

Text Data Preprocessing
Text Data To Numerical Vector Conversion using:
Word2Vec
Pretrained GloVe
Pretrained BERT
We will follow the steps mentioned below:

Data Loading
Exploratory Data Analysis
Data Preparation
Train Test Split
Data Preprocessing (special characters, stop words, lower case, stemming, etc)
Converting text to numerical vector using Word2Vec, Pretrained GloVe and BERT
Preprocessing Test Data
Training on Train Data
Predictions on Test Data
Model Evaluation

Project Hint - Reading the Data from Database

```
In [1]: import sqlite3
import pandas as pd
import numpy as np
```

Step 1 - Reading the Tables from Database file

```
In [2]: # Read the code below and write your observation in the next cell

conn = sqlite3.connect('eng_subtitles_database.db')
cursor = conn.cursor()
cursor.execute("SELECT name FROM sqlite_master WHERE type='table'")
print(cursor.fetchall())

[('zipfiles',)]
```

In the above cell, I am able to read the table inside the database. As mentioned earlier, table name is `zipfiles`. We also know from README.txt that this table contains three columns: 'num', 'name' and 'content'.

Step 2 - Reading the columns of Table

```
In [3]: cursor.execute("PRAGMA table_info('zipfiles')")
cols = cursor.fetchall()
for col in cols:
    print(col[1])
```

```
num
name
content
```

The above code helps in checking the column names in the database table.

Let's now use `SELECT * FROM zipfiles` to read all the data into a `df` variable.

Step 3 - Loading the Database Table inside a Pandas DataFrame

```
In [4]: df = pd.read_sql_query("""SELECT * FROM zipfiles""", conn)
df.head()
```

```
Out[4]:
```

	num	name	content
0	9180533	the.message.(1976).eng.1cd	b'PK\x03\x04\x14\x00\x00\x00\x08\x00\x1c\xa9\x...
1	9180583	here.comes.the.grump.s01.e09.joltin.jack.in.bo...	b'PK\x03\x04\x14\x00\x00\x00\x08\x00\x17\xb9\x...
2	9180592	yumis.cells.s02.e13.episode.2.13.(2022).eng.1cd	b'PK\x03\x04\x14\x00\x00\x00\x08\x00L\xb9\x99V...
3	9180594	yumis.cells.s02.e14.episode.2.14.(2022).eng.1cd	b'PK\x03\x04\x14\x00\x00\x00\x08\x00U\xa9\x99V...
4	9180600	broker.(2022).eng.1cd	b'PK\x03\x04\x14\x00\x00\x00\x08\x001\xa9\x99V...

In [5]: `df.info()`

```
<class 'pandas.core.frame.DataFrame'>
RangeIndex: 82498 entries, 0 to 82497
Data columns (total 3 columns):
#   Column      Non-Null Count  Dtype
---  -
0   num         82498 non-null  int64
1   name        82498 non-null  object
2   content     82498 non-null  object
dtypes: int64(1), object(2)
memory usage: 1.9+ MB
```

Looks like the `content` column donot contain the subtitles text. Instead as mentioned in README.txt, it might be latin-1 encoded.

Step 4 - Printing content of 0th Row

In [6]: `b_data = df.iloc[0, 2]`

```
# here 2 represent the index of content column
# 0 represents the row number
```

```
print(b_data)
```

[illegible]

Worked with something similar earlier.

Step 5 - Unzipping the content of 385th row and decoding using latin-1

```
In [8]: import zipfile
import io

# Assuming 'content' is the binary data from your database
binary_data = df.iloc[385, 2]

# Decompress the binary data using the zipfile module
with io.BytesIO(binary_data) as f:
    with zipfile.ZipFile(f, 'r') as zip_file:
        # Reading only one file in the ZIP archive
        subtitle_content = zip_file.read(zip_file.namelist()[0])

# Now 'subtitle_content' should contain the extracted subtitle content
print(subtitle_content.decode('latin-1')) # Assuming the content is Latin-1 encoded text
```

```
1
00:00:06,000 --> 00:00:12,074
Watch any video online with Open-SUBTITLES
Free Browser extension: osdb.link/ext
```

```
2
00:00:15,370 --> 00:00:16,506
You lose everything, my girl.
```

```
3
00:00:16,530 --> 00:00:19,360
So you've said - four times.
```

```
4
00:00:20,330 --> 00:00:22,120
I definitely had
it on yesterday.
```

```
5
00:00:22,120 --> 00:00:25,700
```

Look's like it worked.

Step 6 - Applying the above Function on the Entire Data

```
In [9]: import zipfile
import io
import pandas as pd

def decode_method(binary_data):
    # Decompress the binary data using the zipfile module
    with io.BytesIO(binary_data) as f:
        with zipfile.ZipFile(f, 'r') as zip_file:
            # Assuming there's only one file in the ZIP archive
            subtitle_name = zip_file.namelist()[0]
            subtitle_content = zip_file.read(subtitle_name)

    # Now 'subtitle_content' should contain the extracted subtitle content
    subtitle_content_str = subtitle_content.decode('latin-1') # Assuming the content is UTF-8 encoded text

    # Create a DataFrame containing the subtitle name and content
    df = pd.DataFrame({'Subtitle_Name': [subtitle_name], 'Subtitle_Content': [subtitle_content_str]})

    return subtitle_content_str, df
```

```
In [10]: df.head()
```

```
Out[10]:
```

	num	name	content
0	9180533	the.message.(1976).eng.1cd	b'PK\x03\x04\x14\x00\x00\x00\x08\x00\x1c\xa9\x...
1	9180583	here.comes.the.grump.s01.e09.joltin.jack.in.bo...	b'PK\x03\x04\x14\x00\x00\x00\x08\x00\x17\xb9\x...
2	9180592	yumis.cells.s02.e13.episode.2.13.(2022).eng.1cd	b'PK\x03\x04\x14\x00\x00\x00\x08\x00L\xb9\x99V...
3	9180594	yumis.cells.s02.e14.episode.2.14.(2022).eng.1cd	b'PK\x03\x04\x14\x00\x00\x00\x08\x00U\xa9\x99V...
4	9180600	broker.(2022).eng.1cd	b'PK\x03\x04\x14\x00\x00\x00\x08\x001\xa9\x99V...

```
In [11]: import zipfile
import io

count = 0

def decode_method(binary_data):
    global count
    # Decompress the binary data using the zipfile module
    # print(count, end=" ")
    count += 1
    with io.BytesIO(binary_data) as f:
        with zipfile.ZipFile(f, 'r') as zip_file:
            # Assuming there's only one file in the ZIP archive
            subtitle_content = zip_file.read(zip_file.namelist()[0])

    # Now 'subtitle_content' should contain the extracted subtitle content
    return subtitle_content.decode('latin-1') # Assuming the content is UTF-8 encoded text
```

```
In [12]: df['file_content'] = df['content'].apply(decode_method)

df.head()
```

```
Out[12]:
```

	num	name	content	file_content
0	9180533	the.message.(1976).eng.1cd	b'PK\x03\x04\x14\x00\x00\x00\x08\x00\x1c\xa9\x...	1\r\n00:00:06,000 --> 00:00:12,074\r\nWatch an...
1	9180583	here.comes.the.grump.s01.e09.joltin.jack.in.bo...	b'PK\x03\x04\x14\x00\x00\x00\x08\x00\x17\xb9\x...	1\r\n00:00:29,359 --> 00:00:32,048\r\nAh! Ther...
2	9180592	yumis.cells.s02.e13.episode.2.13.(2022).eng.1cd	b'PK\x03\x04\x14\x00\x00\x00\x08\x00L\xb9\x99V...	1\r\n00:00:53,200 --> 00:00:56,030\r\n<i>Yumi'...
3	9180594	yumis.cells.s02.e14.episode.2.14.(2022).eng.1cd	b'PK\x03\x04\x14\x00\x00\x00\x08\x00U\xa9\x99V...	1\r\n00:00:06,000 --> 00:00:12,074\r\nWatch an...
4	9180600	broker.(2022).eng.1cd	b'PK\x03\x04\x14\x00\x00\x00\x08\x001\xa9\x99V...	ï»¿1\r\n00:00:06,000 --> 00:00:12,074\r\nWatch...

In [13]: df.info()

```
<class 'pandas.core.frame.DataFrame'>
RangeIndex: 82498 entries, 0 to 82497
Data columns (total 4 columns):
#   Column          Non-Null Count  Dtype
---  -
0   num              82498 non-null  int64
1   name             82498 non-null  object
2   content          82498 non-null  object
3   file_content     82498 non-null  object
dtypes: int64(1), object(3)
memory usage: 2.5+ MB
```

In [14]: df.tail()

Out[14]:

	num	name	content	file_content
82493	9521935	the.prophets.game.(2000).eng.1cd	b'PK\x03\x04\x14\x00\x00\x00\x08\x00\xb8\xa6\x...	ï»¿1\r\n00:01:16,284 --> 00:01:19,537\r\nGod,\...
82494	9521937	west.beirut.(1998).eng.1cd	b'PK\x03\x04\x14\x00\x00\x00\x08\x00\x13\x97\x...	1\r\n00:00:06,000 --> 00:00:12,074\r\napi.Open...
82495	9521938	frankenstein.the.true.story.(1973).eng.1cd	b'PK\x03\x04\x14\x00\x00\x00\x08\x00\$\x97\x9aV...	1\r\n00:00:01,001 --> 00:00:04,630\r\n(Dramati...
82496	9521940	frankenstein.the.true.story.(1973).eng.1cd	b'PK\x03\x04\x14\x00\x00\x00\x08\x00\x00\x97\x...	1\r\n00:00:06,000 --> 00:00:12,074\r\nAdvertis...
82497	9521941	zombie.island.massacre.(1984).eng.1cd	b'PK\x03\x04\x14\x00\x00\x00\x08\x00,\x97\x9aV...	1\r\n00:00:01,919 --> 00:00:03,253\r\n(Sharp w...

Using Random 30k Sample

In [15]: *# Replace df with the name of your DataFrame if it's different*
df = df.sample(n=30000, random_state=1).reset_index(drop=True)

In [16]: df.info()

```
<class 'pandas.core.frame.DataFrame'>
RangeIndex: 30000 entries, 0 to 29999
Data columns (total 4 columns):
#   Column          Non-Null Count  Dtype
---  -
0   num              30000 non-null  int64
1   name             30000 non-null  object
2   content          30000 non-null  object
3   file_content     30000 non-null  object
dtypes: int64(1), object(3)
memory usage: 937.6+ KB
```

In [17]: df.head()

Out[17]:

	num	name	content	file_content
0	9487821	the.client.list.s02.e06.unanswered.prayers.(20...	b'PK\x03\x04\x14\x00\x00\x00\x08\x00\x86}\x9aV...	ĩ»¿1\r\n00:00:01,035 --> 00:00:02,536\r\nRILEY...
1	9183678	the.devil.is.a.woman.(1935).eng.1cd	b'PK\x03\x04\x14\x00\x00\x00\x08\x00o\xa9\x99V...	1\r\n00:00:06,000 --> 00:00:12,074\r\nSupport ...
2	9258855	kingdom.s04.e18.rivercrossing.battle.(2022).en...	b'PK\x03\x04\x14\x00\x00\x00\x08\x00\xd9\x9b\x...	[Script Info]\r\nTitle: English (US)\r\nOrigin...
3	9421238	sorority.babes.in.the.slimeball.bowlorama.2.(2...	b'PK\x03\x04\x14\x00\x00\x00\x08\x0014\x9aV#\x...	1\r\n00:00:00,870 --> 00:00:02,306\r\n[wind wh...
4	9518525	ada.twist.scientist.s03.e03.dadbotinto.the.bra...	b'PK\x03\x04\x14\x00\x00\x00\x08\x00n\x95\x9a...	ĩ»¿1\r\n00:00:06,000 --> 00:00:12,074\r\nWatch...

In [18]: df.head()

Out[18]:

	num	name	content	file_content
0	9487821	the.client.list.s02.e06.unanswered.prayers.(20...	b'PK\x03\x04\x14\x00\x00\x00\x08\x00\x86}\x9aV...	ĩ»¿1\r\n00:00:01,035 --> 00:00:02,536\r\nRILEY...
1	9183678	the.devil.is.a.woman.(1935).eng.1cd	b'PK\x03\x04\x14\x00\x00\x00\x08\x00o\xa9\x99V...	1\r\n00:00:06,000 --> 00:00:12,074\r\nSupport ...
2	9258855	kingdom.s04.e18.rivercrossing.battle.(2022).en...	b'PK\x03\x04\x14\x00\x00\x00\x08\x00\xd9\x9b\x...	[Script Info]\r\nTitle: English (US)\r\nOrigin...
3	9421238	sorority.babes.in.the.slimeball.bowlorama.2.(2...	b'PK\x03\x04\x14\x00\x00\x00\x08\x0014\x9aV#\x...	1\r\n00:00:00,870 --> 00:00:02,306\r\n[wind wh...
4	9518525	ada.twist.scientist.s03.e03.dadbotinto.the.bra...	b'PK\x03\x04\x14\x00\x00\x00\x08\x00n\x95\x9a...	ĩ»¿1\r\n00:00:06,000 --> 00:00:12,074\r\nWatch...

In [19]: df.shape

Out[19]: (30000, 4)

Cleaning the data

Data Prepration - Text Preprocessing

```
In [20]: import re

# Function to clean text
def clean_text(text):
    # Remove timestamps and other patterns (assuming timestamps are like [00:00:00])
    text = re.sub(r'\[\d+:\d+:\d+\]', '', text)

    # Remove special characters, numbers, and extra spaces
    text = re.sub(r'^a-zA-Z\s', '', text)

    # Replace dots with space
    text = re.sub(r'\.', ' ', text)

    # Convert to lowercase
    text = text.lower()

    # Remove unnecessary whitespace characters
    text = re.sub(r'\s+', ' ', text).strip()

    return text
```

```
In [22]: from tqdm import tqdm
tqdm.pandas()
```

Out[24]:	num	name	content	file_content	cleaned_file_content
0	9487821	the.client.list.s02.e06.unanswered.prayers.(20...	b'PK\x03\x04\x14\x00\x00\x00\x08\x00\x86}\x9aV...	ĩ»¿1\r\n00:00:01,035 --> 00:00:02,536\r\nRILEY...	riley ipreviously oni the client list derek yo...
1	9183678	the.devil.is.a.woman.(1935).eng.1cd	b'PK\x03\x04\x14\x00\x00\x00\x08\x00o\xa9\x99V...	1\r\n00:00:06,000 --> 00:00:12,074\r\nSupport ...	support us and become vip member to remove all...
2	9258855	kingdom.s04.e18.rivercrossing.battle.(2022).en...	b'PK\x03\x04\x14\x00\x00\x00\x08\x00\xd9\x9b\x...	[Script Info]\r\nTitle: English (US)\r\nOrigin...	script info title english us original script p...
3	9421238	sorority.babes.in.the.slimeball.bowlorama.2.(2...	b'PK\x03\x04\x14\x00\x00\x00\x08\x0014\x9aV#\x...	1\r\n00:00:00,870 --> 00:00:02,306\r\n[wind wh...	wind whooshing eerie synth music bats chirping...
4	9518525	ada.twist.scientist.s03.e03.dadbotinto.the.bra...	b'PK\x03\x04\x14\x00\x00\x00\x08\x00n\x95\x9a...	ĩ»¿1\r\n00:00:06,000 --> 00:00:12,074\r\nWatch...	watch any video online with opensubtitles free...

```
In [25]: df.drop(columns=['content', 'file_content'], inplace=True)
```

```
In [26]: df.head()
```

```
Out[26]:
```

	num	name	cleaned_file_content
0	9487821	the.client.list.s02.e06.unanswered.prayers.(20...	riley ipreviously oni the client list derek yo...
1	9183678	the.devil.is.a.woman.(1935).eng.1cd	support us and become vip member to remove all...
2	9258855	kingdom.s04.e18.rivercrossing.battle.(2022).en...	script info title english us original script p...
3	9421238	sorority.babes.in.the.slimeball.bowlorama.2.(2...	wind whooshing eerie synth music bats chirping...
4	9518525	ada.twist.scientist.s03.e03.dadbotinto.the.bra...	watch any video online with opensubtitles free...

```
In [27]: df.info()
```

```
<class 'pandas.core.frame.DataFrame'>
RangeIndex: 30000 entries, 0 to 29999
Data columns (total 3 columns):
#   Column                Non-Null Count  Dtype
---  -
0   num                   30000 non-null  int64
1   name                  30000 non-null  object
2   cleaned_file_content  30000 non-null  object
dtypes: int64(1), object(2)
memory usage: 703.3+ KB
```

```
In [28]: df.head()
```

```
Out[28]:
```

	num	name	cleaned_file_content
0	9487821	the.client.list.s02.e06.unanswered.prayers.(20...	riley ipreviously oni the client list derek yo...
1	9183678	the.devil.is.a.woman.(1935).eng.1cd	support us and become vip member to remove all...
2	9258855	kingdom.s04.e18.rivercrossing.battle.(2022).en...	script info title english us original script p...
3	9421238	sorority.babes.in.the.slimeball.bowlorama.2.(2...	wind whooshing eerie synth music bats chirping...
4	9518525	ada.twist.scientist.s03.e03.dadbotinto.the.bra...	watch any video online with opensubtitles free...

```
In [1]: # Save the data and Importing
```

```
In [29]: df.to_csv('clean_data.csv')
```

```
In [ ]: import pandas as pd
```

```
In [30]: df1 = pd.read_csv('clean_data.csv')
```

```
In [31]: df1.shape
```

```
Out[31]: (30000, 4)
```

```
In [32]: df1.head()
```

```
Out[32]:
```

	Unnamed: 0	num	name	cleaned_file_content
0	0	9487821	the.client.list.s02.e06.unanswered.prayers.(20...	riley ipreviously oni the client list derek yo...
1	1	9183678	the.devil.is.a.woman.(1935).eng.1cd	support us and become vip member to remove all...
2	2	9258855	kingdom.s04.e18.rivercrossing.battle.(2022).en...	script info title english us original script p...
3	3	9421238	sorority.babes.in.the.slimeball.bowlorama.2.(2...	wind whooshing eerie synth music bats chirping...
4	4	9518525	ada.twist.scientist.s03.e03.dadbotinto.the.bra...	watch any video online with opensubtitles free...

```
In [33]: df1.drop(columns='Unnamed: 0', inplace =True)
```

```
In [34]: df1.shape
```

```
Out[34]: (30000, 3)
```

```
In [35]: df1.info()
```

```
<class 'pandas.core.frame.DataFrame'>
RangeIndex: 30000 entries, 0 to 29999
Data columns (total 3 columns):
#   Column                Non-Null Count  Dtype
---  -
0   num                    30000 non-null  int64
1   name                   30000 non-null  object
2   cleaned_file_content   30000 non-null  object
dtypes: int64(1), object(2)
memory usage: 703.3+ KB
```

Converting Text to Numerical vectors - Word2Vec Representation

Step 1 - Import Word2Vec module from gensim.models

Step 2 - Convert the sentences to the List of Words (i.e. List of Tokens)

Step 3 - Use Word2Vec to learn numerical vectors for each unique words. Word2Vec uses the list of tokens and generate 300Dimensional numerical vector for each unique word.

Step 4 - Convert the word vectors to document vectors.

```
In [36]: import re
import nltk
from nltk.corpus import stopwords
from nltk.tokenize import word_tokenize
from nltk.stem import WordNetLemmatizer
```

```
In [37]: # Download NLTK resources
nltk.download('punkt')
nltk.download('stopwords')
nltk.download('wordnet')
```

```
[nltk_data] Downloading package punkt to
[nltk_data] C:\Users\manch\AppData\Roaming\nltk_data...
[nltk_data] Package punkt is already up-to-date!
[nltk_data] Downloading package stopwords to
[nltk_data] C:\Users\manch\AppData\Roaming\nltk_data...
[nltk_data] Package stopwords is already up-to-date!
[nltk_data] Downloading package wordnet to
[nltk_data] C:\Users\manch\AppData\Roaming\nltk_data...
[nltk_data] Package wordnet is already up-to-date!
```

Out[37]: True

```
In [38]: from nltk.stem import WordNetLemmatizer, PorterStemmer # Import PorterStemmer

# Initialize the Lemmatizer and PorterStemmer
lemmatizer = WordNetLemmatizer()
stemmer = PorterStemmer()
```

```
In [39]: # Initialize the Lemmatizer
lemmatizer = WordNetLemmatizer()

# Define a function to remove stop words, tokenize, and Lemmatize text
def tokenize_text(text):
    # Tokenize the text
    tokens = word_tokenize(text)

    # Remove stop words
    stop_words = set(stopwords.words('english'))
    tokens = [token for token in tokens if token.lower() not in stop_words]

    # Lemmatize tokens
    tokens = [lemmatizer.lemmatize(token) for token in tokens]

    return tokens
```

```
In [ ]: # 15 mins
```

```
In [40]: # Tokenize, remove stop words, and Lemmatize the 'cleaned_file_content' column
df1['tokenised_content_lemma'] = df1['cleaned_file_content'].progress_apply(tokenize_text)

df1['content_size_lemma'] = df1['tokenised_content_lemma'].progress_apply(len)
```

```
100%|████████████████████████████████████████████████████████████████████████████████| 30000/30000 [13:56<00:00, 35.87it/s]
100%|████████████████████████████████████████████████████████████████████████████████| 30000/30000 [00:00<00:00, 830171.67it/s]
```

```
In [41]: # Display the tokenized content
print(df1[['tokenised_content_lemma', 'content_size_lemma']])
```

	tokenised_content_lemma	content_size_lemma
0	[riley, ipreviously, oni, client, list, derek,...	3244
1	[support, u, become, vip, member, remove, ad, ...	2853
2	[script, info, title, english, u, original, sc...	1732
3	[wind, whooshing, eerie, synth, music, bat, ch...	2948
4	[watch, video, online, opensubtitles, free, br...	1482
...
29995	[please, attempt, perform, stunt, activity, sh...	3055
29996	[cricket, chirring, music, box, music, apiopen...	1805
29997	[previously, oni, velma, last, place, mom, cel...	1797
29998	[one, two, three, four, five, numberblocks, si...	267
29999	[kid, capri, russell, simmons, def, comedy, ja...	1578

[30000 rows x 2 columns]

```
In [42]: df1.head()
```

Out[42]:

	num	name	cleaned_file_content	tokenised_content_lemma	content_size_lemma
0	9487821	the.client.list.s02.e06.unanswered.prayers.(20...	riley ipreviously oni the client list derek yo...	[riley, ipreviously, oni, client, list, derek,...	3244
1	9183678	the.devil.is.a.woman.(1935).eng.1cd	support us and become vip member to remove all...	[support, u, become, vip, member, remove, ad, ...	2853
2	9258855	kingdom.s04.e18.rivercrossing.battle.(2022).en...	script info title english us original script p...	[script, info, title, english, u, original, sc...	1732
3	9421238	sorority.babes.in.the.slimeball.bowlorama.2.(2...	wind whooshing eerie synth music bats chirping...	[wind, whooshing, eerie, synth, music, bat, ch...	2948
4	9518525	ada.twist.scientist.s03.e03.dadbotinto.the.bra...	watch any video online with opensubtitles free...	[watch, video, online, opensubtitles, free, br...	1482

chunking the data


```

In [73]: import pandas as pd

# Assuming df contains a column named 'srt_content' containing the content of each .srt file
# You may need to adjust column names accordingly

# Define a function to create overlapping chunks for each .srt file content
def overlapping_chunks(cleaned_file_content, chunk_size=4, overlap_size=50):
    chunks = []
    lines = cleaned_file_content.split('\n')
    total_lines = len(lines)
    for i in range(0, total_lines - chunk_size + 1, chunk_size - overlap_size):
        chunk = '\n'.join(lines[i:i + chunk_size])
        chunks.append(chunk)
    return chunks

# Apply the overlapping_chunks function to each .srt file content in the DataFrame
df1['overlapping_chunks_file_content'] = df1['cleaned_file_content'].apply(lambda x: overlapping_chunks(x))

# Now df['overlapping_chunks'] contains a List of overlapping chunks for each .srt file
# Each chunk is a string containing a part of the .srt content with overlap
df1.head()

```

Out[73]:

	name	file_content	cleaned_file_content	tokenised_content_lemma	content_size_lemma	overlapping_chunks_file_cont
0	the.client.list.s02.e06.unanswered.prayers. (20...	ĩ»¿ 1\\n00:00:01,035 --> 00:00:02,536\\nRILEY...	riley ipreviously oni the client list derek yo...	[riley, ipreviously, oni, client, list, derek,...	3244	[riley ipreviously oni the client derek
1	the.devil.is.a.woman.(1935).eng.1cd	1\\n00:00:06,000 --> 00:00:12,074\\nSupport ...	support us and become vip member to remove all...	[support, u, become, vip, member, remove, ad, ...	2853	[support us and become member to remove :
2	kingdom.s04.e18.rivercrossing.battle. (2022).en...	[Script Info]\\nTitle: English (US)\\nOrigin...	script info title english us original script p...	[script, info, title, english, u, original, sc...	1732	[script info title english us orig scrip
3	sorority.babes.in.the.slimeball.bowlorama.2. (2...	1\\n00:00:00,870 --> 00:00:02,306\\n\\n[wind wh...	wind whooshing eerie synth music bats chirping...	[wind, whooshing, eerie, synth, music, bat, ch...	2948	[wind whooshing eerie synth m bats chirp
4	ada.twist.scientist.s03.e03.dadbotinto.the.bra...	ĩ»¿ 1\\n00:00:06,000 --> 00:00:12,074\\nWatch...	watch any video online with opensubtitles free...	[watch, video, online, opensubtitles, free, br...	1482	[watch any video online \\n opensubtitles fi

In [43]: df1.shape

Out[43]: (30000, 5)

In [44]: df1.info()

```
<class 'pandas.core.frame.DataFrame'>
RangeIndex: 30000 entries, 0 to 29999
Data columns (total 5 columns):
#   Column                                Non-Null Count  Dtype
---  -
0   num                                    30000 non-null  int64
1   name                                   30000 non-null  object
2   cleaned_file_content                  30000 non-null  object
3   tokenised_content_lemma               30000 non-null  object
4   content_size_lemma                   30000 non-null  int64
dtypes: int64(2), object(3)
memory usage: 1.1+ MB
```

In [45]: df1.to_csv('lemma_data.csv')

In []: df1.info()

In []: import pandas as pd

In [46]: df2 = pd.read_csv('lemma_data.csv')

In [47]: df2.shape

Out[47]: (30000, 6)

In [48]: df2.info()

```
<class 'pandas.core.frame.DataFrame'>
RangeIndex: 30000 entries, 0 to 29999
Data columns (total 6 columns):
#   Column                                Non-Null Count  Dtype
---  -
0   Unnamed: 0                            30000 non-null  int64
1   num                                    30000 non-null  int64
2   name                                   30000 non-null  object
3   cleaned_file_content                  30000 non-null  object
4   tokenised_content_lemma               30000 non-null  object
5   content_size_lemma                   30000 non-null  int64
dtypes: int64(3), object(3)
memory usage: 1.4+ MB
```

In [49]: df2.head()

Out[49]:

	Unnamed: 0	num	name	cleaned_file_content	tokenised_content_lemma	content_size_lemma
0	0	9487821	the.client.list.s02.e06.unanswered.prayers.(20...	riley ipreviously oni the client list derek yo...	['riley', 'ipreviously', 'oni', 'client', 'lis...	3244
1	1	9183678	the.devil.is.a.woman.(1935).eng.1cd	support us and become vip member to remove all...	['support', 'u', 'become', 'vip', 'member', 'r...	2853
2	2	9258855	kingdom.s04.e18.rivercrossing.battle. (2022).en...	script info title english us original script p...	['script', 'info', 'title', 'english', 'u', 'o...	1732
3	3	9421238	sorority.babes.in.the.slimeball.bowlorama.2.(2...	wind whooshing eerie synth music bats chirping...	['wind', 'whooshing', 'eerie', 'synth', 'music...	2948
4	4	9518525	ada.twist.scientist.s03.e03.dadbotinto.the.bra...	watch any video online with opensubtitles free...	['watch', 'video', 'online', 'opensubtitles', ...	1482

In [51]: df2.drop(columns='Unnamed: 0', inplace =True)

```
In [52]: df2.head()
```

```
Out[52]:
```

	num	name	cleaned_file_content	tokenised_content_lemma	content_size_lemma
0	9487821	the.client.list.s02.e06.unanswered.prayers.(20...	riley ipreviously oni the client list derek yo...	['riley', 'ipreviously', 'oni', 'client', 'lis...	3244
1	9183678	the.devil.is.a.woman.(1935).eng.1cd	support us and become vip member to remove all...	['support', 'u', 'become', 'vip', 'member', 'r...	2853
2	9258855	kingdom.s04.e18.rivercrossing.battle.(2022).en...	script info title english us original script p...	['script', 'info', 'title', 'english', 'u', 'o...	1732
3	9421238	sorority.babes.in.the.slimeball.bowlorama.2.(2...	wind whooshing eerie synth music bats chirping...	['wind', 'whooshing', 'eerie', 'synth', 'music...	2948
4	9518525	ada.twist.scientist.s03.e03.dadbotinto.the.bra...	watch any video online with opensubtitles free...	['watch', 'video', 'online', 'opensubtitles', ...	1482

```
In [53]: df2.info()
```

```
<class 'pandas.core.frame.DataFrame'>
RangeIndex: 30000 entries, 0 to 29999
Data columns (total 5 columns):
#   Column                Non-Null Count  Dtype
---  -
0   num                    30000 non-null  int64
1   name                   30000 non-null  object
2   cleaned_file_content   30000 non-null  object
3   tokenised_content_lemma 30000 non-null  object
4   content_size_lemma     30000 non-null  int64
dtypes: int64(2), object(3)
memory usage: 1.1+ MB
```

```
In [54]: import gensim
```

```
print(gensim.__version__)
from gensim.models import Word2Vec
```

4.3.2

```
In [57]: # 15 mins
```

```
In [56]: model = Word2Vec(list(df2.tokenised_content_lemma), vector_size=300, min_count=1)

print(model)
```

Word2Vec<vocab=31, vector_size=300, alpha=0.025>

In [64]: df2.info()

```
<class 'pandas.core.frame.DataFrame'>
RangeIndex: 30000 entries, 0 to 29999
Data columns (total 6 columns):
 #   Column                Non-Null Count  Dtype
---  -
 0   num                   30000 non-null  int64
 1   name                  30000 non-null  object
 2   cleaned_file_content  30000 non-null  object
 3   tokenised_content_lemma 30000 non-null  object
 4   content_size_lemma     30000 non-null  int64
 5   doc_vector            30000 non-null  object
dtypes: int64(2), object(4)
memory usage: 1.4+ MB
```

In [65]: df2.to_csv('w2v_data.csv')

In []:

Pretrained BERT for Sentence Vectors

In [1]: `import` pandas `as` pd

In [2]: df3 = pd.read_csv('w2v_data.csv')

In [3]: df3.shape

Out[3]: (30000, 7)

```
In [4]: df3.head()
```

Out[4]:

	Unnamed: 0	num	name	cleaned_file_content	tokenised_content_lemma	content_size_lemma	doc_vector
0	0	9487821	the.client.list.s02.e06.unanswered.prayers. (20...	riley ipreviously oni the client list derek yo...	['riley', 'ipreviously', 'oni', 'client', 'lis...	3244	[-1.09343290e-01 -7.68003017e-02 -1.80567741e-...
1	1	9183678	the.devil.is.a.woman.(1935).eng.1cd	support us and become vip member to remove all...	['support', 'u', 'become', 'vip', 'member', 'r...	2853	[-9.83889028e-02 -7.53047839e-02 -1.65519699e-...
2	2	9258855	kingdom.s04.e18.rivercrossing.battle. (2022).en...	script info title english us original script p...	['script', 'info', 'title', 'english', 'u', 'o...	1732	[-4.89264429e-02 -1.39095232e-01 -1.84526831e-...
3	3	9421238	sorority.babes.in.the.slimeball.bowlorama.2. (2...	wind whooshing eerie synth music bats chirping...	['wind', 'whooshing', 'eerie', 'synth', 'music...	2948	[-9.89964455e-02 -6.88564703e-02 -1.65541591e-...
4	4	9518525	ada.twist.scientist.s03.e03.dadbotinto.the.bra...	watch any video online with opensubtitles free...	['watch', 'video', 'online', 'opensubtitles', ...	1482	[-0.10334199 -0.06256191 -0.15584423 -0.213279...

```
In [5]: df3.drop(columns='Unnamed: 0', inplace =True)
```

```
In [6]: df3.info()

<class 'pandas.core.frame.DataFrame'>
RangeIndex: 30000 entries, 0 to 29999
Data columns (total 6 columns):
#   Column                                Non-Null Count  Dtype
---  -
0   num                                    30000 non-null  int64
1   name                                  30000 non-null  object
2   cleaned_file_content                  30000 non-null  object
3   tokenised_content_lemma               30000 non-null  object
4   content_size_lemma                    30000 non-null  int64
5   doc_vector                            30000 non-null  object
dtypes: int64(2), object(4)
memory usage: 1.4+ MB
```

```
In [7]: df3.head()
```

```
Out[7]:
```

	num	name	cleaned_file_content	tokenised_content_lemma	content_size_lemma	doc_vector
0	9487821	the.client.list.s02.e06.unanswered.prayers. (20...	riley ipreviously oni the client list derek yo...	['riley', 'ipreviously', 'oni', 'client', 'lis...	3244	[-1.09343290e-01 -7.68003017e-02 -1.80567741e-...
1	9183678	the.devil.is.a.woman.(1935).eng.1cd	support us and become vip member to remove all...	['support', 'u', 'become', 'vip', 'member', 'r...	2853	[-9.83889028e-02 -7.53047839e-02 -1.65519699e-...
2	9258855	kingdom.s04.e18.rivercrossing.battle. (2022).en...	script info title english us original script p...	['script', 'info', 'title', 'english', 'u', 'o...	1732	[-4.89264429e-02 -1.39095232e-01 -1.84526831e-...
3	9421238	sorority.babes.in.the.slimeball.bowlorama.2.(2...	wind whooshing eerie synth music bats chirping...	['wind', 'whooshing', 'eerie', 'synth', 'music...	2948	[-9.89964455e-02 -6.88564703e-02 -1.66541591e-...
4	9518525	ada.twist.scientist.s03.e03.dadbotinto.the.bra...	watch any video online with opensubtitles free...	['watch', 'video', 'online', 'opensubtitles', ...	1482	[-0.10334199 -0.06256191 -0.15584423 -0.213279...

```
In [8]: from sentence_transformers import SentenceTransformer, util  
  
model = SentenceTransformer('all-MiniLM-L6-v2')
```

```
In [9]: from tqdm import tqdm  
tqdm.pandas()
```

```
In [10]: # from sentence_transformers import SentenceTransformer, util  
  
# model = SentenceTransformer('all-MiniLM-L6-v2')  
def encode_text(text):  
    embedding = model.encode(text)  
    return embedding
```



```
In [11]: df11 = df3.loc[:2000,:]
df11.head()
```

Out[11]:

	num	name	cleaned_file_content	tokenised_content_lemma	content_size_lemma	doc_vector
0	9487821	the.client.list.s02.e06.unanswered.prayers. (20...	riley ipreviously oni the client list derek yo...	['riley', 'ipreviously', 'oni', 'client', 'lis...	3244	[-1.09343290e-01 -7.68003017e-02 -1.80567741e-...
1	9183678	the.devil.is.a.woman.(1935).eng.1cd	support us and become vip member to remove all...	['support', 'u', 'become', 'vip', 'member', 'r...	2853	[-9.83889028e-02 -7.53047839e-02 -1.65519699e-...
2	9258855	kingdom.s04.e18.rivercrossing.battle. (2022).en...	script info title english us original script p...	['script', 'info', 'title', 'english', 'u', 'o...	1732	[-4.89264429e-02 -1.39095232e-01 -1.84526831e-...
3	9421238	sorority.babes.in.the.slimeball.bowlorama.2.(2...	wind whooshing eerie synth music bats chirping...	['wind', 'whooshing', 'eerie', 'synth', 'music...	2948	[-9.89964455e-02 -6.88564703e-02 -1.66541591e-...
4	9518525	ada.twist.scientist.s03.e03.dadbotinto.the.bra...	watch any video online with opensubtitles free...	['watch', 'video', 'online', 'opensubtitles', ...	1482	[-0.10334199 -0.06256191 -0.15584423 -0.213279...

```
In [12]: df11['sentance_content_bert'] = df11['cleaned_file_content'].progress_apply(encode_text)
df11.head()
```

```
100%|████████████████████████████████████████████████████████████████████████████████| 2001/2001 [11:07<00:00, 3.00it/s]
C:\Users\manch\AppData\Local\Temp\ipykernel_9688\712963940.py:1: SettingWithCopyWarning:
A value is trying to be set on a copy of a slice from a DataFrame.
Try using .loc[row_indexer,col_indexer] = value instead
```

See the caveats in the documentation: https://pandas.pydata.org/pandas-docs/stable/user_guide/indexing.html#returning-a-view-versus-a-copy (https://pandas.pydata.org/pandas-docs/stable/user_guide/indexing.html#returning-a-view-versus-a-copy)

```
df11['sentance_content_bert'] = df11['cleaned_file_content'].progress_apply(encode_text)
```

Out[12]:

	num	name	cleaned_file_content	tokenised_content_lemma	content_size_lemma	doc_vector	sentance_content_bert
0	9487821	the.client.list.s02.e06.unanswered.prayers. (20...	riley ipreviously oni the client list derek yo...	['riley', 'ipreviously', 'oni', 'client', 'lis...	3244	[-1.09343290e-01 -7.68003017e-02 -1.80567741e-...	[-0.11649825, -0.044796105, 0.026207969, -0.10...
1	9183678	the.devil.is.a.woman.(1935).eng.1cd	support us and become vip member to remove all...	['support', 'u', 'become', 'vip', 'member', 'r...	2853	[-9.83889028e-02 -7.53047839e-02 -1.65519699e-...	[-0.017634636, -0.06649098, 0.008615398, -0.07...
2	9258855	kingdom.s04.e18.rivercrossing.battle. (2022).en...	script info title english us original script p...	['script', 'info', 'title', 'english', 'u', 'o...	1732	[-4.89264429e-02 -1.39095232e-01 -1.84526831e-...	[-0.017091548, -0.044068865, 0.04417335, -0.04...
3	9421238	sorority.babes.in.the.slimeball.bowlorama.2. (2...	wind whooshing eerie synth music bats chirping...	['wind', 'whooshing', 'eerie', 'synth', 'music...	2948	[-9.89964455e-02 -6.88564703e-02 -1.66541591e-...	[-0.06590253, -0.02547965, 0.048389476, -0.016...
4	9518525	ada.twist.scientist.s03.e03.dadbotinto.the.bra...	watch any video online with opensubtitles free...	['watch', 'video', 'online', 'opensubtitles', ...	1482	[-0.10334199 -0.06256191 -0.15584423 -0.213279...	[-0.068130724, -0.024480091, -0.025265323, -0....

```
In [13]: df12 = df3.loc[2001:4000,:]
df12.head()
```

Out[13]:

	num	name	cleaned_file_content	tokenised_content_lemma	content_size_lemma	doc_vector
2001	9242706	the.yakuzas.guide.to.babysitting.s01.e11.firew...	script info title english us original script b...	['script', 'info', 'title', 'english', 'u', 'o...	1529	[-3.07599939e-02 -1.56279713e-01 -1.56834319e-...
2002	9428559	the.lord.of.the.rings.the.fellowship.of.the.ri...	watch any video online with opensubtitles free...	['watch', 'video', 'online', 'opensubtitles', ...	5994	[-0.10344059 -0.08081925 -0.16428243 -0.213812...
2003	9271256	bullet.train.(2022).eng.1cd	apiopensubtitlesorg is deprecated please imple...	['apiopensubtitlesorg', 'deprecated', 'please'...	5755	[-8.83942544e-02 -7.16175139e-02 -1.83909342e-...
2004	9297763	blockbuster.s01.e09.thimble.(2022).eng.1cd	all right gang huddle up this is for you guys ...	['right', 'gang', 'huddle', 'guy', 'oh', 'didn...	1989	[-8.56300741e-02 -7.75319487e-02 -1.63196027e-...
2005	9344265	the.most.beautiful.flower.(2022).eng.1cd	netflix presents biasa high school green beans...	['netflix', 'present', 'biasa', 'high', 'schoo...	1513	[-1.04019031e-01 -7.96759427e-02 -1.75259918e-...

```
In [14]: df12['sentance_content_bert'] = df12['cleaned_file_content'].progress_apply(encode_text)
df12.head()
```

```
100%|████████████████████████████████████████████████████████████████████████████████| 2000/2000 [24:19<00:00, 1.37it/s]
C:\Users\manch\AppData\Local\Temp\ipykernel_9688\2228528983.py:1: SettingWithCopyWarning:
A value is trying to be set on a copy of a slice from a DataFrame.
Try using .loc[row_indexer,col_indexer] = value instead
```

See the caveats in the documentation: https://pandas.pydata.org/pandas-docs/stable/user_guide/indexing.html#returning-a-view-versus-a-copy (https://pandas.pydata.org/pandas-docs/stable/user_guide/indexing.html#returning-a-view-versus-a-copy)

```
df12['sentance_content_bert'] = df12['cleaned_file_content'].progress_apply(encode_text)
```

Out[14]:

	num	name	cleaned_file_content	tokenised_content_lemma	content_size_lemma	doc_vector	sentance_content_bert
2001	9242706	the.yakuzas.guide.to.babysitting.s01.e11.firew...	script info title english us original script b...	['script', 'info', 'title', 'english', 'u', 'o...	1529	[-3.07599939e-02 -1.56279713e-01 -1.56834319e-...	[-0.0034559942, -0.014797437, 0.07932088, -0.0...
2002	9428559	the.lord.of.the.rings.the.fellowship.of.the.ri...	watch any video online with opensubtitles free...	['watch', 'video', 'online', 'opensubtitles', ...	5994	[-0.10344059 -0.08081925 -0.16428243 -0.213812...	[-0.06378279, 0.0143626975, -0.05514103, -0.05...
2003	9271256	bullet.train.(2022).eng.1cd	apiopensubtitlesorg is deprecatd please imple...	['apiopensubtitlesorg', 'deprecatd', 'please'...	5755	[-8.83942544e-02 -7.16175139e-02 -1.83909342e-...	[-0.03831797, -0.07563398, 0.009899806, 0.0112...
2004	9297763	blockbuster.s01.e09.thimble.(2022).eng.1cd	all right gang huddle up this is for you guys ...	['right', 'gang', 'huddle', 'guy', 'oh', 'didn...	1989	[-8.56300741e-02 -7.75319487e-02 -1.63196027e-...	[-0.13685279, 0.02369652, -0.0012757194, -0.01...
2005	9344265	the.most.beautiful.flower.(2022).eng.1cd	netflix presents biasa high school green beans...	['netflix', 'present', 'biasa', 'high', 'schoo...	1513	[-1.04019031e-01 -7.96759427e-02 -1.75259918e-...	[-0.0951997, -0.08066843, 0.031748097, -0.0581...

```
In [15]: new_df = pd.concat([df11,df12],ignore_index=True, axis=0)
new_df.head()
```

```
Out[15]:
```

	num	name	cleaned_file_content	tokenised_content_lemma	content_size_lemma	doc_vector	sentence_content_bert
0	9487821	the.client.list.s02.e06.unanswered.prayers. (20...	riley ipreviously oni the client list derek yo...	['riley', 'ipreviously', 'oni', 'client', 'lis...	3244	[-1.09343290e- 01 -7.68003017e-02 -1.80567741e-...	[-0.11649825, -0.044796105, 0.026207969, -0.10...
1	9183678	the.devil.is.a.woman.(1935).eng.1cd	support us and become vip member to remove all...	['support', 'u', 'become', 'vip', 'member', 'r...	2853	[-9.83889028e- 02 -7.53047839e-02 -1.65519699e-...	[-0.017634636, -0.06649098, 0.008615398, -0.07...
2	9258855	kingdom.s04.e18.rivercrossing.battle. (2022).en...	script info title english us original script p...	['script', 'info', 'title', 'english', 'u', 'o...	1732	[-4.89264429e- 02 -1.39095232e-01 -1.84526831e-...	[-0.017091548, -0.044068865, 0.04417335, -0.04...
3	9421238	sorority.babes.in.the.slimeball.bowlorama.2. (2...	wind whooshing eerie synth music bats chirping...	['wind', 'whooshing', 'eerie', 'synth', 'music...	2948	[-9.89964455e- 02 -6.88564703e-02 -1.66541591e-...	[-0.06590253, -0.02547965, 0.048389476, -0.016...
4	9518525	ada.twist.scientist.s03.e03.dadbotinto.the.bra...	watch any video online with opensubtitles free...	['watch', 'video', 'online', 'opensubtitles', ...	1482	[-0.10334199 -0.06256191 -0.15584423 -0.213279...	[-0.068130724, -0.024480091, -0.025265323, -0....


```
In [17]: new_df = pd.concat([new_df,df13],ignore_index=True, axis=0)
new_df.head()
```

```
Out[17]:
```

	num	name	cleaned_file_content	tokenised_content_lemma	content_size_lemma	doc_vector	sentence_content_bert
0	9487821	the.client.list.s02.e06.unanswered.prayers. (20...	riley ipreviously oni the client list derek yo...	['riley', 'ipreviously', 'oni', 'client', 'lis...	3244	[-1.09343290e- 01 -7.68003017e-02 -1.80567741e-...	[-0.11649825, -0.044796105, 0.026207969, -0.10...
1	9183678	the.devil.is.a.woman.(1935).eng.1cd	support us and become vip member to remove all...	['support', 'u', 'become', 'vip', 'member', 'r...	2853	[-9.83889028e- 02 -7.53047839e-02 -1.65519699e-...	[-0.017634636, -0.06649098, 0.008615398, -0.07...
2	9258855	kingdom.s04.e18.rivercrossing.battle. (2022).en...	script info title english us original script p...	['script', 'info', 'title', 'english', 'u', 'o...	1732	[-4.89264429e- 02 -1.39095232e-01 -1.84526831e-...	[-0.017091548, -0.044068865, 0.04417335, -0.04...
3	9421238	sorority.babes.in.the.slimeball.bowlorama.2. (2...	wind whooshing eerie synth music bats chirping...	['wind', 'whooshing', 'eerie', 'synth', 'music...	2948	[-9.89964455e- 02 -6.88564703e-02 -1.66541591e-...	[-0.06590253, -0.02547965, 0.048389476, -0.016...
4	9518525	ada.twist.scientist.s03.e03.dadbotinto.the.bra...	watch any video online with opensubtitles free...	['watch', 'video', 'online', 'opensubtitles', ...	1482	[-0.10334199 -0.06256191 -0.15584423 -0.213279...	[-0.068130724, -0.024480091, -0.025265323, -0....


```
In [19]: new_df = pd.concat([new_df,df14],ignore_index=True, axis=0)
new_df.head()
```

```
Out[19]:
```

	num	name	cleaned_file_content	tokenised_content_lemma	content_size_lemma	doc_vector	sentence_content_bert
0	9487821	the.client.list.s02.e06.unanswered.prayers. (20...	riley ipreviously oni the client list derek yo...	['riley', 'ipreviously', 'oni', 'client', 'lis...	3244	[-1.09343290e- 01 -7.68003017e-02 -1.80567741e-...	[-0.11649825, -0.044796105, 0.026207969, -0.10...
1	9183678	the.devil.is.a.woman.(1935).eng.1cd	support us and become vip member to remove all...	['support', 'u', 'become', 'vip', 'member', 'r...	2853	[-9.83889028e- 02 -7.53047839e-02 -1.65519699e-...	[-0.017634636, -0.06649098, 0.008615398, -0.07...
2	9258855	kingdom.s04.e18.rivercrossing.battle. (2022).en...	script info title english us original script p...	['script', 'info', 'title', 'english', 'u', 'o...	1732	[-4.89264429e- 02 -1.39095232e-01 -1.84526831e-...	[-0.017091548, -0.044068865, 0.04417335, -0.04...
3	9421238	sorority.babes.in.the.slimeball.bowlorama.2. (2...	wind whooshing eerie synth music bats chirping...	['wind', 'whooshing', 'eerie', 'synth', 'music...	2948	[-9.89964455e- 02 -6.88564703e-02 -1.66541591e-...	[-0.06590253, -0.02547965, 0.048389476, -0.016...
4	9518525	ada.twist.scientist.s03.e03.dadbotinto.the.bra...	watch any video online with opensubtitles free...	['watch', 'video', 'online', 'opensubtitles', ...	1482	[-0.10334199 -0.06256191 -0.15584423 -0.213279...	[-0.068130724, -0.024480091, -0.025265323, -0....

```
In [20]: df15 = df3.loc[8001:10000,:]  
df15['sentence_content_bert'] = df15['cleaned_file_content'].progress_apply(encode_text)  
df15.head()
```

```
100%|████████████████████████████████████████████████████████████████████████████████| 2000/2000 [08:24<00:00, 3.97it/s]
```

```
C:\Users\manch\AppData\Local\Temp\ipykernel_9688\845179494.py:2: SettingWithCopyWarning:
```

```
A value is trying to be set on a copy of a slice from a DataFrame.
```

```
Try using .loc[row_indexer,col_indexer] = value instead
```

See the caveats in the documentation: https://pandas.pydata.org/pandas-docs/stable/user_guide/indexing.html#returning-a-view-versus-a-copy (https://pandas.pydata.org/pandas-docs/stable/user_guide/indexing.html#returning-a-view-versus-a-copy)

```
df15['sentence_content_bert'] = df15['cleaned_file_content'].progress_apply(encode_text)
```

Out[20]:

	num		name	cleaned_file_content	tokenised_content_lemma	content_size_lemma	doc_vector	sentence_content_bert
8001	9351448		the.banshees.of.inisherin. (2022).eng.1cd	apiopensubtitlesorg is deprecated please imple...	['apiopensubtitlesorg', 'deprecated', 'please'...	4459	[-8.64537656e-02 -9.04118344e-02 -1.79184362e-...	[0.0032120578, -0.09480547, 0.052504268, -0.06...
8002	9401841		american.dad.s04.e15.wife.insurance. (2009).eng...	good morning usa i got a feeling that its gonn...	['good', 'morning', 'usa', 'got', 'feeling', '...	1585	[-0.08580665 -0.07550405 -0.16680738 -0.235261...	[-0.013176337, -0.031332716, 0.04096018, 0.041...
8003	9256158		house.of.wax.(1953).eng.1cd	watch any video online with opensubtitles free...	['watch', 'video', 'online', 'opensubtitles', ...	3869	[-9.80368853e-02 -6.34661317e-02 -1.57282323e-...	[-0.02410146, -0.031204488, 0.03174951, -0.112...
8004	9508645		kojak.s02.e03.hush.now.dont.you.die. (1974).eng...	music music music music music skip class run a...	['music', 'music', 'music', 'music', 'music', ...	2861	[-0.09006649 -0.06592161 -0.17272265 -0.235484...	[-0.050546892, -0.014168973, 0.061474536, -0.0...
8005	9219633		this.is.us.s02.e16.vegas.baby. (2018).eng.1cd	i previously oni this is us you know when i wa...	['previously', 'oni', 'u', 'know', 'kid', 'obs...	3002	[-1.10897392e-01 -7.63068497e-02 -1.73733607e-...	[-0.110696316, -0.10000415, 0.0813859, -0.0829...

```
In [21]: new_df = pd.concat([new_df,df15],ignore_index=True, axis=0)
new_df.head()
```

```
Out[21]:
```

	num	name	cleaned_file_content	tokenised_content_lemma	content_size_lemma	doc_vector	sentence_content_bert
0	9487821	the.client.list.s02.e06.unanswered.prayers. (20...	riley ipreviously oni the client list derek yo...	['riley', 'ipreviously', 'oni', 'client', 'lis...	3244	[-1.09343290e- 01 -7.68003017e-02 -1.80567741e-...	[-0.11649825, -0.044796105, 0.026207969, -0.10...
1	9183678	the.devil.is.a.woman.(1935).eng.1cd	support us and become vip member to remove all...	['support', 'u', 'become', 'vip', 'member', 'r...	2853	[-9.83889028e- 02 -7.53047839e-02 -1.65519699e-...	[-0.017634636, -0.06649098, 0.008615398, -0.07...
2	9258855	kingdom.s04.e18.rivercrossing.battle. (2022).en...	script info title english us original script p...	['script', 'info', 'title', 'english', 'u', 'o...	1732	[-4.89264429e- 02 -1.39095232e-01 -1.84526831e-...	[-0.017091548, -0.044068865, 0.04417335, -0.04...
3	9421238	sorority.babes.in.the.slimeball.bowlorama.2. (2...	wind whooshing eerie synth music bats chirping...	['wind', 'whooshing', 'eerie', 'synth', 'music...	2948	[-9.89964455e- 02 -6.88564703e-02 -1.66541591e-...	[-0.06590253, -0.02547965, 0.048389476, -0.016...
4	9518525	ada.twist.scientist.s03.e03.dadbotinto.the.bra...	watch any video online with opensubtitles free...	['watch', 'video', 'online', 'opensubtitles', ...	1482	[-0.10334199 -0.06256191 -0.15584423 -0.213279...	[-0.068130724, -0.024480091, -0.025265323, -0....


```
In [23]: new_df = pd.concat([new_df,df16],ignore_index=True, axis=0)
new_df.head()
```

```
Out[23]:
```

	num	name	cleaned_file_content	tokenised_content_lemma	content_size_lemma	doc_vector	sentence_content_bert
0	9487821	the.client.list.s02.e06.unanswered.prayers. (20...	riley ipreviously oni the client list derek yo...	['riley', 'ipreviously', 'oni', 'client', 'lis...	3244	[-1.09343290e-01 -7.68003017e-02 -1.80567741e-...	[-0.11649825, -0.044796105, 0.026207969, -0.10...
1	9183678	the.devil.is.a.woman.(1935).eng.1cd	support us and become vip member to remove all...	['support', 'u', 'become', 'vip', 'member', 'r...	2853	[-9.83889028e-02 -7.53047839e-02 -1.65519699e-...	[-0.017634636, -0.06649098, 0.008615398, -0.07...
2	9258855	kingdom.s04.e18.rivercrossing.battle. (2022).en...	script info title english us original script p...	['script', 'info', 'title', 'english', 'u', 'o...	1732	[-4.89264429e-02 -1.39095232e-01 -1.84526831e-...	[-0.017091548, -0.044068865, 0.04417335, -0.04...
3	9421238	sorority.babes.in.the.slimeball.bowlorama.2. (2...	wind whooshing eerie synth music bats chirping...	['wind', 'whooshing', 'eerie', 'synth', 'music...	2948	[-9.89964455e-02 -6.88564703e-02 -1.66541591e-...	[-0.06590253, -0.02547965, 0.048389476, -0.016...
4	9518525	ada.twist.scientist.s03.e03.dadbotinto.the.bra...	watch any video online with opensubtitles free...	['watch', 'video', 'online', 'opensubtitles', ...	1482	[-0.10334199 -0.06256191 -0.15584423 -0.213279...	[-0.068130724, -0.024480091, -0.025265323, -0....


```
In [25]: new_df = pd.concat([new_df,df17],ignore_index=True, axis=0)
new_df.head()
```

```
Out[25]:
```

	num	name	cleaned_file_content	tokenised_content_lemma	content_size_lemma	doc_vector	sentence_content_bert
0	9487821	the.client.list.s02.e06.unanswered.prayers. (20...	riley ipreviously oni the client list derek yo...	['riley', 'ipreviously', 'oni', 'client', 'lis...	3244	[-1.09343290e- 01 -7.68003017e-02 -1.80567741e-...	[-0.11649825, -0.044796105, 0.026207969, -0.10...
1	9183678	the.devil.is.a.woman.(1935).eng.1cd	support us and become vip member to remove all...	['support', 'u', 'become', 'vip', 'member', 'r...	2853	[-9.83889028e- 02 -7.53047839e-02 -1.65519699e-...	[-0.017634636, -0.06649098, 0.008615398, -0.07...
2	9258855	kingdom.s04.e18.rivercrossing.battle. (2022).en...	script info title english us original script p...	['script', 'info', 'title', 'english', 'u', 'o...	1732	[-4.89264429e- 02 -1.39095232e-01 -1.84526831e-...	[-0.017091548, -0.044068865, 0.04417335, -0.04...
3	9421238	sorority.babes.in.the.slimeball.bowlorama.2. (2...	wind whooshing eerie synth music bats chirping...	['wind', 'whooshing', 'eerie', 'synth', 'music...	2948	[-9.89964455e- 02 -6.88564703e-02 -1.66541591e-...	[-0.06590253, -0.02547965, 0.048389476, -0.016...
4	9518525	ada.twist.scientist.s03.e03.dadbotinto.the.bra...	watch any video online with opensubtitles free...	['watch', 'video', 'online', 'opensubtitles', ...	1482	[-0.10334199 -0.06256191 -0.15584423 -0.213279...	[-0.068130724, -0.024480091, -0.025265323, -0....


```
In [27]: new_df = pd.concat([new_df,df18],ignore_index=True, axis=0)
new_df.head()
```

```
Out[27]:
```

	num	name	cleaned_file_content	tokenised_content_lemma	content_size_lemma	doc_vector	sentence_content_bert
0	9487821	the.client.list.s02.e06.unanswered.prayers. (20...	riley ipreviously oni the client list derek yo...	['riley', 'ipreviously', 'oni', 'client', 'lis...	3244	[-1.09343290e-01 -7.68003017e-02 -1.80567741e-...	[-0.11649825, -0.044796105, 0.026207969, -0.10...
1	9183678	the.devil.is.a.woman.(1935).eng.1cd	support us and become vip member to remove all...	['support', 'u', 'become', 'vip', 'member', 'r...	2853	[-9.83889028e-02 -7.53047839e-02 -1.65519699e-...	[-0.017634636, -0.06649098, 0.008615398, -0.07...
2	9258855	kingdom.s04.e18.rivercrossing.battle. (2022).en...	script info title english us original script p...	['script', 'info', 'title', 'english', 'u', 'o...	1732	[-4.89264429e-02 -1.39095232e-01 -1.84526831e-...	[-0.017091548, -0.044068865, 0.04417335, -0.04...
3	9421238	sorority.babes.in.the.slimeball.bowlorama.2. (2...	wind whooshing eerie synth music bats chirping...	['wind', 'whooshing', 'eerie', 'synth', 'music...	2948	[-9.89964455e-02 -6.88564703e-02 -1.66541591e-...	[-0.06590253, -0.02547965, 0.048389476, -0.016...
4	9518525	ada.twist.scientist.s03.e03.dadbotinto.the.bra...	watch any video online with opensubtitles free...	['watch', 'video', 'online', 'opensubtitles', ...	1482	[-0.10334199 -0.06256191 -0.15584423 -0.213279...	[-0.068130724, -0.024480091, -0.025265323, -0....


```
In [29]: new_df = pd.concat([new_df,df19],ignore_index=True, axis=0)
new_df.head()
```

```
Out[29]:
```

	num	name	cleaned_file_content	tokenised_content_lemma	content_size_lemma	doc_vector	sentence_content_bert
0	9487821	the.client.list.s02.e06.unanswered.prayers. (20...	riley ipreviously oni the client list derek yo...	['riley', 'ipreviously', 'oni', 'client', 'lis...	3244	[-1.09343290e- 01 -7.68003017e-02 -1.80567741e-...	[-0.11649825, -0.044796105, 0.026207969, -0.10...
1	9183678	the.devil.is.a.woman.(1935).eng.1cd	support us and become vip member to remove all...	['support', 'u', 'become', 'vip', 'member', 'r...	2853	[-9.83889028e- 02 -7.53047839e-02 -1.65519699e-...	[-0.017634636, -0.06649098, 0.008615398, -0.07...
2	9258855	kingdom.s04.e18.rivercrossing.battle. (2022).en...	script info title english us original script p...	['script', 'info', 'title', 'english', 'u', 'o...	1732	[-4.89264429e- 02 -1.39095232e-01 -1.84526831e-...	[-0.017091548, -0.044068865, 0.04417335, -0.04...
3	9421238	sorority.babes.in.the.slimeball.bowlorama.2. (2...	wind whooshing eerie synth music bats chirping...	['wind', 'whooshing', 'eerie', 'synth', 'music...	2948	[-9.89964455e- 02 -6.88564703e-02 -1.66541591e-...	[-0.06590253, -0.02547965, 0.048389476, -0.016...
4	9518525	ada.twist.scientist.s03.e03.dadbotinto.the.bra...	watch any video online with opensubtitles free...	['watch', 'video', 'online', 'opensubtitles', ...	1482	[-0.10334199 -0.06256191 -0.15584423 -0.213279...	[-0.068130724, -0.024480091, -0.025265323, -0....


```
In [31]: new_df = pd.concat([new_df,df20],ignore_index=True, axis=0)
new_df.head()
```

```
Out[31]:
```

	num	name	cleaned_file_content	tokenised_content_lemma	content_size_lemma	doc_vector	sentence_content_bert
0	9487821	the.client.list.s02.e06.unanswered.prayers. (20...	riley ipreviously oni the client list derek yo...	['riley', 'ipreviously', 'oni', 'client', 'lis...	3244	[-1.09343290e- 01 -7.68003017e-02 -1.80567741e-...	[-0.11649825, -0.044796105, 0.026207969, -0.10...
1	9183678	the.devil.is.a.woman.(1935).eng.1cd	support us and become vip member to remove all...	['support', 'u', 'become', 'vip', 'member', 'r...	2853	[-9.83889028e- 02 -7.53047839e-02 -1.65519699e-...	[-0.017634636, -0.06649098, 0.008615398, -0.07...
2	9258855	kingdom.s04.e18.rivercrossing.battle. (2022).en...	script info title english us original script p...	['script', 'info', 'title', 'english', 'u', 'o...	1732	[-4.89264429e- 02 -1.39095232e-01 -1.84526831e-...	[-0.017091548, -0.044068865, 0.04417335, -0.04...
3	9421238	sorority.babes.in.the.slimeball.bowlorama.2. (2...	wind whooshing eerie synth music bats chirping...	['wind', 'whooshing', 'eerie', 'synth', 'music...	2948	[-9.89964455e- 02 -6.88564703e-02 -1.66541591e-...	[-0.06590253, -0.02547965, 0.048389476, -0.016...
4	9518525	ada.twist.scientist.s03.e03.dadbotinto.the.bra...	watch any video online with opensubtitles free...	['watch', 'video', 'online', 'opensubtitles', ...	1482	[-0.10334199 -0.06256191 -0.15584423 -0.213279...	[-0.068130724, -0.024480091, -0.025265323, -0....

```
In [32]: df21 = df3.loc[20001:22000,:]
df21['sentance_content_bert'] = df21['cleaned_file_content'].progress_apply(encode_text)
df21.head()
```

```
100%|████████████████████████████████████████████████████████████████████████████████| 2000/2000 [04:45<00:00, 7.01it/s]
```

```
C:\Users\manch\AppData\Local\Temp\ipykernel_9688\90145499.py:2: SettingWithCopyWarning:
```

```
A value is trying to be set on a copy of a slice from a DataFrame.
```

```
Try using .loc[row_indexer,col_indexer] = value instead
```

See the caveats in the documentation: https://pandas.pydata.org/pandas-docs/stable/user_guide/indexing.html#returning-a-view-versus-a-copy (https://pandas.pydata.org/pandas-docs/stable/user_guide/indexing.html#returning-a-view-versus-a-copy)

```
df21['sentance_content_bert'] = df21['cleaned_file_content'].progress_apply(encode_text)
```

Out[32]:

	num	name	cleaned_file_content	tokenised_content_lemma	content_size_lemma	doc_vector	sentance_content_bert
20001	9231080	flatland.s01.e05.a.gentle.rain. (2002).eng.1cd	support us and become vip member to remove all...	['support', 'u', 'become', 'vip', 'member', 'r...	1622	[-0.10122931 -0.07119595 -0.17184035 -0.226168...	[-0.09338511, -0.04040897, -0.010522487, -0.05...
20002	9332116	battle.for.saipan.(2022).eng.1cd	use the free code joinnow at wwwplayshipseu dr...	['use', 'free', 'code', 'joinnow', 'wwwplayshi...	2936	[-0.08494363 -0.06882316 -0.16988894 -0.232197...	[-0.08687571, -0.07450393, 0.060404934, -0.074...
20003	9362436	the.great.s01.e09.love.hurts. (2020).eng.1cd	apiopensubtitlesorg is deprecated please imple...	['apiopensubtitlesorg', 'deprecated', 'please'...	2764	[-0.08708873 -0.07407045 -0.16646565 -0.235027...	[-0.025401864, -0.029697467, 0.023453519, 0.00...
20004	9495498	tiny.beautiful.things.s01.e07.go. (2023).eng.1cd	advertise your product or brand here contact w...	['advertise', 'product', 'brand', 'contact', '...	1903	[-0.11263889 -0.07305092 -0.17696753 -0.229608...	[-0.1058409, -0.09009985, 0.022006258, -0.0389...
20005	9210526	survivor.s20.e10.going.down.in.flames. (2010).e...	jeff probst previously on isurvivori with thei...	['jeff', 'probst', 'previously', 'isurvivori',...	2928	[-0.10461924 -0.07693632 -0.1731256 -0.227250...	[-0.07107766, -0.028268559, 0.009724177, -0.07...

```
In [33]: new_df = pd.concat([new_df,df21],ignore_index=True, axis=0)
new_df.head()
```

```
Out[33]:
```

	num	name	cleaned_file_content	tokenised_content_lemma	content_size_lemma	doc_vector	sentence_content_bert
0	9487821	the.client.list.s02.e06.unanswered.prayers. (20...	riley ipreviously oni the client list derek yo...	['riley', 'ipreviously', 'oni', 'client', 'lis...	3244	[-1.09343290e-01 -7.68003017e-02 -1.80567741e-...	[-0.11649825, -0.044796105, 0.026207969, -0.10...
1	9183678	the.devil.is.a.woman.(1935).eng.1cd	support us and become vip member to remove all...	['support', 'u', 'become', 'vip', 'member', 'r...	2853	[-9.83889028e-02 -7.53047839e-02 -1.65519699e-...	[-0.017634636, -0.06649098, 0.008615398, -0.07...
2	9258855	kingdom.s04.e18.rivercrossing.battle. (2022).en...	script info title english us original script p...	['script', 'info', 'title', 'english', 'u', 'o...	1732	[-4.89264429e-02 -1.39095232e-01 -1.84526831e-...	[-0.017091548, -0.044068865, 0.04417335, -0.04...
3	9421238	sorority.babes.in.the.slimeball.bowlorama.2. (2...	wind whooshing eerie synth music bats chirping...	['wind', 'whooshing', 'eerie', 'synth', 'music...	2948	[-9.89964455e-02 -6.88564703e-02 -1.66541591e-...	[-0.06590253, -0.02547965, 0.048389476, -0.016...
4	9518525	ada.twist.scientist.s03.e03.dadbotinto.the.bra...	watch any video online with opensubtitles free...	['watch', 'video', 'online', 'opensubtitles', ...	1482	[-0.10334199 -0.06256191 -0.15584423 -0.213279...	[-0.068130724, -0.024480091, -0.025265323, -0....


```
In [35]: new_df = pd.concat([new_df,df22],ignore_index=True, axis=0)
new_df.head()
```

```
Out[35]:
```

	num	name	cleaned_file_content	tokenised_content_lemma	content_size_lemma	doc_vector	sentence_content_bert
0	9487821	the.client.list.s02.e06.unanswered.prayers. (20...	riley ipreviously oni the client list derek yo...	['riley', 'ipreviously', 'oni', 'client', 'lis...	3244	[-1.09343290e- 01 -7.68003017e-02 -1.80567741e-...	[-0.11649825, -0.044796105, 0.026207969, -0.10...
1	9183678	the.devil.is.a.woman.(1935).eng.1cd	support us and become vip member to remove all...	['support', 'u', 'become', 'vip', 'member', 'r...	2853	[-9.83889028e- 02 -7.53047839e-02 -1.65519699e-...	[-0.017634636, -0.06649098, 0.008615398, -0.07...
2	9258855	kingdom.s04.e18.rivercrossing.battle. (2022).en...	script info title english us original script p...	['script', 'info', 'title', 'english', 'u', 'o...	1732	[-4.89264429e- 02 -1.39095232e-01 -1.84526831e-...	[-0.017091548, -0.044068865, 0.04417335, -0.04...
3	9421238	sorority.babes.in.the.slimeball.bowlorama.2. (2...	wind whooshing eerie synth music bats chirping...	['wind', 'whooshing', 'eerie', 'synth', 'music...	2948	[-9.89964455e- 02 -6.88564703e-02 -1.66541591e-...	[-0.06590253, -0.02547965, 0.048389476, -0.016...
4	9518525	ada.twist.scientist.s03.e03.dadbotinto.the.bra...	watch any video online with opensubtitles free...	['watch', 'video', 'online', 'opensubtitles', ...	1482	[-0.10334199 -0.06256191 -0.15584423 -0.213279...	[-0.068130724, -0.024480091, -0.025265323, -0....


```
In [37]: new_df = pd.concat([new_df,df23],ignore_index=True, axis=0)
new_df.head()
```

```
Out[37]:
```

	num	name	cleaned_file_content	tokenised_content_lemma	content_size_lemma	doc_vector	sentence_content_bert
0	9487821	the.client.list.s02.e06.unanswered.prayers. (20...	riley ipreviously oni the client list derek yo...	['riley', 'ipreviously', 'oni', 'client', 'lis...	3244	[-1.09343290e-01 -7.68003017e-02 -1.80567741e-...	[-0.11649825, -0.044796105, 0.026207969, -0.10...
1	9183678	the.devil.is.a.woman.(1935).eng.1cd	support us and become vip member to remove all...	['support', 'u', 'become', 'vip', 'member', 'r...	2853	[-9.83889028e-02 -7.53047839e-02 -1.65519699e-...	[-0.017634636, -0.06649098, 0.008615398, -0.07...
2	9258855	kingdom.s04.e18.rivercrossing.battle. (2022).en...	script info title english us original script p...	['script', 'info', 'title', 'english', 'u', 'o...	1732	[-4.89264429e-02 -1.39095232e-01 -1.84526831e-...	[-0.017091548, -0.044068865, 0.04417335, -0.04...
3	9421238	sorority.babes.in.the.slimeball.bowlorama.2. (2...	wind whooshing eerie synth music bats chirping...	['wind', 'whooshing', 'eerie', 'synth', 'music...	2948	[-9.89964455e-02 -6.88564703e-02 -1.66541591e-...	[-0.06590253, -0.02547965, 0.048389476, -0.016...
4	9518525	ada.twist.scientist.s03.e03.dadbotinto.the.bra...	watch any video online with opensubtitles free...	['watch', 'video', 'online', 'opensubtitles', ...	1482	[-0.10334199 -0.06256191 -0.15584423 -0.213279...	[-0.068130724, -0.024480091, -0.025265323, -0....


```
In [39]: new_df = pd.concat([new_df,df24],ignore_index=True, axis=0)
new_df.head()
```

```
Out[39]:
```

	num	name	cleaned_file_content	tokenised_content_lemma	content_size_lemma	doc_vector	sentence_content_bert
0	9487821	the.client.list.s02.e06.unanswered.prayers. (20...	riley ipreviously oni the client list derek yo...	['riley', 'ipreviously', 'oni', 'client', 'lis...	3244	[-1.09343290e- 01 -7.68003017e-02 -1.80567741e-...	[-0.11649825, -0.044796105, 0.026207969, -0.10...
1	9183678	the.devil.is.a.woman.(1935).eng.1cd	support us and become vip member to remove all...	['support', 'u', 'become', 'vip', 'member', 'r...	2853	[-9.83889028e- 02 -7.53047839e-02 -1.65519699e-...	[-0.017634636, -0.06649098, 0.008615398, -0.07...
2	9258855	kingdom.s04.e18.rivercrossing.battle. (2022).en...	script info title english us original script p...	['script', 'info', 'title', 'english', 'u', 'o...	1732	[-4.89264429e- 02 -1.39095232e-01 -1.84526831e-...	[-0.017091548, -0.044068865, 0.04417335, -0.04...
3	9421238	sorority.babes.in.the.slimeball.bowlorama.2. (2...	wind whooshing eerie synth music bats chirping...	['wind', 'whooshing', 'eerie', 'synth', 'music...	2948	[-9.89964455e- 02 -6.88564703e-02 -1.66541591e-...	[-0.06590253, -0.02547965, 0.048389476, -0.016...
4	9518525	ada.twist.scientist.s03.e03.dadbotinto.the.bra...	watch any video online with opensubtitles free...	['watch', 'video', 'online', 'opensubtitles', ...	1482	[-0.10334199 -0.06256191 -0.15584423 -0.213279...	[-0.068130724, -0.024480091, -0.025265323, -0....


```
In [41]: new_df = pd.concat([new_df,df25],ignore_index=True, axis=0)
new_df.head()
```

```
Out[41]:
```

	num	name	cleaned_file_content	tokenised_content_lemma	content_size_lemma	doc_vector	sentence_content_bert
0	9487821	the.client.list.s02.e06.unanswered.prayers. (20...	riley ipreviously oni the client list derek yo...	['riley', 'ipreviously', 'oni', 'client', 'lis...	3244	[-1.09343290e- 01 -7.68003017e-02 -1.80567741e-...	[-0.11649825, -0.044796105, 0.026207969, -0.10...
1	9183678	the.devil.is.a.woman.(1935).eng.1cd	support us and become vip member to remove all...	['support', 'u', 'become', 'vip', 'member', 'r...	2853	[-9.83889028e- 02 -7.53047839e-02 -1.65519699e-...	[-0.017634636, -0.06649098, 0.008615398, -0.07...
2	9258855	kingdom.s04.e18.rivercrossing.battle. (2022).en...	script info title english us original script p...	['script', 'info', 'title', 'english', 'u', 'o...	1732	[-4.89264429e- 02 -1.39095232e-01 -1.84526831e-...	[-0.017091548, -0.044068865, 0.04417335, -0.04...
3	9421238	sorority.babes.in.the.slimeball.bowlorama.2. (2...	wind whooshing eerie synth music bats chirping...	['wind', 'whooshing', 'eerie', 'synth', 'music...	2948	[-9.89964455e- 02 -6.88564703e-02 -1.66541591e-...	[-0.06590253, -0.02547965, 0.048389476, -0.016...
4	9518525	ada.twist.scientist.s03.e03.dadbotinto.the.bra...	watch any video online with opensubtitles free...	['watch', 'video', 'online', 'opensubtitles', ...	1482	[-0.10334199 -0.06256191 -0.15584423 -0.213279...	[-0.068130724, -0.024480091, -0.025265323, -0....

```
In [42]: new_df.shape
```

```
Out[42]: (30000, 7)
```

```
In [43]: new_df.info()
```

```
<class 'pandas.core.frame.DataFrame'>
RangeIndex: 30000 entries, 0 to 29999
Data columns (total 7 columns):
 #   Column                Non-Null Count  Dtype
---  -
 0   num                   30000 non-null  int64
 1   name                  30000 non-null  object
 2   cleaned_file_content  30000 non-null  object
 3   tokenised_content_lemma  30000 non-null  object
 4   content_size_lemma     30000 non-null  int64
 5   doc_vector            30000 non-null  object
 6   sentance_content_bert  30000 non-null  object
dtypes: int64(2), object(5)
memory usage: 1.6+ MB
```

```
In [44]: new_df.to_csv('bert_data.csv')
```

chunking the data

```
In [45]: # Example function to chunk a document
def chunk_document(document, chunk_size, overlap):
    chunks = []
    start_idx = 0
    while start_idx < len(document):
        end_idx = min(start_idx + chunk_size, len(document))
        chunks.append(document[start_idx:end_idx])
        start_idx += chunk_size - overlap
    return chunks

# Parameters
chunk_size = 500 # Set the chunk size
overlap = 100    # Set the overlap size

# Apply chunking to content
new_df['content_chunks'] = new_df['cleaned_file_content'].apply(lambda x: chunk_document(x, chunk_size, overlap))
```

Final data set like below

In [46]: new_df.head()

Out[46]:

	num	name	cleaned_file_content	tokenised_content_lemma	content_size_lemma	doc_vector	sentence_content_bert	co
0	9487821	the.client.list.s02.e06.unanswered.prayers. (20...	riley ipreviously oni the client list derek yo...	['riley', 'ipreviously', 'oni', 'client', 'lis...	3244	[-1.09343290e-01 -7.68003017e-02 -1.80567741e-...	[-0.11649825, -0.044796105, 0.026207969, -0.10...	[r o
1	9183678	the.devil.is.a.woman.(1935).eng.1cd	support us and become vip member to remove all...	['support', 'u', 'become', 'vip', 'member', 'r...	2853	[-9.83889028e-02 -7.53047839e-02 -1.65519699e-...	[-0.017634636, -0.06649098, 0.008615398, -0.07...	[
2	9258855	kingdom.s04.e18.rivercrossing.battle. (2022).en...	script info title english us original script p...	['script', 'info', 'title', 'english', 'u', 'o...	1732	[-4.89264429e-02 -1.39095232e-01 -1.84526831e-...	[-0.017091548, -0.044068865, 0.04417335, -0.04...	c
3	9421238	sorority.babes.in.the.slimeball.bowlorama.2. (2...	wind whooshing eerie synth music bats chirping...	['wind', 'whooshing', 'eerie', 'synth', 'music...	2948	[-9.89964455e-02 -6.88564703e-02 -1.66541591e-...	[-0.06590253, -0.02547965, 0.048389476, -0.016...	[w
4	9518525	ada.twist.scientist.s03.e03.dadbotinto.the.bra...	watch any video online with opensubtitles free...	['watch', 'video', 'online', 'opensubtitles', ...	1482	[-0.10334199 -0.06256191 -0.15584423 -0.213279...	[-0.068130724, -0.024480091, -0.025265323, -0....	[w



In [47]: new_df.info()

```
<class 'pandas.core.frame.DataFrame'>
RangeIndex: 30000 entries, 0 to 29999
Data columns (total 8 columns):
#   Column                Non-Null Count  Dtype
---  -
0   num                    30000 non-null  int64
1   name                   30000 non-null  object
2   cleaned_file_content   30000 non-null  object
3   tokenised_content_lemma 30000 non-null  object
4   content_size_lemma     30000 non-null  int64
5   doc_vector             30000 non-null  object
6   sentence_content_bert   30000 non-null  object
7   content_chunks         30000 non-null  object
dtypes: int64(2), object(6)
memory usage: 1.8+ MB
```

In [48]: new_df.to_csv('final_data.csv')

Using the final data set

In [2]: `import pandas as pd`

In [3]: `final_df = pd.read_csv('final_data.csv')`

```
In [4]: final_df.head
```

```

Out[4]: <bound method NDFrame.head of          Unnamed: 0      num      name \
0          0  9487821  the.client.list.s02.e06.unanswered.prayers.(20...
1          1  9183678          the.devil.is.a.woman.(1935).eng.1cd
2          2  9258855  kingdom.s04.e18.rivercrossing.battle.(2022).en...
3          3  9421238  sorority.babes.in.the.slimeball.bowlorama.2.(2...
4          4  9518525  ada.twist.scientist.s03.e03.dadbotinto.the.bra...
...          ...          ...          ...
29995      29995  9473109  king.of.the.road.s03.e03.chained.up.and.trippe...
29996      29996  9487995          badoet.(2015).eng.1cd
29997      29997  9400464          velma.s01.e03.velma.kai.(2023).eng.1cd
29998      29998  9489287          numberblocks.s03.e26.thirteen.(2019).eng.1cd
29999      29999  9305337  def.comedy.jam.s02.e09.episode.2.9.(1992).eng.1cd

          cleaned_file_content \
0  riley i previously oni the client list derek yo...
1  support us and become vip member to remove all...
2  script info title english us original script p...
3  wind whooshing eerie synth music bats chirping...
4  watch any video online with opensubtitles free...
...          ...
29995  please do not attempt to perform any of these ...
29996  crickets chirring music box music apiopensubti...
29997  i previously oni velma the last place my moms ...
29998  one two three four five numberblocks six seven...
29999  kid capri its the russell simmons def comedy j...

          tokenised_content_lemma  content_size_lemma \
0  ['riley', 'i previously', 'oni', 'client', 'lis...  3244
1  ['support', 'u', 'become', 'vip', 'member', 'r...  2853
2  ['script', 'info', 'title', 'english', 'u', 'o...  1732
3  ['wind', 'whooshing', 'eerie', 'synth', 'music...  2948
4  ['watch', 'video', 'online', 'opensubtitles', ...  1482
...          ...
29995  ['please', 'attempt', 'perform', 'stunt', 'act...  3055
29996  ['cricket', 'chirring', 'music', 'box', 'music...  1805
29997  ['previously', 'oni', 'velma', 'last', 'place'...  1797
29998  ['one', 'two', 'three', 'four', 'five', 'numbe...  267
29999  ['kid', 'capri', 'russell', 'simmons', 'def', ...  1578

          doc_vector \
0  [-1.09343290e-01 -7.68003017e-02 -1.80567741e-...
1  [-9.83889028e-02 -7.53047839e-02 -1.65519699e-...
2  [-4.89264429e-02 -1.39095232e-01 -1.84526831e-...
3  [-9.89964455e-02 -6.88564703e-02 -1.66541591e-...
4  [-0.10334199 -0.06256191 -0.15584423 -0.213279...
...          ...

```

```

29995 [-1.06857479e-01 -7.70704076e-02 -1.95597991e-...
29996 [-0.06744428 -0.07895508 -0.16554289 -0.238525...
29997 [-0.09154191 -0.08225142 -0.17166042 -0.234753...
29998 [-9.71302763e-02 -2.97208223e-02 -1.45906299e-...
29999 [-0.08471369 -0.07977783 -0.18513532 -0.242066...

```

```

                                sentence_content_bert \
0      [-1.16498247e-01 -4.47961055e-02  2.62079686e-...
1      [-1.76346358e-02 -6.64909780e-02  8.61539785e-...
2      [-1.70915481e-02 -4.40688655e-02  4.41733487e-...
3      [-6.59025311e-02 -2.54796501e-02  4.83894758e-...
4      [-6.81307241e-02 -2.44800914e-02 -2.52653230e-...
...
29995 [ 7.56419310e-03 -5.54080568e-02  2.06312984e-...
29996 [ 1.17038190e-02 -2.66452506e-02  2.01791693e-...
29997 [-8.32793862e-02 -6.34732619e-02  5.84286172e-...
29998 [-8.33337158e-02 -1.16495080e-01 -2.15825532e-...
29999 [-4.05492522e-02 -7.91372731e-02  1.41583364e-...

```

```

                                content_chunks
0      ['riley ipreviously oni the client list derek ...
1      ['support us and become vip member to remove a...
2      ['script info title english us original script...
3      ['wind whooshing eerie synth music bats chirpi...
4      ['watch any video online with opensubtitles fr...
...
29995 ['please do not attempt to perform any of thes...
29996 ['crickets chirring music box music apiopensub...
29997 ['i previously oni velma the last place my mom...
29998 ['one two three four five numberblocks six sev...
29999 ['kid capri its the russell simmons def comedy...

```

```
[30000 rows x 9 columns]>
```

```
In [5]: final_df.info()
```

```
<class 'pandas.core.frame.DataFrame'>
RangeIndex: 30000 entries, 0 to 29999
Data columns (total 9 columns):
#   Column                                Non-Null Count  Dtype
---  -
0   Unnamed: 0                            30000 non-null  int64
1   num                                    30000 non-null  int64
2   name                                  30000 non-null  object
3   cleaned_file_content                  30000 non-null  object
4   tokenised_content_lemma               30000 non-null  object
5   content_size_lemma                   30000 non-null  int64
6   doc_vector                            30000 non-null  object
7   sentance_content_bert                 30000 non-null  object
8   content_chunks                        30000 non-null  object
dtypes: int64(3), object(6)
memory usage: 2.1+ MB
```

```
In [6]: final_df = final_df.rename(columns={'sentance_content_bert': 'embeddings'})
```

```
In [7]: final_df.head()
```

Out[7]:

	Unnamed: 0	num	name	cleaned_file_content	tokenised_content_lemma	content_size_lemma	doc_vector	embeddings
0	0	9487821	the.client.list.s02.e06.unanswered.prayers. (20...	riley ipreviously oni the client list derek yo...	['riley', 'ipreviously', 'oni', 'client', 'lis...	3244	[-1.09343290e-01 -7.68003017e-02 -1.80567741e-...	[-1.16498247e-01 -4.47961055e-02 2.62079686e-..
1	1	9183678	the.devil.is.a.woman.(1935).eng.1cd	support us and become vip member to remove all...	['support', 'u', 'become', 'vip', 'member', 'r...	2853	[-9.83889028e-02 -7.53047839e-02 -1.65519699e-...	[-1.76346358e-02 -6.64909780e-02 8.61539785e-..
2	2	9258855	kingdom.s04.e18.rivercrossing.battle. (2022).en...	script info title english us original script p...	['script', 'info', 'title', 'english', 'u', 'o...	1732	[-4.89264429e-02 -1.39095232e-01 -1.84526831e-...	[-1.70915481e-02 -4.40688655e-02 4.41733487e-..
3	3	9421238	sorority.babes.in.the.slimeball.bowlorama.2. (2...	wind whooshing eerie synth music bats chirping...	['wind', 'whooshing', 'eerie', 'synth', 'music...	2948	[-9.89964455e-02 -6.88564703e-02 -1.66541591e-...	[-6.59025311e-02 -2.54796501e-02 4.83894758e-..
4	4	9518525	ada.twist.scientist.s03.e03.dadbotinto.the.bra...	watch any video online with opensubtitles free...	['watch', 'video', 'online', 'opensubtitles', ...	1482	[-0.10334199 -0.06256191 -0.15584423 -0.213279...	[-6.81307241e-02 -2.44800914e-02 -2.52653230e-..

```
In [8]: # Assuming 'df' is your DataFrame
emb_df = final_df[['num', 'name', 'content_chunks', 'embeddings']].copy()
```

```
In [9]: # Assuming 'df' is your DataFrame
emb_df1 = final_df[['num', 'name', 'content_chunks', 'embeddings']].copy()
```

In [10]: `emb_df.info()`

```
<class 'pandas.core.frame.DataFrame'>
RangeIndex: 30000 entries, 0 to 29999
Data columns (total 4 columns):
#   Column          Non-Null Count  Dtype
---  -
0   num              30000 non-null  int64
1   name             30000 non-null  object
2   content_chunks   30000 non-null  object
3   embeddings        30000 non-null  object
dtypes: int64(1), object(3)
memory usage: 937.6+ KB
```

In [11]: `emb_df.head()`

Out[11]:

	num	name	content_chunks	embeddings
0	9487821	the.client.list.s02.e06.unanswered.prayers.(20...	['riley i previously oni the client list derek ...	[-1.16498247e-01 -4.47961055e-02 2.62079686e-...
1	9183678	the.devil.is.a.woman.(1935).eng.1cd	['support us and become vip member to remove a...	[-1.76346358e-02 -6.64909780e-02 8.61539785e-...
2	9258855	kingdom.s04.e18.rivercrossing.battle.(2022).en...	['script info title english us original script...	[-1.70915481e-02 -4.40688655e-02 4.41733487e-...
3	9421238	sorority.babes.in.the.slimeball.bowlorama.2.(2...	['wind whooshing eerie synth music bats chirpi...	[-6.59025311e-02 -2.54796501e-02 4.83894758e-...
4	9518525	ada.twist.scientist.s03.e03.dadbotinto.the.bra...	['watch any video online with opensubtitles fr...	[-6.81307241e-02 -2.44800914e-02 -2.52653230e-...

In [12]: `import sqlite3`

```
# Specify the path to your desired location (e.g., /content/drive/MyDrive)
db_path = 'chromadb.db'

# Connect to the SQLite database (creates the file if it doesn't exist)
conn = sqlite3.connect(db_path)
conn.close() # Close the connection

print(f"SQLite database '{db_path}' created successfully.")
```

SQLite database 'chromadb.db' created successfully.


```
In [13]: import sqlite3

try:
    # Specify the path to your desired location (e.g., /content/drive/MyDrive)
    db_path = r'R:\Data Science\ZDS\Internship\Innomatics\Assignment8_Search Engine Project\chromadb.db'

    # Connect to the SQLite database (creates the file if it doesn't exist)
    conn = sqlite3.connect(db_path)
    conn.close() # Close the connection

    print(f"SQLite database '{db_path}' created successfully.")
except Exception as e:
    print("An error occurred:", e)
```

SQLite database 'R:\Data Science\ZDS\Internship\Innomatics\Assignment8_Search Engine Project\chromadb.db' created successfully.

```
In [14]: import chromadb
chroma_client = chromadb.Client()
```

```
In [15]: collection = chroma_client.create_collection(name="my_collection")
```

```
In [16]: emb_df.head()
```

```
Out[16]:
```

	num	name	content_chunks	embeddings
0	9487821	the.client.list.s02.e06.unanswered.prayers.(20...	['riley ipreviously oni the client list derek ...	[-1.16498247e-01 -4.47961055e-02 2.62079686e-...
1	9183678	the.devil.is.a.woman.(1935).eng.1cd	['support us and become vip member to remove a...	[-1.76346358e-02 -6.64909780e-02 8.61539785e-...
2	9258855	kingdom.s04.e18.rivercrossing.battle.(2022).en...	['script info title english us original script...	[-1.70915481e-02 -4.40688655e-02 4.41733487e-...
3	9421238	sorority.babes.in.the.slimeball.bowlorama.2.(2...	['wind whooshing eerie synth music bats chirpi...	[-6.59025311e-02 -2.54796501e-02 4.83894758e-...
4	9518525	ada.twist.scientist.s03.e03.dadbotinto.the.bra...	['watch any video online with opensubtitles fr...	[-6.81307241e-02 -2.44800914e-02 -2.52653230e-...

Create a file in Cromadb

```
In [18]: import sqlite3
import pandas as pd
import json

# SQLite database path
db_path = 'chromadb1136.db'

# Connect to the SQLite database
conn = sqlite3.connect(db_path)

# Convert lists to strings
emb_df['content_chunks'] = emb_df['content_chunks'].apply(json.dumps)
emb_df['embeddings'] = emb_df['embeddings'].apply(json.dumps)

# Save DataFrame to SQLite
emb_df.to_sql("dd", conn, if_exists="replace", index=False)

print(f"DataFrame from {emb_df} is successfully saved to ChromaDB using SQLite: {db_path}")
```

	num	name \
0	9487821	the.client.list.s02.e06.unanswered.prayers.(20...
1	9183678	the.devil.is.a.woman.(1935).eng.1cd
2	9258855	kingdom.s04.e18.rivercrossing.battle.(2022).en...
3	9421238	sorority.babes.in.the.slimeball.bowlorama.2.(2...
4	9518525	ada.twist.scientist.s03.e03.dadbotinto.the.bra...
...
29995	9473109	king.of.the.road.s03.e03.chained.up.and.trippe...
29996	9487995	badoet.(2015).eng.1cd
29997	9400464	velma.s01.e03.velma.kai.(2023).eng.1cd
29998	9489287	numberblocks.s03.e26.thirteen.(2019).eng.1cd
29999	9305337	def.comedy.jam.s02.e09.episode.2.9.(1992).eng.1cd

	content_chunks \
0	"\"['riley ipreviously oni the client list der...
1	"\"['support us and become vip member to remov...
2	"\"['script info title english us original scr...
3	"\"['wind whooshing eerie synth music bats chi...
4	"\"['watch any video online with opensubtitles...
...	...
29995	"\"['please do not attempt to perform any of t...
29996	"\"['crickets chirring music box music apiopen...
29997	"\"['i previously oni velma the last place my ...
29998	"\"['one two three four five numberblocks six ...
29999	"\"['kid capri its the russell simmons def com...

	embeddings
0	"[-1.16498247e-01 -4.47961055e-02 2.62079686e...
1	"[-1.76346358e-02 -6.64909780e-02 8.61539785e...
2	"[-1.70915481e-02 -4.40688655e-02 4.41733487e...
3	"[-6.59025311e-02 -2.54796501e-02 4.83894758e...
4	"[-6.81307241e-02 -2.44800914e-02 -2.52653230e...
...	...
29995	"[7.56419310e-03 -5.54080568e-02 2.06312984e...
29996	"[1.17038190e-02 -2.66452506e-02 2.01791693e...
29997	"[-8.32793862e-02 -6.34732619e-02 5.84286172e...
29998	"[-8.33337158e-02 -1.16495080e-01 -2.15825532e...
29999	"[-4.05492522e-02 -7.91372731e-02 1.41583364e...

[30000 rows x 4 columns] is successfully saved to ChromaDB using SQLite: chromadb1136.db

In []:

In []:

Experiment on 5k data set

In [34]: `new_df.head()`

Out[34]:

	num	name	cleaned_file_content	tokenised_content_lemma	con
0	9487821	the.client.list.s02.e06.unanswered.prayers. (20...	riley ipreviously oni the client list derek yo...	['riley', 'ipreviously', 'oni', 'client', 'lis...	
1	9183678	the.devil.is.a.woman.(1935).eng.1cd	support us and become vip member to remove all...	['support', 'u', 'become', 'vip', 'member', 'r...	
2	9258855	kingdom.s04.e18.rivercrossing.battle. (2022).en...	script info title english us original script p...	['script', 'info', 'title', 'english', 'u', 'o...	
3	9421238	sorority.babes.in.the.slimeball.bowlorama.2. (2...	wind whooshing eerie synth music bats chirping...	['wind', 'whooshing', 'eerie', 'synth', 'music...	
4	9518525	ada.twist.scientist.s03.e03.dadbotinto.the.bra...	watch any video online with opensubtitles free...	['watch', 'video', 'online', 'opensubtitles', ...	

In [35]: `final_df = new_df.copy()`

In [52]: `df1 = new_df.copy()`

Create a file in Cromadb

In [38]: `import sqlite3`

```
# Specify the path to your desired location (e.g., /content/drive/MyDrive)
db_path = 'chromadb.db'

# Connect to the SQLite database (creates the file if it doesn't exist)
conn = sqlite3.connect(db_path)
conn.close() # Close the connection

print(f"SQLite database '{db_path}' created successfully.")
```

SQLite database 'chromadb.db' created successfully.

```
In [39]: import sqlite3

try:
    # Specify the path to your desired location (e.g., /content/drive/MyDrive)
    db_path = r'R:\Data Science\ZDS\Internship\Innomatics\Assignment8_Search Engine Pro

    # Connect to the SQLite database (creates the file if it doesn't exist)
    conn = sqlite3.connect(db_path)
    conn.close() # Close the connection

    print(f"SQLite database '{db_path}' created successfully.")
except Exception as e:
    print("An error occurred:", e)
```

SQLite database 'R:\Data Science\ZDS\Internship\Innomatics\Assignment8_Search Engine Project\chromadb.db' created successfully.

```
In [40]: import chromadb
chroma_client = chromadb.Client()
```

```
In [41]: collection = chroma_client.create_collection(name="my_collection")
```

```
In [53]: df1 = df1.rename(columns={'sentance_content_bert': 'embeddings'})
```

```
In [69]: emb_df1 = df1.copy()
```

```
In [54]: # Assuming 'df' is your DataFrame
emb_df1 = df1[['num', 'name', 'content_chunks', 'embeddings']].copy()
```

```
In [56]: emb_df1.head(2)
```

```
Out[56]:
```

	num	name	content_chunks	embeddings
0	9487821	the.client.list.s02.e06.unanswered.prayers. (20...	[riley ipreviously oni the client list derek y...	[-0.022161696, -0.111041605, 0.02133926, 0.050...
1	9183678	the.devil.is.a.woman.(1935).eng.1cd	[support us and become vip member to remove al...	[-0.0005046579, -0.060602482, -0.03618558, -0....

```
In [ ]:
```

```
In [ ]: collection = chroma_client.create_collection(name="my_collection")
```

```
In [ ]: df1 = df1.rename(columns={'sentance_content_bert': 'embeddings'})
```

```
In [66]: import sqlite3
import pandas as pd
import json

# SQLite database path
db_path = 'chromadb1133.db'

# Connect to the SQLite database
conn = sqlite3.connect(db_path)

# Convert lists to strings
emb_df1['content_chunks'] = emb_df1['content_chunks'].apply(json.dumps)
emb_df1['embeddings'] = emb_df1['embeddings'].apply(lambda x: json.dumps(x.tolist()))

# Save DataFrame to SQLite
emb_df1.to_sql("dd", conn, if_exists="replace", index=False)

print(f"DataFrame from {emb_df1} is successfully saved to ChromaDB using SQLite: {db_p
```

```
DataFrame from          num                                name \
0      9487821  the.client.list.s02.e06.unanswered.prayers.(20...
1      9183678                the.devil.is.a.woman.(1935).eng.1cd
2      9258855  kingdom.s04.e18.rivercrossing.battle.(2022).en...
3      9421238  sorority.babes.in.the.slimeball.bowlorama.2.(2...
4      9518525  ada.twist.scientist.s03.e03.dadbotinto.the.bra...
...      ...
4996  9237926                blue.gender.s01.e16.a.sign.(2000).eng.1cd
4997  9375637  the.equalizer.s01.e10.reckoning.(2021).eng.1cd
4998  9226928  rap.sht.s01.e08.something.for.the.road.(2022)...
4999  9316079  dead.to.me.s03.e08.well.find.a.way.(2022).eng.1cd
5000  9305324  def.comedy.jam.s01.e04.episode.1.4.(1992).eng.1cd

                                content_chunks \
0      "\"\\\"\\\"[\\\"\\\"\\\"\\\"riley ipreviously oni the clie...
1      "\"\\\"\\\"[\\\"\\\"\\\"\\\"support us and become vip memb...
2      "\"\\\"\\\"[\\\"\\\"\\\"\\\"script info title english us o...
3      "\"\\\"\\\"[\\\"\\\"\\\"\\\"wind whooshing eerie synth mus...
4      "\"\\\"\\\"[\\\"\\\"\\\"\\\"watch any video online with op...
...      ...
4996  "\"\\\"\\\"[\\\"\\\"\\\"\\\"watch any video online with op...
4997  "\"\\\"\\\"[\\\"\\\"\\\"\\\"iim the one you call when you ...
4998  "\"\\\"\\\"[\\\"\\\"\\\"\\\"rap music playing wondering ho...
4999  "\"\\\"\\\"[\\\"\\\"\\\"\\\"wwwopensubtitlescom hey mama h...
5000  "\"\\\"\\\"[\\\"\\\"\\\"\\\"theme music playing kid capri ...

                                embeddings
0      "[-0.02216169610619545, -0.1110416054725647, 0...
1      "[-0.0005046579171903431, -0.06060248240828514...
2      "[-0.03646128252148628, -0.030952945351600647,...
3      "[0.02794903703033924, -0.07734321057796478, 0...
4      "[0.0011420863447710872, -0.11879019439220428,...
...      ...
4996  "[0.007031263317912817, -0.10100991278886795, ...
4997  "[0.006331723649054766, -0.08844021707773209, ...
4998  "[0.010032282210886478, -0.14277908205986023, ...
4999  "[-0.04041879624128342, -0.07212220132350922, ...
5000  "[0.014225239865481853, -0.08701223880052567, ...

[5001 rows x 4 columns] is successfully saved to ChromaDB using SQLite: chromadb113
3.db
```

In []:

In [68]: `emb_df1.content_chunks[66]`

```
witn wasnington but it does not look promising nmm whats tne big wnoop we can keep
the talismans here in one of uncles magic boxes the risk is too great the talisman
s will be safest from dark forces on sacred ground we must take them to the benshu
i temple but but isnt that like near china jade the talismans are unpredictable da
ngerous you should not be playing with \\\\", \\\\", that like
near china jade the talismans are unpredictable dangerous you should not be playin
g with them in the first place blows raspberry snip testing im in im gonna miss yo
u guys how about one last hurrah for old times sake hmm hydraulic humming uhoh huh
aah ooh groaning man section s a long way from washington captain black black i pr
omise your trip will be worthwhile sir gasps black follow me once you see just how
powerful these talismans are im sure youll agree its imperative we keep them
\\\\\", \\\\", once you see just how powerful these talismans ar
e im sure youll agree its imperative we keep them here under lock and laser whew t
hose guys must be blind observe the power of invisibility astonishing isnt it i co
uld be over here or perhaps im over here hmm must not be doing it right never mind
lets try levitation grunting uhh i dont understand its running a section can be ve
ry stressful captain black im authorizing weeks of r and r effective immediately i
suggest you get as far away from here \\\\", \\\\",black im aut
horizing weeks of r and r effective immediately i suggest you get as far away from
here as possible im very sorry about this you should be sorry we are not used to t
```

In [63]: `import sqlite3`

```
# Connect to the SQLite database
conn = sqlite3.connect('chromadb1133.db')
cursor = conn.cursor()

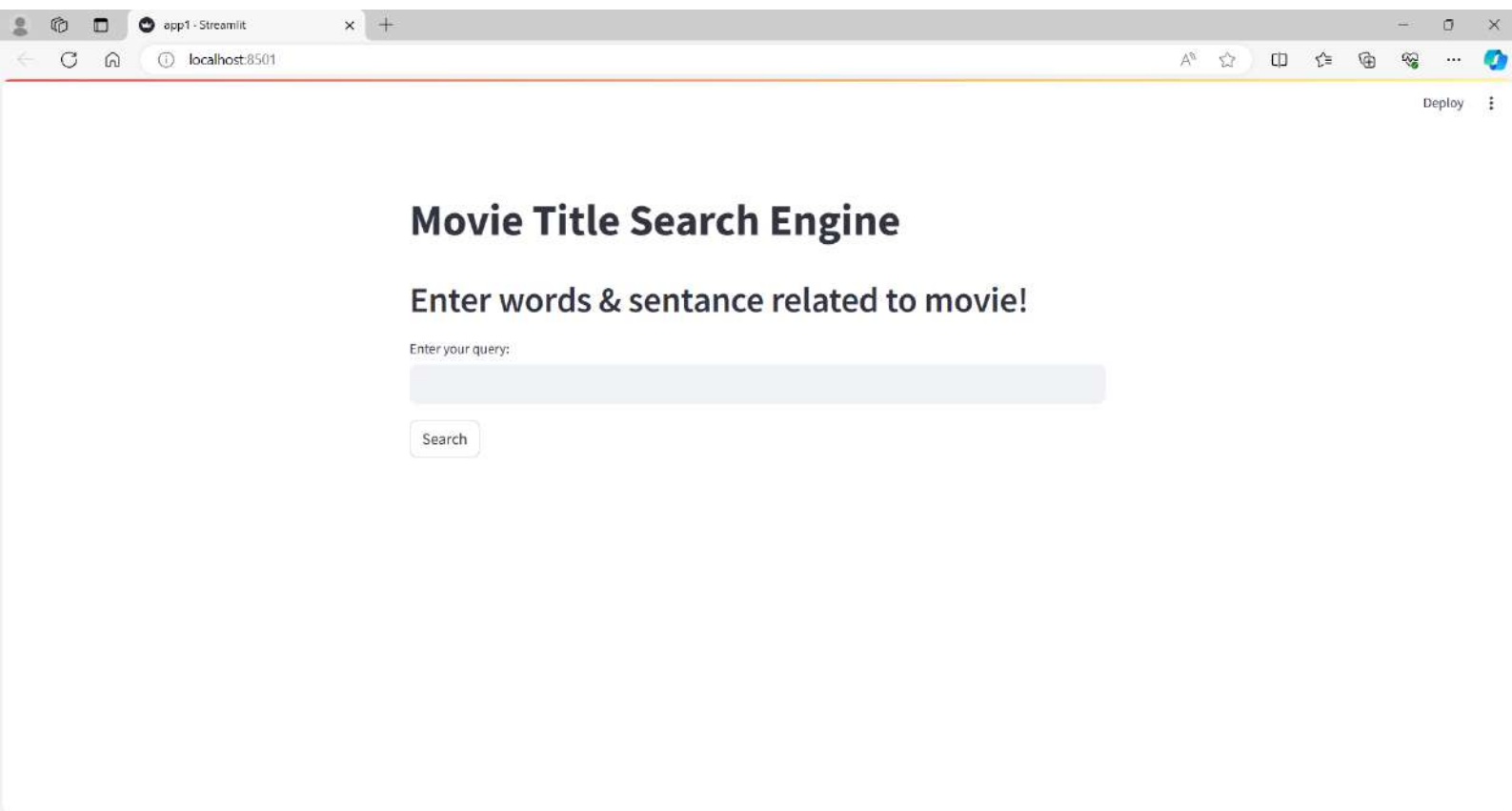
# Fetch the table schema
cursor.execute("PRAGMA table_info(dd)")
columns = cursor.fetchall()

# Print the table schema
for column in columns:
    print(column)

# Close the cursor and connection
cursor.close()
conn.close()
```

```
(0, 'num', 'INTEGER', 0, None, 0)
(1, 'name', 'TEXT', 0, None, 0)
(2, 'content_chunks', 'TEXT', 0, None, 0)
(3, 'embeddings', 'TEXT', 0, None, 0)
```

In []:



Movie Title Search Engine

Enter words & sentence related to movie!

Enter your query:

Search

Movie Title Search Engine

Enter words & sentences related to movies!

Enter your query:

watch

Search

Top 10 Unique Names:

1. rollerball.(1975).eng.1cd
2. americas.got.talent.s17.e20.qualifiers.5.results.(2022).eng.1cd
3. unforgettable.s04.e10.game.on.(2016).eng.1cd
4. detective.knight.rogue.(2022).eng.1cd
5. jared.from.subway.catching.a.monster.s01.e03.part.3.(2023).eng.1cd
6. full.swing.s01.e02.win.or.go.home.(2023).eng.1cd
7. star.trek.lower.decks.s03.e01.grounded.(2022).eng.1cd

Movie Title Search Engine

Enter words & sentences related to movies!

Enter your query:

criminal conspiracy guess

Search

Top 10 Unique Names:

1. shetland.s07.e06.episode.7.6.(2022).eng.1cd
2. rabbit.hole.s01.e03.the.algorithms.of.control.(2023).eng.1cd
3. big.sky.s03.e10.a.thin.layer.of.rock.(2022).eng.1cd
4. the.gold.s01.e06.ill.be.remembered.(2023).eng.1cd
5. trafficked.with.mariana.van.zeller.s01.e06.cocaine.(2020).eng.1cd
6. harina.s01.e07.episode.1.7.(2022).eng.1cd
7. the.mentalist.s06.e15.white.as.the.driven.snow.(2014).eng.1cd
8. 24.s06.e14.day.6.700.p.m.800.p.m.(2007).eng.1cd
9. project.wolf.hunting.(2022).eng.1cd

Movie Title Search Engine

Enter words & sentences related to movies!

Enter your query:

music playing kid capri its the

Search

Top 10 Unique Names:

1. def.comedy.jam.s06.e11.episode.6.11.(1997).eng.1cd
2. def.comedy.jam.s06.e07.episode.6.7.(1997).eng.1cd
3. def.comedy.jam.s01.e06.episode.1.6.(1992).eng.1cd
4. lady.voyeur.s01.e07.ela.esta.entre.nos.(2023).eng.1cd
5. def.comedy.jam.s01.e04.episode.1.4.(1992).eng.1cd
6. koala.man.s01.e05.ode.to.a.koala.bear.(2023).eng.1cd
7. what.we.do.in.the.shadows.s04.e10.sunrise.sunset.(2022).eng.1cd
8. queer.eye.s06.e05.crawzaddy.(2021).eng.1cd

[illegible]

Movie Title Search Engine

Enter words & sentences related to movies!

Enter your query:

it going to be no peeking laughing gasping its a boy jessica wait theres something else in there gasping i

Search

Top 10 Unique Names:

1. great.news.s02.e05.night.of.the.living.screen.(2017).eng.1cd
2. the.xfiles.s02.e21.the.calusari.(1995).eng.1cd
3. accused.s01.e07.brendas.story.(.eng.1cd
4. friends.s01.e24.the.one.where.rachel.finds.out.(1995).eng.1cd
5. fantasy.island.s02.e02.hurricane.helenethe.bachelor.party.(2023).eng.1cd
6. beverly.hills.90210.s07.e27.i.only.have.eyes.for.you.(1997).eng.1cd
7. the.nanny.s02.e19.a.fine.friendship.(1995).eng.1cd
8. nevsu.s02.e03.nitzas.diet.(2020).eng.1cd
9. being.human.s03.e03.the.teens.they.are.a.changin.(2013).eng.1cd

Movie Title Search Engine

Enter words & sentences related to movies!

Enter your query:

that like near china jade the talismans are unpredictable dangerous you should not be playing with ther

Search

Top 10 Unique Names:

1. rush.s03.e22.episode.3.22.(2010).eng.1cd
2. the.rookie.s05.e06.the.reckoning.(2022).eng.1cd
3. detectorists.s01.e01.episode.1.(2014).eng.1cd
4. the.middle.s01.e12.the.neighbor.(2010).eng.1cd
5. hai.shi.shen.lou.(1987).eng.1cd
6. detectorists.s03.e02.episode.3.2.(2017).eng.1cd
7. my.dad.the.bounty.hunter.s01.e05.chillion5.(2023).eng.1cd
8. jackie.chan.adventures.s02.e39.the.amazing.girl.(2002).eng.1cd

EXPLORER

GENAI_LLMS_APPS

> .env

> bug_fix_2

> chatbot_3

> genai_app_1

> search_enginee

data

chromadb1.db

chromadb666.db

chromadb1111.db

chromadb1122.db

app.py

app1.py

app2.py

app3.py

app1122.py

app1133.py

chroma.sqlite3

chromadb666.db

chat_key.pem

scp

OUTLINE

TIMELINE

app1122.py

search_enginee > app1122.py > main

```
1 import streamlit as st
2 import sqlite3
3 import numpy as np
4 from sklearn.metrics.pairwise import cosine_similarity
5 from sentence_transformers import SentenceTransformer
6
7 # Connect to the SQLite database
8 def connect_db(database_path):
9     conn = sqlite3.connect("data/chromadb1122.db")
10    return conn
11
12 def get_top_10_unique_names(query, conn, model):
13     c = conn.cursor()
14     query_embedding = model.encode([query])[0] # Use [0] to get the single embedding from the list
15     similarities = []
16     c.execute("SELECT * FROM dd")
17     for row in c.fetchall():
18         dd_num, dd_name, dd_content_chunks, dd_embeddings = row # Assuming the fourth column is not needed
19         try:
20             embeddings = np.fromstring(dd_embeddings[1:-1], sep=', ') # Parse as numpy array
21             # Check if embeddings is a list of valid numbers
22             if embeddings.size == 0: # Skip empty embeddings
23                 continue
24             similarity = cosine_similarity(query_embedding.reshape(1, -1), embeddings.reshape(1, -1))[0][0]
25             similarities.append((dd_name, similarity))
26         except Exception as e:
27             print(f"Error processing embeddings for {dd_name}: {e}")
28             continue
29     c.close()
30
31     sorted_names = [name for name, _ in sorted(similarities, key=lambda x: x[1], reverse=True)]
32     unique_names = []
33     for name in sorted_names:
34         if name not in unique_names:
35             unique_names.append(name)
36             if len(unique_names) == 10:
37                 break
```

EXPLORER

GENAI_LLMS_APPS

> .env

> bug_fix_2

> chatbot_3

> genai_app_1

search_enginee

data

chromadb1.db

chromadb666.db

chromadb1111.db

chromadb1122.db

app.py

app1.py

app2.py

app3.py

app1122.py

app1133.py

chroma.sqlite3

chromadb666.db

chat_key.pem

scp

OUTLINE

TIMELINE

app1122.py

search_enginee > app1122.py > main

```
12 def get_top_10_unique_names(query, conn, model):
31     sorted_names = [name for name, _ in sorted(similarities, key=lambda x: x[1], reverse=True)]
32     unique_names = []
33     for name in sorted_names:
34         if name not in unique_names:
35             unique_names.append(name)
36             if len(unique_names) == 10:
37                 break
38     return unique_names
40
41 # Main function to run the Streamlit app
42 def main():
43     st.title('Movie Title Search Engine')
44     st.header('Enter words & sentences related to movies!')
45     query = st.text_input('Enter your query:')
46     database_path = 'data/chromadb1122.db'
47     model = SentenceTransformer('all-MiniLM-L6-v2')
48
49     if st.button('Search'):
50         if query:
51             try:
52                 conn = connect_db(database_path)
53                 top_10_unique_names = get_top_10_unique_names(query, conn, model)
54                 if not top_10_unique_names:
55                     st.write('No matches found.')
56                 else:
57                     st.write('Top 10 Unique Names:')
58                     for i, name in enumerate(top_10_unique_names, start=1):
59                         st.write(f"{i}. {name}")
60             except Exception as e:
61                 st.error(f"An error occurred: {e}")
62             finally:
63                 conn.close()
64
65 if __name__ == '__main__':
66     main()
```

Ln 50, Col 18 Spaces: 4 UTF-8 CRLF Python 3.11.3 (.env/.venv)