

Evaluasi CPMK 1 dan

Soal 1

1.CPMK1- Sub CPMK1.1(bobot : 10)

Lakukan crawling data teks dari media sosial/ web dan simpan hasilnya dalam bentuk excel/csv.

Setiap mahasiswa harus melakukan crawling dengan kata kunci tertentu.

Kata kunci tidak boleh sama dengan mahasiswa lainnya.

```
1 # Import required Python package
2 !pip install pandas
3
4 # Install Node.js (because tweet-harvest built using Node.js)
5 !sudo apt-get update
6 !sudo apt-get install -y ca-certificates curl gnupg
7 !sudo mkdir -p /etc/apt/keyrings
8 !curl -fsSL https://deb.nodesource.com/gpgkey/nodesource-repo.gpg.key | sudo gpg --
9
10 !NODE_MAJOR=20 && echo "deb [signed-by=/etc/apt/keyrings/nodesource.gpg] https://de
11
12 !sudo apt-get update
13 !sudo apt-get install nodejs -y
14
15 !node -v
```

```
Requirement already satisfied: pandas in /usr/local/lib/python3.10/dist-packages (
Requirement already satisfied: python-dateutil>=2.8.1 in /usr/local/lib/python3.10
Requirement already satisfied: pytz>=2020.1 in /usr/local/lib/python3.10/dist-pack
Requirement already satisfied: numpy>=1.21.0 in /usr/local/lib/python3.10/dist-pac
Requirement already satisfied: six>=1.5 in /usr/local/lib/python3.10/dist-packages
Hit:1 https://cloud.r-project.org/bin/linux/ubuntu jammy-cran40/ InRelease
Hit:2 https://developer.download.nvidia.com/compute/cuda/repos/ubuntu2204/x86_64
Hit:3 https://deb.nodesource.com/node_20.x nodistro InRelease
Get:4 http://security.ubuntu.com/ubuntu jammy-security InRelease [110 kB]
Hit:5 http://archive.ubuntu.com/ubuntu jammy InRelease
Get:6 http://archive.ubuntu.com/ubuntu jammy-updates InRelease [119 kB]
Hit:7 https://ppa.launchpadcontent.net/c2d4u.team/c2d4u4.0+/ubuntu jammy InRelease
Get:8 http://archive.ubuntu.com/ubuntu jammy-backports InRelease [109 kB]
Hit:9 https://ppa.launchpadcontent.net/deadsnakes/ppa/ubuntu jammy InRelease
Hit:10 https://ppa.launchpadcontent.net/graphics-drivers/ppa/ubuntu jammy InRelease
Hit:11 https://ppa.launchpadcontent.net/ubuntugis/ppa/ubuntu jammy InRelease
Fetched 338 kB in 3s (132 kB/s)
Reading package lists... Done
Reading package lists... Done
Building dependency tree... Done
Reading state information... Done
ca-certificates is already the newest version (20230311ubuntu0.22.04.1).
curl is already the newest version (7.81.0-1ubuntu1.14).
gnupg is already the newest version (2.2.27-3ubuntu2.1).
0 upgraded, 0 newly installed, 0 to remove and 20 not upgraded.
gpg: cannot open '/dev/tty': No such device or address
curl: (23) Failed writing body
deb [signed-by=/etc/apt/keyrings/nodesource.gpg] https://deb.nodesource.com/node_2
Hit:1 https://cloud.r-project.org/bin/linux/ubuntu jammy-cran40/ InRelease
Hit:2 https://deb.nodesource.com/node_20.x nodistro InRelease
Hit:3 https://developer.download.nvidia.com/compute/cuda/repos/ubuntu2204/x86_64
Hit:4 http://security.ubuntu.com/ubuntu jammy-security InRelease
Hit:5 http://archive.ubuntu.com/ubuntu jammy InRelease
Hit:6 http://archive.ubuntu.com/ubuntu jammy-updates InRelease
Hit:7 http://archive.ubuntu.com/ubuntu jammy-backports InRelease
Hit:8 https://ppa.launchpadcontent.net/c2d4u.team/c2d4u4.0+/ubuntu jammy InRelease
Hit:9 https://ppa.launchpadcontent.net/deadsnakes/ppa/ubuntu jammy InRelease
Hit:10 https://ppa.launchpadcontent.net/graphics-drivers/ppa/ubuntu jammy InRelease
Hit:11 https://ppa.launchpadcontent.net/ubuntugis/ppa/ubuntu jammy InRelease
Reading package lists... Done
Reading package lists... Done
Building dependency tree... Done
Reading state information... Done
nodejs is already the newest version (20.8.1-1nodesource1).
0 upgraded, 0 newly installed, 0 to remove and 20 not upgraded.
v20.8.1
```

```
1 # Crawl Data
2
```



```
..
node_modules
sample_data
tweets-data
cleaned_result.csv
cleaned_twitter_data.csv
package-lock.json
package.json
parsed_result.csv
result.csv
slang_tokens.csv
slang_twitter_data.csv
slangword.csv
stemmed_twitter_data.csv
stopword_twitter_data.csv
tokenized_twitter_data.csv
```

```

3 filename = 'result.csv'
4 search_keyword = 'ganjar pranowo'
5 limit = 500
6
7 Inpx --yes tweet-harvest@latest -o "{filename}" -s "{search_keyword}" -l {limit}

```

```

Your tweets saved to: /content/tweets-data/result.csv
Total tweets saved: 340

```

```

Got some tweets, saving to file...
Your tweets saved to: /content/tweets-data/result.csv
Total tweets saved: 351

```

```

Got some tweets, saving to file...
Your tweets saved to: /content/tweets-data/result.csv
Total tweets saved: 367

```

```

Got some tweets, saving to file...
Your tweets saved to: /content/tweets-data/result.csv
Total tweets saved: 382

```

```

Got some tweets, saving to file...
Your tweets saved to: /content/tweets-data/result.csv
Total tweets saved: 397

```

```

Got some tweets, saving to file...
Your tweets saved to: /content/tweets-data/result.csv
Total tweets saved: 408

```

```

Got some tweets, saving to file...
Your tweets saved to: /content/tweets-data/result.csv
Total tweets saved: 420

```

```

Got some tweets, saving to file...
Your tweets saved to: /content/tweets-data/result.csv
Total tweets saved: 431

```

```

--Taking a break, waiting for 10 seconds...

```

```

Got some tweets, saving to file...
Your tweets saved to: /content/tweets-data/result.csv
Total tweets saved: 440

```

```

Got some tweets, saving to file...
Your tweets saved to: /content/tweets-data/result.csv
Total tweets saved: 454

```

```

Got some tweets, saving to file...
Your tweets saved to: /content/tweets-data/result.csv
Total tweets saved: 465

```

```

Got some tweets, saving to file...
Your tweets saved to: /content/tweets-data/result.csv
Total tweets saved: 480

```

```

Got some tweets, saving to file...
Your tweets saved to: /content/tweets-data/result.csv
Total tweets saved: 497

```

```

Got some tweets, saving to file...
Your tweets saved to: /content/tweets-data/result.csv
Total tweets saved: 512
Already got 512 tweets, done scrolling...

```

```

1 import pandas as pd
2
3 # Specify the path to your CSV file
4 file_path = f"tweets-data/{filename}"
5
6 # Read the CSV file into a pandas DataFrame
7 df = pd.read_csv(file_path, delimiter=";")
8
9 # Display the DataFrame
10 display(df)

```

	created_at	id_str	full_text	quote_count	reply_count	re
0	Fri Oct 20 11:03:30 +0000 2023	1715322846833205374	@syahirularif Yang pasti @jokowi @gibran_tweet...	0	0	
1	Fri Oct 20 11:03:25 +0000 2023	1715322827472335187	@Nikmatul_Sg @ganjarpranowo @mohmahfudmd Visi ...	0	0	
2	Fri Oct 20 11:03:23 +0000 2023	1715322818022592572	@Anase1985 @Muhayya__ @ganjarpranowo @mohmahfu...	0	0	
3	Fri Oct 20 11:03:22 +0000 2023	1715322813689843824	@Melihat_Indo @ganjarpranowo @mohmahfudmd maki...	0	0	
4	Fri Oct 20 11:03:22 +0000 2023	1715322813408837699	@Prihati_utami Dimana pun berada pilihan tetap...	0	0	
...
507	Fri Oct 20 10:34:07 +0000 2023	1715315452325642461	@ganjarpranowo Pak Ganjar memang pantas di duk...	0	0	
508	Fri Oct 20 10:34:07 +0000 2023	1715315450840551783	@Melihat_Indo @ganjarpranowo @mohmahfudmd Mere...	0	0	
509	Fri Oct 20 10:34:05 +0000 2023	1715315443194605775	@imadya @ganjarpranowo @basuki_btp Ditunggu re...	0	0	

1 # Cek jumlah data yang didapatkan

2

3 num_tweets = len(df)

4 print(f"Jumlah tweet dalam dataframe adalah {num_tweets}.")

Jumlah tweet dalam dataframe adalah 512.

+0000 2023

sudah terbukti

Soal 2

Lakukan pembersihan meliputi penghapusan punctuation, angka dan karakter yang tidak penting menggunakan menggunakan regex.
Simpan hasilnya menjadi file csv/excel.

```
1 import regex as re
2 import pandas as pd

1 def clean_text(text):
2     # Menghapus tanda baca dan karakter non-alfanumerik
3     cleaned_text = re.sub(r'^a-zA-Z\s', '', text)
4     return cleaned_text

1 input_file = "result.csv"
2 output_file = "cleaned_result.csv"

1 file_path = f"tweets-data/{filename}"
2 df = pd.read_csv(file_path, delimiter=";")

1 df['full_text'] = df['full_text'].apply(clean_text)

1 df.to_csv(output_file, index=False, encoding='utf-8')
2
3 print('Data sudah dibersihkan dengan nama file cleaned_result.csv')

Data sudah dibersihkan dengan nama file cleaned_result.csv
```

1 pd.read_csv('/content/cleaned_result.csv')

	created_at	id_str	full_text	quote_count	reply_count	retw
0	Fri Oct 20 11:03:30 +0000 2023	1715322846833205374	syahirularif Yang pasti jokowi gibrantweet duk...	0	0	
1	Fri Oct 20 11:03:25 +0000 2023	1715322827472335187	NikmatulSg ganjarpranowo mohmahfudmd Visi misi...	0	0	
2	Fri Oct 20 11:03:23 +0000 2023	1715322818022592572	Anase Muhayya ganjarpranowo mohmahfudmd chchot...	0	0	
3	Fri Oct 20 11:03:22 +0000 2023	1715322813689843824	MelihatIndo ganjarpranowo mohmahfudmd makin ma...	0	0	
4	Fri Oct 20 11:03:22 +0000 2023	1715322813408837699	Prihatitutami Dimana pun berada pilihan tetap g...	0	0	
...
507	Fri Oct 20 10:34:07 +0000 2023	1715315452325642461	ganjarpranowo Pak Ganjar memang pantas di duku...	0	0	
508	Fri Oct 20 10:34:07 +0000 2023	1715315450840551783	MelihatIndo ganjarpranowo mohmahfudmd Mereka m...	0	0	
509	Fri Oct 20 10:34:05 +0000 2023	1715315443194605775	imadya ganjarpranowo basukibtp Ditunggu rekama...	0	0	
510	Fri Oct 20 10:33:58 +0000 2023	1715315416174842157	MelihatIndo ganjarpranowo mohmahfudmd Kita yak...	0	0	
511	Fri Oct 20 10:33:57 +0000 2023	1715315411368210936	ganjarpranowo Pak ganjar yang sudah terbukti k...	0	0	

512 rows × 12 columns

▼ Soal 3a

a. Lakukan parsing dan simpan hasilnya dalam bentuk excel/csv (bobot:10)

```
1 # Buka file CSV asal dan buat file CSV baru untuk hasil parsing
2 with open('cleaned_result.csv', 'r', encoding='utf-8') as csvfile:
3     reader = csv.DictReader(csvfile)
4
5     with open('parsed_result.csv', 'w', newline='', encoding='utf-8') as outputfile
6         fieldnames = ['Word'] # Kolom dalam file CSV baru
7         writer = csv.DictWriter(outputfile, fieldnames=fieldnames)
8         writer.writeheader()
9
10        for row in reader:
11            # Akses data dari kolom "full_text"
12            full_text = row['full_text']
13
14            # Memisahkan teks menjadi kata-kata
15            words = full_text.split()
16
17            # Menulis kata-kata ke dalam file CSV baru
```

```

18         for word in words:
19             writer.writerow({'Word': word})
20
21 print('Hasil parsing telah disimpan dalam "parsed_result.csv"')

```

Hasil parsing telah disimpan dalam "parsed_result.csv"

▼ Soal 3b

b.Carilah kata-kata slangword yang ada dalam dataset Anda, dengan cara mencocokkan dengan kamus KBBI (terlampir).

Simpan hasilnya dalam bentuk csv/excel.

Tampilkan 100 kata slang yang Anda dapatkan dan tampilkan dalam bentuk Gunakan (bobot:25)

```

1 # Baca file CSV dengan dataset Twitter (cleaned_result.csv)
2 twitter_data = pd.read_csv('cleaned_result.csv', encoding='utf-8')
3
4 # Baca file CSV dengan daftar kata-kata slang (slangword.csv)
5 slang_data = pd.read_csv('slangword.csv', encoding='utf-8', names=['Slang Word'])
6
7 # Konversi daftar kata-kata slang ke dalam bentuk set untuk pencocokan lebih efisien
8 slang_set = set(slang_data['Slang Word'])
9
10 # Buat kolom baru dalam dataset Twitter untuk kata-kata yang bukan slang
11 twitter_data['Cleaned Text'] = twitter_data['full_text'].apply(
12     lambda text: ' '.join([word for word in text.split() if word not in slang_set])
13 )
14
15 # Simpan hasilnya dalam file CSV baru
16 twitter_data.to_csv('slang_twitter_data.csv', index=False, encoding='utf-8')
17
18 # Tampilkan 100 kata slang
19 print('100 Kata Slang:')
20 print(slang_data['Slang Word'][:100])
21
22 # Tampilkan hasilnya dalam bentuk data frame
23 print('Hasil Data Frame Setelah Menghapus Slang:')
24 print(twitter_data.head())

```

```

100 Kata Slang:
0      Slangwords
1      Adam
2      Aga
3      Agustus
4      Ahad
...
95      Jogi
96      Johar
97      Juja
98      Juli
99      Jumadilakhir
Name: Slang Word, Length: 100, dtype: object
Hasil Data Frame Setelah Menghapus Slang:

```

	created_at	id_str \
0	Fri Oct 20 11:03:30 +0000 2023	1715322846833205374
1	Fri Oct 20 11:03:25 +0000 2023	1715322827472335187
2	Fri Oct 20 11:03:23 +0000 2023	1715322818022592572
3	Fri Oct 20 11:03:22 +0000 2023	1715322813689843824
4	Fri Oct 20 11:03:22 +0000 2023	1715322813408837699

```


```

	full_text	quote_count \
0	syahirularif Yang pasti jokowi gibrantweet duk...	0
1	NikmatulSg ganjarpranowo mohmahfudmd Visi misi...	0
2	Anase Muhayya ganjarpranowo mohmahfudmd chchot...	0
3	MelihatIndo ganjarpranowo mohmahfudmd makin ma...	0
4	Prihatiutami Dimana pun berada pilihan tetap g...	0

```


```

	reply_count	retweet_count	favorite_count	lang	user_id_str \
0	0	0	0	in	3244905644
1	0	0	0	in	1707053631005212672
2	0	0	0	in	1260061504767836162
3	0	0	0	in	1659432246464622594
4	0	0	0	in	1641586661233692678

```


```

	conversation_id_str	username \
0	1715320765036495140	Mi73Hel
1	1715286319017103792	Anita33_
2	1715305293226348847	Dharma_tc
3	1715293032529203504	bougenwill4
4	1715026015330914478	fadlian_syah29

```

                                tweet_url \
0  https://twitter.com/Mi73Hel/status/17153228468...
1  https://twitter.com/Anita33/status/1715322827...
2  https://twitter.com/Dharma_tc/status/171532281...
3  https://twitter.com/bougenwill4/status/1715322...
4  https://twitter.com/fadlian_syah29/status/1715...

                                Cleaned Text
0  syahirularif Yang jokowi gibrantweet ganjarpra...
1  NikmatulSg ganjarpranowo mohmahfudmd Visi dila...
2  Anase Muhayya ganjarpranowo mohmahfudmd chchot...
3  MelihatIndo ganjarpranowo mohmahfudmd VisiMisi...
4  Prihatiutami Dimana pranowo gaspol menangkan

1  ## Membaca file "slangword.csv" untuk mengumpulkan daftar slangword
2  # slangwords = set()
3  # with open('slangword.csv', 'r', encoding='utf-8') as slangfile:
4  #     reader = csv.reader(slangfile)
5  #     for row in reader:
6  #         slangwords.add(row[0].strip().lower()) # Memastikan semua slangword dal
7
8  ## Membaca file "cleaned_result.csv" untuk mencari kata-kata slang
9  # slang_found = set()
10 # with open('cleaned_result.csv', 'r', encoding='utf-8') as csvfile:
11 #     reader = csv.DictReader(csvfile)
12
13 #     for row in reader:
14 #         full_text = row['full_text']
15
16 #         # Memisahkan teks menjadi kata-kata
17 #         words = full_text.split()
18
19 #         # Mencari kata-kata slang
20 #         for word in words:
21 #             if word.lower() in slangwords:
22 #                 slang_found.add(word)
23
24 # ## Menyimpan hasil pencarian dalam file CSV
25 # with open('slang_found.csv', 'w', newline='', encoding='utf-8') as outputfile:
26 #     fieldnames = ['Slang Word']
27 #     writer = csv.DictWriter(outputfile, fieldnames=fieldnames)
28 #     writer.writeheader()
29 #     for word in slang_found:
30 #         writer.writerow({'Slang Word': word})
31
32 # ## Menampilkan 100 kata slang pertama
33 # slang_found_list = list(slang_found)
34 # print('kata slang yang ditemukan:')
35 # for i in range(400):
36 #     print(slang_found_list[i])

```

▼ Soal 3c

c.Lakukan tokenizing berdasarkan hasil 3b, simpan hasilnya dalam bentuk csv/excel dan tampilkan 100 token pertama. (bobot:10)

```

1 data = []
2 for slang in slang_found:
3     tokens = re.findall(r'\b\w+\b', slang) # Tokenisasi dengan menghilangkan karak
4     data.append([slang, ' '.join(tokens)])
5
6 df = pd.DataFrame(data, columns=['Slang Word', 'Tokens'])
7
8 # Menyimpan hasil dalam file CSV
9 df.to_csv('slang_tokens.csv', index=False, encoding='utf-8')
10
11 # Menampilkan 100 token pertama dalam bentuk data frame
12 print('100 token pertama:')
13 print(df.head(100))

```

```

100 token pertama:
   Slang Word      Tokens
0         akal         akal
1         Tuhan         Tuhan
2         dulu         dulu
3         peran         peran
4         Ketua         Ketua
..         ...         ...

```

```

95     bangsa     bangsa
96     pernyataan pernyataan
97     diskusi     diskusi
98     Kabupaten  Kabupaten
99     Kebangsaan  Kebangsaan

```

```
[100 rows x 2 columns]
```

```

1 import nltk
2 nltk.download('punkt')
3 from nltk.tokenize import word_tokenize

[nltk_data] Downloading package punkt to /root/nltk_data...
[nltk_data] Unzipping tokenizers/punkt.zip.

1 from nltk.tokenize import word_tokenize
2
3 # Baca file CSV dengan data Twitter yang telah dibersihkan (cleaned_twitter_data.csv)
4 twitter_data = pd.read_csv('cleaned_twitter_data.csv', encoding='utf-8')
5
6 # Tokenisasi teks dalam kolom 'Cleaned Text'
7 twitter_data['Tokens'] = twitter_data['Cleaned Text'].apply(lambda text: word_tokenize(text))
8
9 # Simpan hasil tokenisasi dalam file CSV
10 twitter_data.to_csv('tokenized_twitter_data.csv', index=False, encoding='utf-8')
11
12 # Tampilkan 100 token pertama dalam bentuk data frame
13 print('100 Token Pertama:')
14 print(twitter_data['Tokens'].apply(lambda tokens: tokens[:100]))
15
16 # Tampilkan hasilnya dalam bentuk data frame
17 print('Hasil Data Frame Setelah Tokenisasi:')
18 print(twitter_data.head())

```

```

3      [MelihatIndo, ganjarpranowo, mohmahfudmd, Visi...
4      [Prihatiutami, Dimana, pranowo, gaspol, menang...
      ...
507      [ganjarpranowo, Pak, Ganjar, RI]
508      [MelihatIndo, ganjarpranowo, mohmahfudmd, Mere...
509      [imadya, ganjarpranowo, basukibtp, Ditunggu, p...
510      [MelihatIndo, ganjarpranowo, mohmahfudmd, Kita...
511      [ganjarpranowo, Pak, kinerjanya]
Name: Tokens, Length: 512, dtype: object
Hasil Data Frame Setelah Tokenisasi:
      created_at      id_str \
0  Fri Oct 20 11:03:30 +0000 2023 1715322846833205374
1  Fri Oct 20 11:03:25 +0000 2023 1715322827472335187
2  Fri Oct 20 11:03:23 +0000 2023 1715322818022592572
3  Fri Oct 20 11:03:22 +0000 2023 1715322813689843824
4  Fri Oct 20 11:03:22 +0000 2023 1715322813408837699

      full_text  quote_count \
0  syahirularif Yang pasti jokowi gibrantweet duk...      0
1  NikmatulSg ganjarpranowo mohmahfudmd Visi misi...      0
2  Anase Muhayya ganjarpranowo mohmahfudmd chhot...      0
3  MelihatIndo ganjarpranowo mohmahfudmd makin ma...      0
4  Prihatiutami Dimana pun berada pilihan tetap g...      0

      reply_count  retweet_count  favorite_count  lang  user_id_str \
0              0              0              0  in      3244905644
1              0              0              0  in 1707053631005212672
2              0              0              0  in 1260061504767836162
3              0              0              0  in 1659432246464622594
4              0              0              0  in 1641586661233692678

      conversation_id_str  username \
0  1715320765036495140      Mi73Hel
1  1715286319017103792      Anita33_
2  1715305293226348847      Dharma_tc
3  1715293032529203504      bougenvill4
4  1715026015330914478  fadlian_syah29

      tweet_url \
0  https://twitter.com/Mi73Hel/status/17153228468...
1  https://twitter.com/Anita33_/status/1715322827...
2  https://twitter.com/Dharma_tc/status/171532281...
3  https://twitter.com/bougenvill4/status/1715322...
4  https://twitter.com/fadlian_syah29/status/1715...

      Cleaned Text \
0  syahirularif Yang jokowi gibrantweet ganjarpra...
1  NikmatulSg ganjarpranowo mohmahfudmd Visi dila...

```

```

0 [syahirularif, Yang, jokowi, gibrantweet, ganj...
1 [NikmatulSg, ganjarpranowo, mohmahfudmd, Visi,...
2 [Anase, Muhayya, ganjarpranowo, mohmahfudmd, c...
3 [MelihatIndo, ganjarpranowo, mohmahfudmd, Visi...
4 [Prihatiutami, Dimana, pranowo, gaspol, menang...

```

▼ Soal 3d

d.Lakukan stopwords removing berdasarkan hasil 3c, simpan hasilnya dalam bentuk csv/excel. (bobot:10)!

```

1 from nltk.corpus import stopwords
2 from nltk.tokenize import word_tokenize
3
4 # Download dataset stopwords untuk bahasa Indonesia
5 nltk.download('stopwords')
6 nltk.download('punkt')

[nltk_data] Downloading package stopwords to /root/nltk_data...
[nltk_data] Package stopwords is already up-to-date!
[nltk_data] Downloading package punkt to /root/nltk_data...
[nltk_data] Package punkt is already up-to-date!
True

1 # Baca file CSV dengan dataset Twitter (tokenized_twitter_data.csv)
2 twitter_data = pd.read_csv('tokenized_twitter_data.csv', encoding='utf-8')
3
4 # Fungsi untuk menghapus stopwords
5 def remove_stopwords(text):
6     stop_words = set(stopwords.words('indonesian'))
7     words = word_tokenize(text)
8     return [word for word in words if word.lower() not in stop_words]
9
10 # Melakukan penghapusan stopwords pada kolom "Tokens"
11 twitter_data['Tokens'] = twitter_data['Tokens'].apply(remove_stopwords)
12
13 # Simpan hasilnya dalam file CSV baru
14 twitter_data.to_csv('stopword_twitter_data.csv', index=False, encoding='utf-8')
15
16 print('Stopwords telah dihapus dan hasilnya disimpan dalam "stopword_twitter_data.csv".

Stopwords telah dihapus dan hasilnya disimpan dalam "stopword_twitter_data.csv".

```

▼ Soal 3e

e.Lakukan stemming berdasarkan hasil 3d dan tampilkan 100 stem pertama. (bobot:10)

```

1 !pip install Sastrawi
2 from Sastrawi.Stemmer.StemmerFactory import StemmerFactory

Requirement already satisfied: Sastrawi in /usr/local/lib/python3.10/dist-packages

1 # Baca file CSV dengan dataset Twitter yang sudah dibersihkan dari stopwords
2 twitter_data = pd.read_csv('stopword_twitter_data.csv', encoding='utf-8')
3
4 # Inisialisasi Stemmer Bahasa Indonesia
5 factory = StemmerFactory()
6 stemmer = factory.create_stemmer()
7
8 # Fungsi untuk melakukan stemming
9 def stem_text(text):
10     return stemmer.stem(text)
11
12 # Melakukan stemming pada kolom "Cleaned Text"
13 twitter_data['Stemmed Text'] = twitter_data['Cleaned Text'].apply(stem_text)
14
15 # Simpan hasil stemming dalam file CSV baru
16 twitter_data.to_csv('stemmed_twitter_data.csv', index=False, encoding='utf-8')
17
18 # Tampilkan 100 kata hasil stemming pertama
19 print('100 Kata Hasil Stemming Pertama:')
20 stemmed_tokens = twitter_data['Stemmed Text'].apply(lambda text: ' '.join(text.split()))
21 print(stemmed_tokens)

```



```
100 Kata Hasil Stemming Pertama:
0      syahirularif yang jokowi gibrantweet ganjarpra...
1      nikmatulsg ganjarpranowo mohmahfudmd visi laks...
2      anase muhayya ganjarpranowo mohmahfudmd chhot...
3      melihatindo ganjarpranowo mohmahfudmd visimisi...
4      prihatiutami mana pranowo gaspol menang
      ...
507      ganjarpranowo pak ganjar ri
508      melihatindo ganjarpranowo mohmahfudmd mereka m...
509      imadya ganjarpranowo basukibtp tunggu podcast ...
510      melihatindo ganjarpranowo mohmahfudmd kita vis...
511      ganjarpranowo pak kerja
Name: Stemmed Text, Length: 512, dtype: object
```

Disk  79.78 GB available