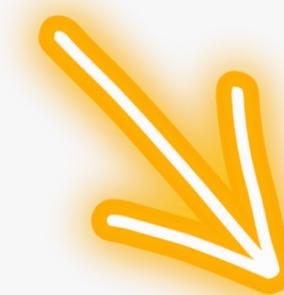
```
✓ [4] model = BertForSequenceClassification.from_pretrained('bert-base-multilingual-cased', num_labels=2, output_attentions=False, output_hidden_states=False)
      # Set device to GPU if available
      device = torch.device('cuda' if torch.cuda.is_available() else 'cpu')
      model.to(device)

      # Set up optimizer and scheduler
      epochs = 5
      optimizer = AdamW(model.parameters(), lr=2e-5, eps=1e-8)
      total_steps = len(train_dataloader) * epochs
      scheduler = get_linear_schedule_with_warmup(optimizer, num_warmup_steps=0, num_training_steps=total_steps)

      from tqdm import tqdm

      # Train the model
      for epoch in range(epochs):
          model.train()
          total_loss = 0
          progress_bar = tqdm(train_dataloader, desc=f'Epoch {epoch+1}', leave=False)
          for step, batch in enumerate(progress_bar):
              batch_inputs = batch[0].to(device)
              batch_masks = batch[1].to(device)
              batch_labels = batch[2].to(device)
              optimizer.zero_grad()
              outputs = model(batch_inputs, attention_mask=batch_masks, labels=batch_labels)
              loss = outputs[0]
              total_loss += loss.item()
              loss.backward()
              torch.nn.utils.clip_grad_norm_(model.parameters(), 1.0)
              optimizer.step()
              scheduler.step()
              avg_train_loss = total_loss / (step+1)
              progress_bar.set_postfix({'training_loss': avg_train_loss})
              print(f'Epoch: {epoch + 1}, Training Loss: {avg_train_loss}')

```



### ▼ Save the trained model

```
✓ [5] model_path = 'bert_sa_sentiment_model.pt'
      torch.save(model.state_dict(), model_path)
```

## d. Evaluasi model

- Evaluate the model

```
[6] # Load the saved model
17s   model = BertForSequenceClassification.from_pretrained('bert-base-multilingual-cased', num_labels=2, output_attentions=False, output_hidden_states=False)
        model.load_state_dict(torch.load(model_path))
        model.to(device)

# Evaluate the model on test data
model.eval()
total_preds = []
total_labels = []
with torch.no_grad():
    for batch in test_dataloader:
        batch_inputs = batch[0].to(device)
        batch_masks = batch[1].to(device)
        batch_labels = batch[2].to('cpu')
        outputs = model(batch_inputs, attention_mask=batch_masks)
        logits = outputs[0]
        preds = torch.argmax(logits, axis=1)
        total_preds.extend(preds)
        total_labels.extend(batch_labels)

    total_labels = torch.tensor(total_labels)
    total_preds = torch.tensor(total_preds)

# Calculate F1 score
f1 = f1_score(total_labels.detach().cpu().numpy(), total_preds.detach().cpu().numpy(), average='weighted')
print(f'Weighted F1 score: {f1:.4f}')
```



Some weights of BertForSequenceClassification  
You should probably TRAIN this model  
Weighted F1 score: 0.8491

## e. Memasukkan text nya untuk memprediksi hasil sentiment

```
[7] def predict_sentiment(text):
    # Remove stop words from input text

    text = stop_words.remove(text)
    # Tokenize input text using BERT tokenizer
    encoding = tokenizer.encode_plus(text, max_length=64, truncation=True, padding='max_length', add_special_tokens=True, return_attention_mask=True, return_tensors='pt')
    inputs = encoding['input_ids'].to(device)
    masks = encoding['attention_mask'].to(device)
    # Pass input through the trained model and return predicted sentiment
    with torch.no_grad():
        outputs = model(inputs, attention_mask=masks)
        logits = outputs[0]
        probs = torch.softmax(logits, dim=1)
        _, pred_label = torch.max(probs, dim=1)
    return 'positive' if pred_label == 1 else 'negative'
```

# Result:

```
[8] text = "Bodoh kamu"
sentiment = predict_sentiment(text)
print(sentiment)

negative

[9] text = "Ditengah masyarakat yg gak begitu suka baca dan lebih tertarik dengan hal2 visual, Ganjar dan Prabowo lebih ahlinya, sedikit membicarakan data, jarang mengkritiki, tapi lebih sering tampil di medi
sentiment = predict_sentiment(text)
print(sentiment)

positive

[10] text = "Insya Allah Pak Anies Presiden 2024 ini akun bayaran dan survei bayaran"
sentiment = predict_sentiment(text)
print(sentiment)

positive

[11] text = "Sampai skrg saya bingung apa ide dan gagasan perubahan dari Anies? Sampai skrg ini, Anies masih berlutut menyinggung program program Jokowi belum ada suatu hal yang besar ide perubahan yang selalu
sentiment = predict_sentiment(text)
print(sentiment)

positive
```

## f. Memprediksi proces.csv (hasil data preprocessing)

			Unnamed: 0	Caption	tokenized	freqdist	tokens_NSW	normalized	tokens_stemmed	merged
0	5	karena anies bisa menang di jakarta saja karen...	['karena', 'anies', 'bisa', 'menang', 'di', 'j...']	<FreqDist with 45 samples and 54 outcomes>	['anies', 'menang', 'jakarta', 'alasanya', 'ah...	anies menang jakarta alasanya ahok kena beliau...				
1	6	posisi ketiga terus tapi kayaknya dijegal teru...	['posisi', 'ketiga', 'terus', 'tapi', 'kayakny...']	<FreqDist with 8 samples and 9 outcomes>	['posisi', 'ketiga', 'kayaknya', 'dijegal', 'k...']	['posisi', 'ketiga', 'kayaknya', 'dijegal', 'k...']	['posisi', 'tiga', 'kayak', 'jegal', 'kayak', ...]	['posisi', 'tiga', 'kayak', 'jegal', 'kayak', ...]	posisi ketiga kayaknya dijegal kayak why	
2	7	hasil survei litbang kompas sebelum pilkada dk...	['hasil', 'survey', 'litbang', 'kompas', 'sebe...']	<FreqDist with 17 samples and 18 outcomes>	['hasil', 'survey', 'litbang', 'kompas', 'pilk...']	hasil survei litbang kompas pilkada dkj ahysil...				
3	8	hehwowkwowk waktu pilgub jakarta juga peringka...	['hehwowkwowk', 'waktu', 'pilgub', 'jakarta', ...]	<FreqDist with 9 samples and 10 outcomes>	['hehwowkwowk', 'pilgub', 'jakarta', 'peringka...']	hehwowkwowk pilgub jakarta peringkat mulu				
4	9	menang survei kok bangga liat aja dilapangan m...	['menang', 'survei', 'kok', 'bangga', 'liat', ...]	<FreqDist with 11 samples and 11 outcomes>	['menang', 'survei', 'bangga', 'liat', 'dilapa...']	['menang', 'survei', 'bangga', 'lihat', 'dilap...']	['menang', 'survei', 'bangga', 'lihat', 'lapan...']	['menang', 'survei', 'bangga', 'lihat', 'lapan...']	menang survei bangga lihat dilapangan masyarakat...	

## g. Menghapus data yang tidak diperlukan

```
[14] predict.drop(['Unnamed: 0', 'tokenized', 'freqdist', 'tokens_NSW', 'tokens_stemmed', 'normalized'], inplace=True, axis=1)
```

## **h. Memasukkan kolom hasil sentimen analisis ke dalam dataframe**

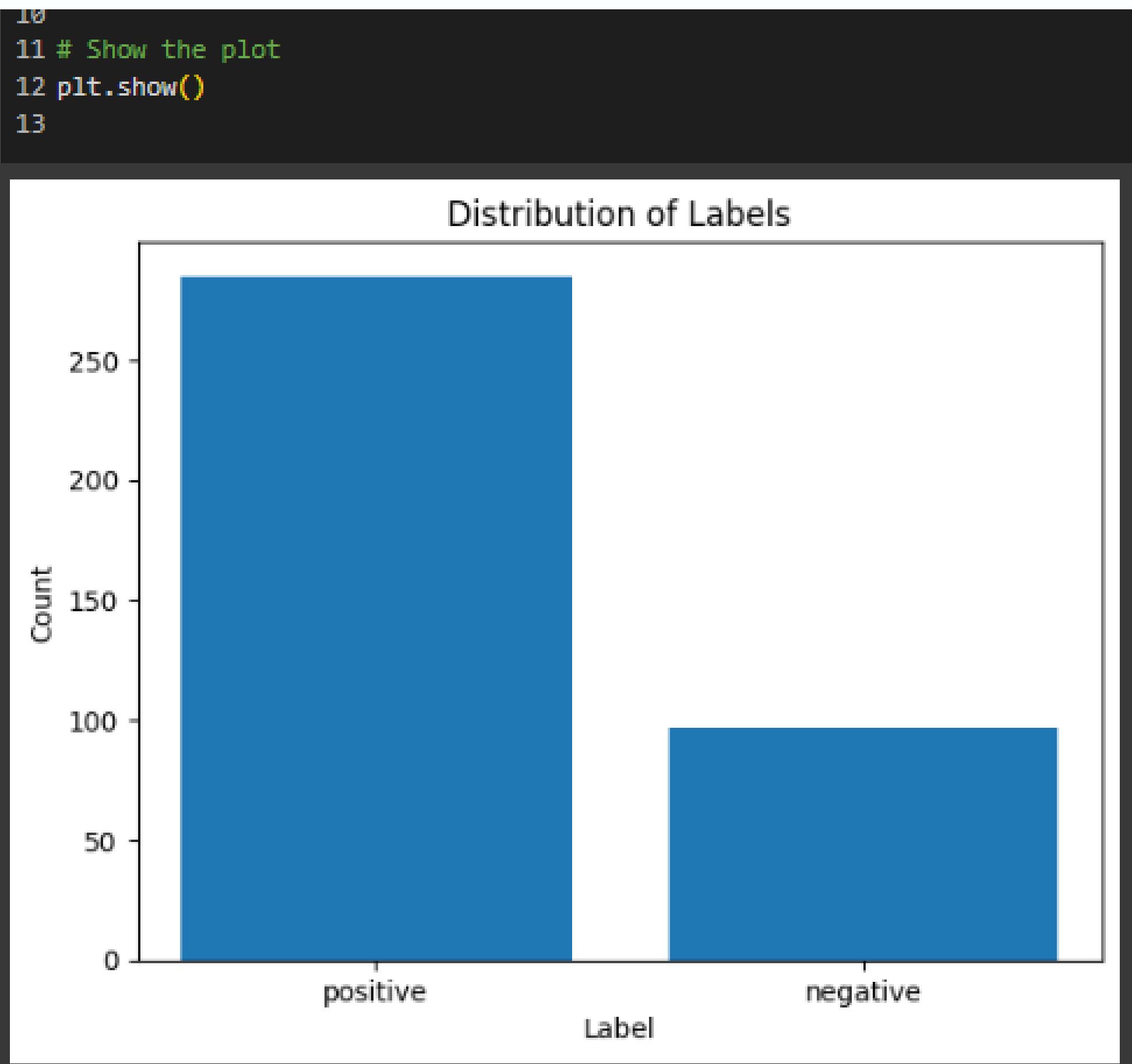
```
[17] # Add the predicted_sentiments list as a new column to the DataFrame  
predict['predicted_sentiments'] = predicted_sentiments  
  
# Write the updated DataFrame to the CSV file  
predict.to_csv('predict2.csv', index=False)
```

# print hasil:

```
[✓] [16] # Assuming you have already imported the necessary libraries and have the DataFrame `data` with the desired column  
2s  
# Create an empty list to store the predicted sentiments  
predicted_sentiments = []  
  
# Iterate through the column values using a loop  
for text in predict['merged']:  
    # Check if the value is a string  
    if isinstance(text, str):  
        # Apply the predict_sentiment() function to each text value  
        sentiment = predict_sentiment(text)  
    else:  
        # Handle non-string values (e.g., NaN or float)  
        sentiment = 'Unknown'  
    # Append the predicted sentiment to the list  
    predicted_sentiments.append(sentiment)  
  
# Print or manipulate the list of predicted sentiments  
for sentiment in predicted_sentiments:  
    print(sentiment)
```

positive  
negative  
negative  
positive  
negative  
positive  
Unknown  
positive  
positive  
negative  
positive  
positive  
positive  
positive  
negative  
positive

## i. Histogram hasil perbandingan analisis sentimen





# KESIMPULAN

# HASIL :



**predict2**

# KESIMPULAN

Hasil analisis sentimen untuk topik "Kenapa Elektabilitas Anies Tak Kunjung Naik?" membawaikan dua bagian sentimen, positif dan negatif. Namun, dari hasil analisis, terlihat bahwa lebih banyak terdapatnya sentimen positif dari komen netizen di instagram.

## Sentimen positif:

- bismillah anis menang pilpres tak peduli apa kata bazer
- saya pilih anis

## Sentimen negatif:

- anies jarang mandi kyny



# TERIMA KASIH