```
1 from google.colab import drive
2 drive.mount('/content/drive')
   Mounted at /content/drive
1 !pip install contractions
   Looking in indexes: <a href="https://pypi.org/simple">https://us-python.pkg.dev/colab-</a>
   Collecting contractions
     Downloading contractions-0.1.72-py2.py3-none-any.whl (8.3 kB)
   Collecting textsearch>=0.0.21
     Downloading textsearch-0.0.24-py2.py3-none-any.whl (7.6 kB)
   Collecting pyahocorasick
     Downloading pyahocorasick-1.4.4-cp37-cp37m-manylinux 2 17 x86 64.manylinux2
                                         | 106 kB 5.2 MB/s
   Collecting anyascii
     Downloading anyascii-0.3.1-py3-none-any.whl (287 kB)
                                          | 287 kB 37.1 MB/s
   Installing collected packages: pyahocorasick, anyascii, textsearch, contracti
   Successfully installed anyascii-0.3.1 contractions-0.1.72 pyahocorasick-1.4.4
  4
1 import pandas as pd
2 import contractions
3 import re
4 import string
5 # Importing wordcloud for plotting word clouds and textwrap for wrapping longe
6 from wordcloud import WordCloud, STOPWORDS
7 from textwrap import wrap
8 # For visualizations
9 import matplotlib.pyplot as plt
```

0. Loading and preparing data

```
1 # import pandas as pd
2 # jsonObj = pd.read_json(path_or_buf= "/content/drive/MyDrive/DATA SCIENCE/CAS
1 # jsonObj

Automatic saving failed. This file was updated remotely or in another tab. Show
```

	id	conversation_id	text	speaker	meta	reply- to	timesta
0	L1045	L1044	They do not!	u0	{'movie_id': 'm0', 'parsed': [{'rt': 1, 'toks'	L1044	N
1	L1044	L1044	They do to!	u2	{'movie_id': 'm0', 'parsed':	None	N
text_data text_data		nObj["text"]					
0 1 2 3 4				The I S	do not! y do to! hope so. he okay? et's go.		
304708 304709 304710 304711 304712 Name: tex	I'm to	Chelmsford seems take the Sikaliones, yes, Mr Verel Durnford Wigth: 304713, dtyp	with the r Your ord eker. Gent lliam Verel	main colu ders, Mr lemen who	mn to Vereker? can		

1 # text_data.to_csv("/content/drive/MyDrive/DATA SCIENCE/CASE STUDY 2/Movie_quo

1 movie_quotes = pd.read_csv("/content/drive/MyDrive/DATA SCIENCE/CASE STUDY 2/M
2 movie_quotes

Automatic saving failed. This file was updated remotely or in another tab.

Show diff

- 1. EDA

```
They do to!

1 df = movie_quotes
2 df.head()

text

0 They do not!

1 They do to!

2 I hope so.

3 She okay?

4 Let's go.

Colonel Durnford William Vereker I hear you
```

1.1 Checking for null values and removing them

```
1 df.isnull().sum()
    text    267
    dtype: int64

1 df.dropna(subset=['text'],inplace = True)
2 df.isnull().sum()
    text    0
    dtype: int64
```

There were 267 dialogues, that were null, so we removed them from our dataset.

1.2 Removing contractions

```
1 ' '.join(df['text'].tolist())

'They do not! They do to! I hope so. She okay? Let\'s qo. Wow Okay -- you\'re
Automatic saving failed. This file was updated remotely or in another tab. Show

diff
wca.lng pastes. The rear you. Mind good stuff eventually. Thank God! If I had to hear one more story about your coiffure. Me. This endless ...hlonde habble. T\'m like horing myself. What I df['text_new'] = df.apply(lambda row : contractions.fix(row['text']), axis = 1

1 ' '.join(df['text_new'].tolist())
```

'They do not! They do to! I hope so. She okay? Let us go. Wow Okay -- you are going to need to learn how to lie. No I am kidding. You know how sometimes y ou just become this "persona"? And you do not know how to quit? Like my fear of wearing pastels? The "real you". What good stuff? I figured you would get to the good stuff eventually. Thank God! If I had to hear one more story about your coiffure. Me. This endless a blonde habble. I am like boring mys

1 df.head()

	text	text_new
0	They do not!	They do not!
1	They do to!	They do to!
2	I hope so.	I hope so.
3	She okay?	She okay?
4	Let's go.	Let us go.

1 df = df.drop(columns=['text'])
2 df.head()



Let us go.

1

1.3 Turning to lower case

```
1 df['cleaned']=df['text_new'].apply(lambda x: x.lower())
2 ' '.join(df['text_new'].tolist())

'They do not! They do to! I hope so. She okay? Let us qo. Wow Okay -- you are sometimes y Automatic saving failed. This file was updated remotely or in another tab.

Show

diff

wearing pastets. The reat year must good start. I righted you would get to the good stuff eventually. Thank God! If I had to hear one more story about your coiffure. Me. This endless a blonde habble. I am like boring mys
```

1.4 Removing punctuations

```
1 df['cleaned']=df['cleaned'].apply(lambda x: re.sub('[%s]' % re.escape(string.p
2 ' '.join(df['text_new'].tolist())
```

'They do not! They do to! I hope so. She okay? Let us go. Wow Okay -- you are going to need to learn how to lie. No I am kidding. You know how sometimes y ou just become this "persona"? And you do not know how to quit? Like my fear of wearing pastels? The "real you". What good stuff? I figured you would get to the good stuff eventually. Thank God! If I had to hear one more story about your coiffure... Me. This endless a blonde habble. I am like boring mys

1.5 Removing extra spaces

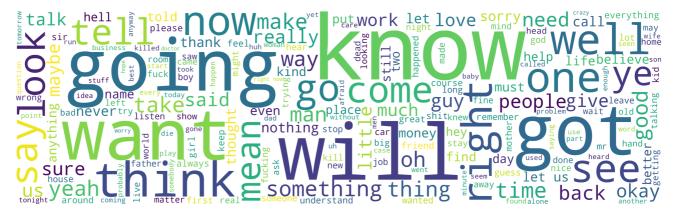
```
df['cleaned']=df['cleaned'].apply(lambda x: re.sub(' +',' ',x))
' '.join(df['cleaned'].tolist())
```

'they do not they do to i hope so she okay let us go wow okay you are going to need to learn how to lie no i am kidding you know how sometimes you just be come this persona and you do not know how to quit like my fear of wearing pastels the real you what good stuff i figured you would get to the good stuff eventually thank god if i had to hear one more story about your coiffure me this endless blonde babble i am like boring myself what crap do you listen to the

1.5 Stopwords

```
1 all words = ""
 2 for quote in df.cleaned:
   # split the value
 4 tokens = quote.split()
    all words += " ".join(tokens)+" "
 1
    stopwords = set(STOPWORDS)
 2
 3
 4
    wordcloud = WordCloud(width = 4000, height = 1200,
5
                     background color ='white',
                     stopwords = stopwords,
6
7
                     min_font_size = 10).generate(all_words)
8
9
    # plot the WordCloud image
    plt.figure(figsize = (40, 12), facecolor = None)
10
    plt.imshow(wordcloud)
11
12
    plt.axis("off")
 Automatic saving failed. This file was updated remotely or in another tab.
                                                             Show
```

 \Box



1

Automatic saving failed. This file was updated remotely or in another tab. Show diff

Os completed at 20:05