# DATA SCIENCE REPORT

WAJAHAT MASOOD      SP19-BSE-003
ABDULLAH RAUF       SP19-BSE-007
SAIF UR REHMAN      SP19-BSE-013

**Assignment :**
**Text Data Preprocessing & Analysis**

**Group 1**

# Workflow

## Data Extraction

Data is been entreated from twitter about the president of Pakistan.

## Cleaning of data

Manual cleaning of bogus data

## Tokenization of the data

By using NLTK python library we tokenize our data.

## Stop Word Removal

By Using GENSIM.PARSING & REMOVE_STOPWORDS we remove the stop words to make analysis better

## Stemming

By Using NLTK STEMMER & PORTER STEMMER we Do STEMMING

# Workflow Cont...

### Lemmatization

By using NLTK STEMMER & WORD NET LEMATIZER we do Lemmatization.

### WORD CLOUD

By Using WORD CLOUD LIB, PIL IMAGE, NUMPY, URL LIB, MATPLOT we draw a word cloud

### VSM

By using COSIN_SIMILARITY LIB from SK learn. Metrices

### N Gram

By Using NKTK we import N gram we take UNI gram, bi gram, Tri Gram up to n gram
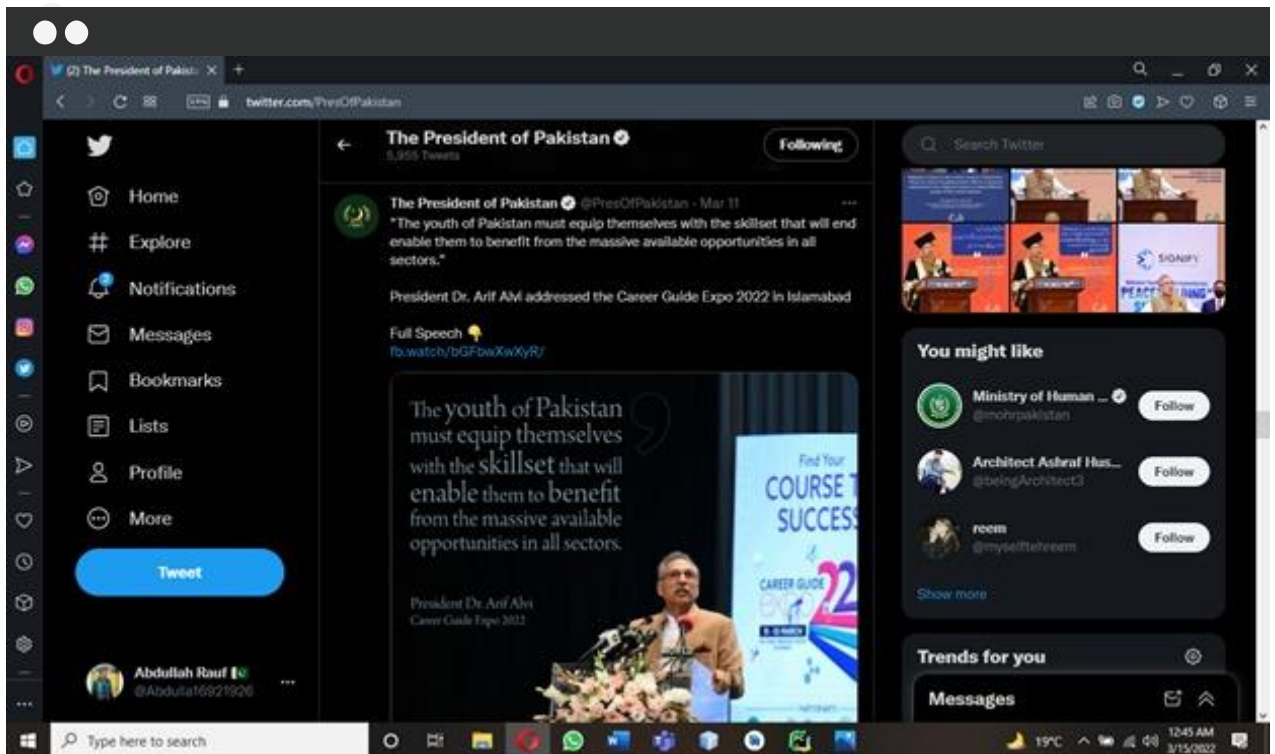
# TOOLS WE USE

**PyCharm**

1



**Google COLAB**

2



# DETAIL OF DATA

We gather Data from twitter (President of Pakistan). We Extract the data and clean it first and than we start our work and its is mentioned in workflow that how we did our wok.

# Screen Shot # 00



# DATA SET

---

We gather Data from twitter (President of Pakistan). We Extract the data and clean it first and than we start our work and its is mentioned in workflow that how we did our wok.

# Screen Shot # 01



# TOKENIZATION

Tokenization is the process of replacing sensitive data with unique identification symbols that retain all the essential information about the data without compromising its security.

# Screen Shot # 02



# STOP WORD REMOVAL

Stop word removal is one of the most used preprocessing steps across different NLP applications. The idea is simply removing the words that occur commonly across all the documents in the corpus. Typically, articles and pronouns are generally classified as stop words.

# Screen Shot # 03



# STEMMING

In processing unstructured text, stemming is the process of converting multiple forms of the same word into one stem, to simplify the task of analyzing the processed text. For example, in the previous sentence, "processing," "process," and "processed" would all be converted to the single stem "process."
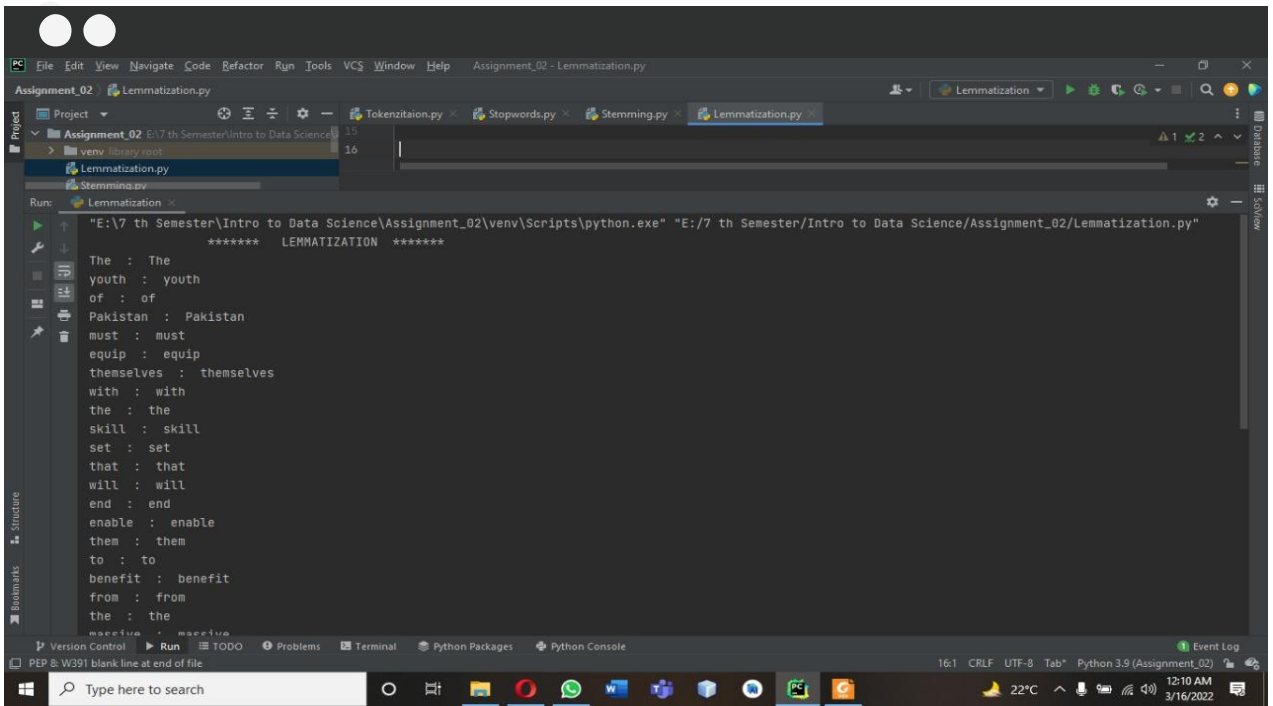
# Screen Shot # 04



# STEMMING RESULT

In processing unstructured text, stemming is the process of converting multiple forms of the same word into one stem, to simplify the task of analyzing the processed text. For example, in the previous sentence, "processing," "process," and "processed" would all be converted to the single stem "process."

# Screen Shot # 05



```python
import nltk
from nltk.stem import WordNetLemmatizer
wordnet_lemmatizer = WordNetLemmatizer()

text = """The youth of Pakistan must equip themselves with the skill set that will end enable them to benefit from
massive available opportunities in all sectors.
President Dr. Arif Alvi addressed the Career Guide Expo 2022 in Islamabad"""

tokenization = nltk.word_tokenize(text)

print("          *******  LEMMATIZATION  *******          ")

for w in tokenization:
    print("{}  :  {}".format(w, wordnet_lemmatizer.lemmatize(w)))
```

# LEMMETIZER

Lemmatization is the method to normalize the text documents. The main goal of the text normalization is to keep the vocabulary small and remove the noise(unwanted stuff) which helps to improve the accuracy of many language modeling tasks.
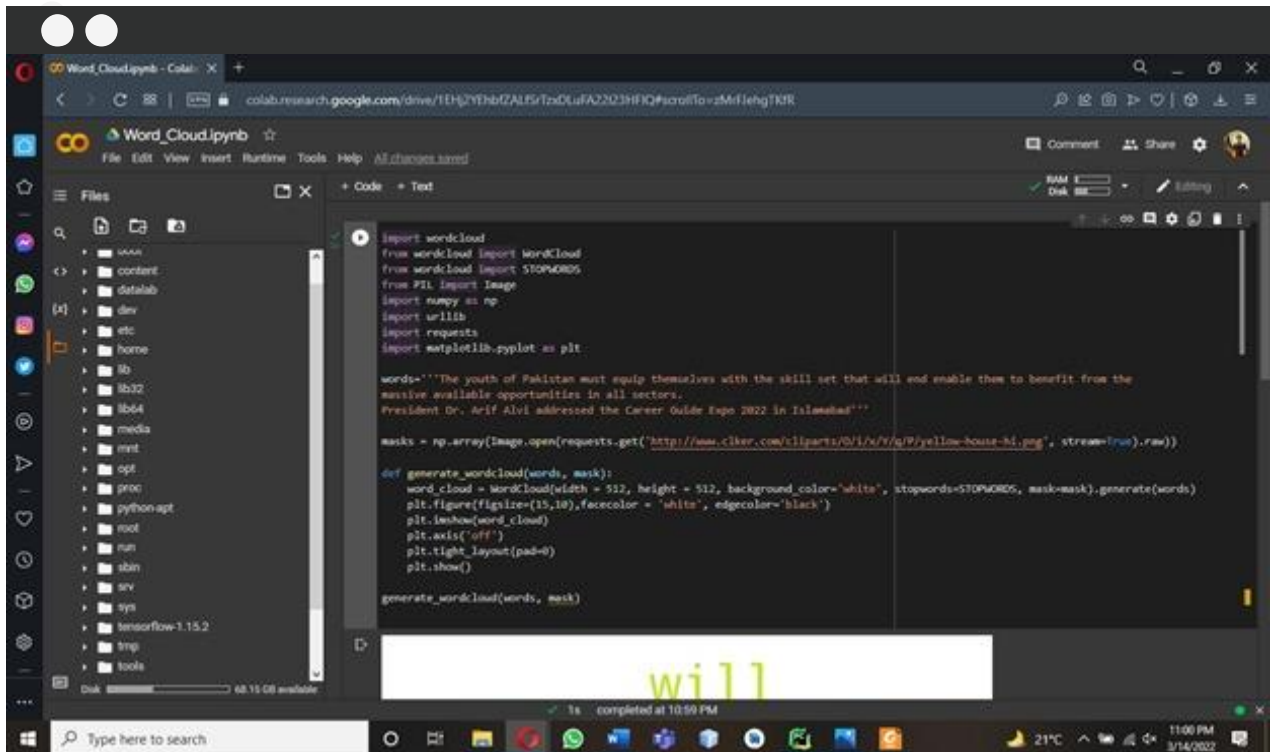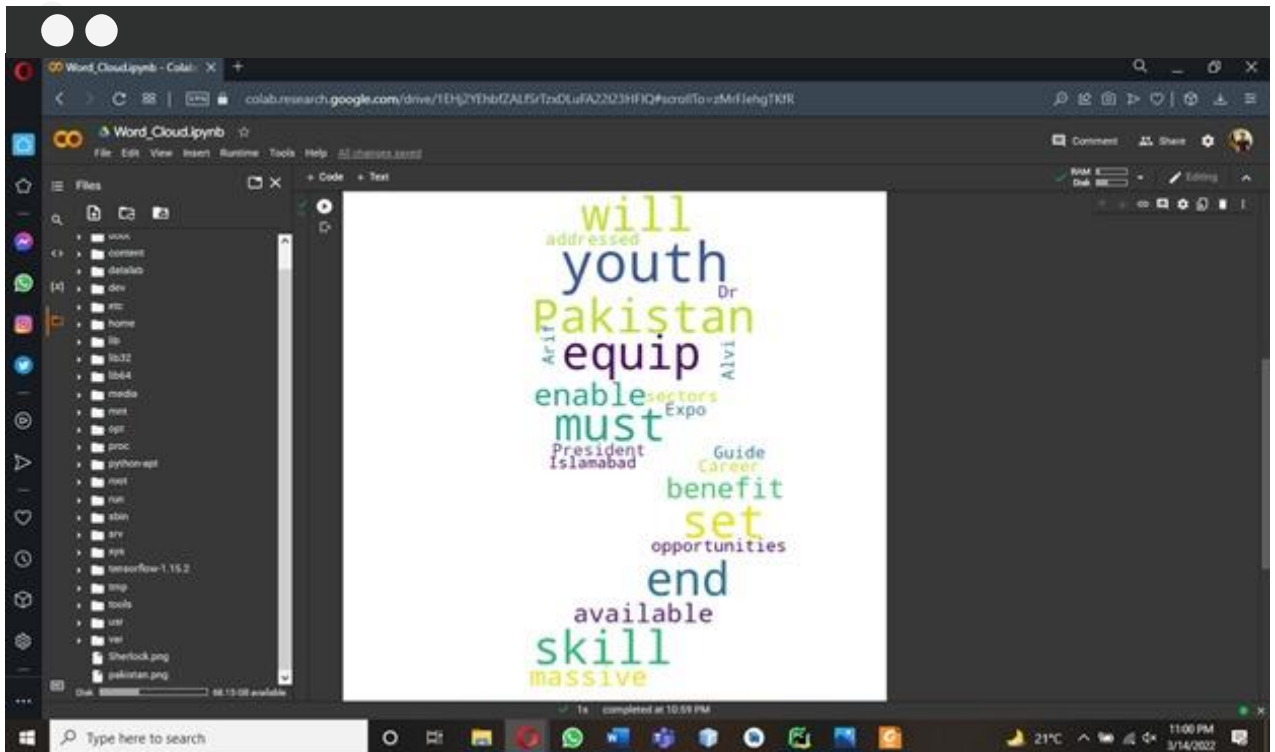
# Screen Shot # 06



# LEMMETIZER OUTPUT

Lemmatization is the method to normalize the text documents. The main goal of the text normalization is to keep the vocabulary small and remove the noise(unwanted stuff) which helps to improve the accuracy of many language modeling tasks.
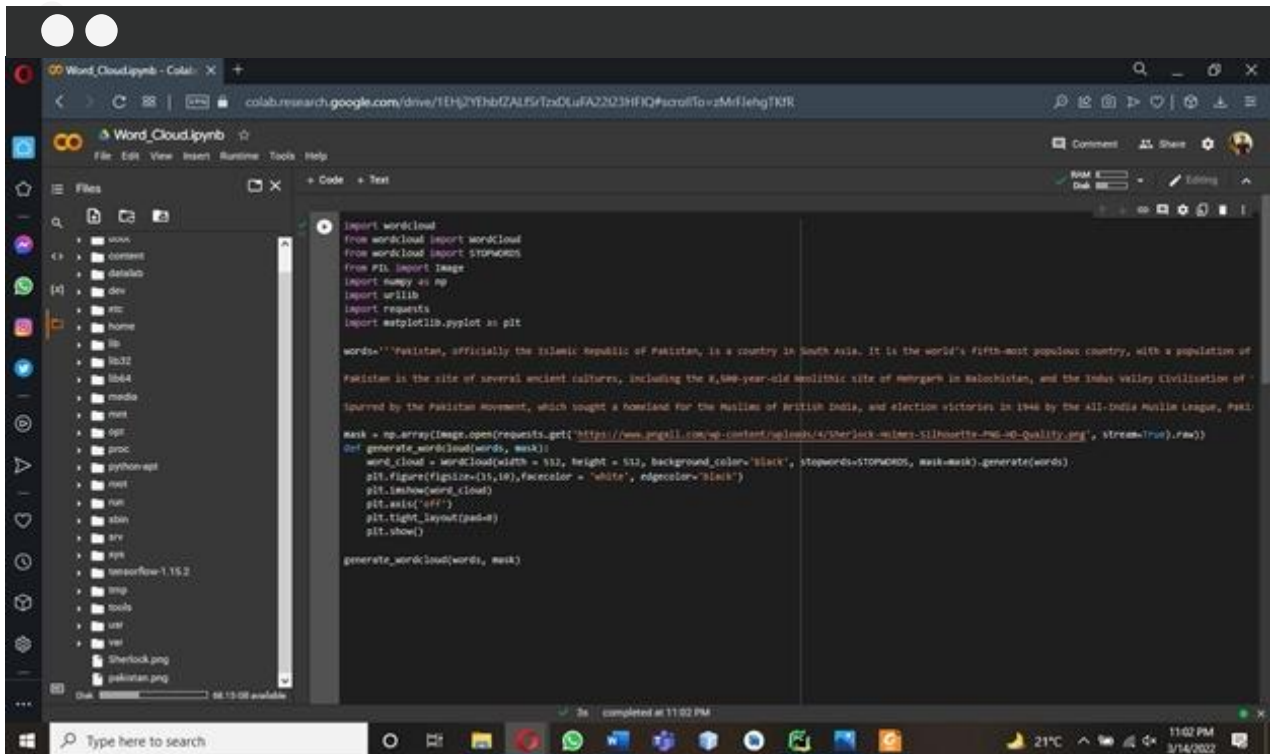
# Screen Shot # 07



# WORD CLOUD OF DATA SET

Word clouds or tag clouds are graphical representations of word frequency that give greater prominence to words that appear more frequently in a source text.
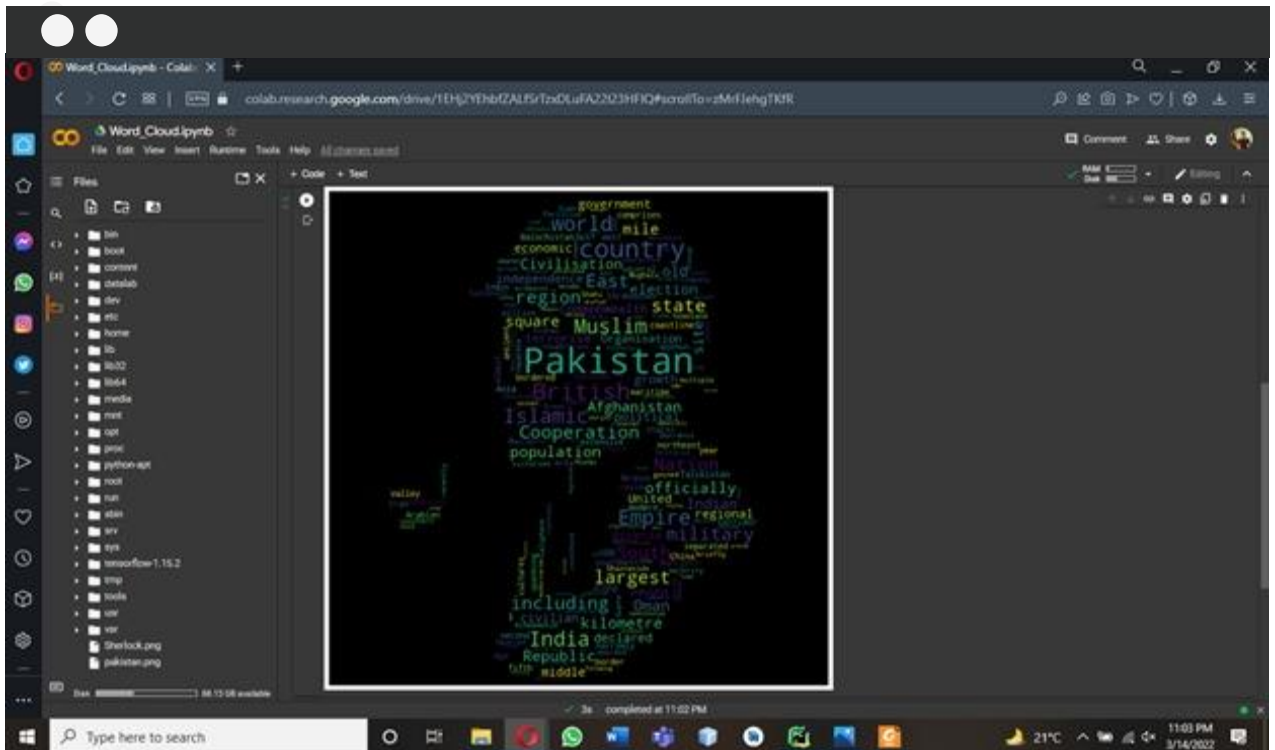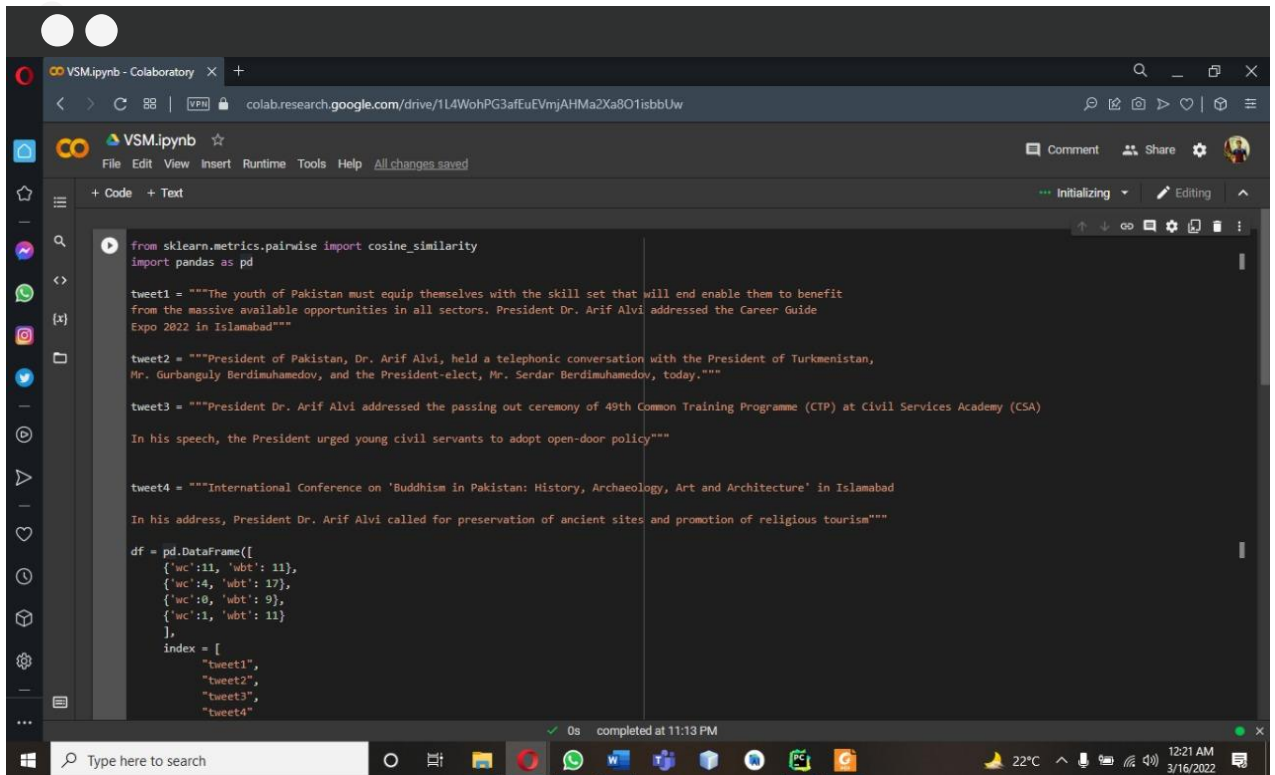
# Screen Shot # 08



# WORD CLOUD OF DATA SET OUTPUT

Word clouds or tag clouds are graphical representations of word frequency that give greater prominence to words that appear more frequently in a source text.

# Screen Shot # 09



# WORD CLOUD OF DATA SET

Word clouds or tag clouds are graphical representations of word frequency that give greater prominence to words that appear more frequently in a source text. DATA FROM WIKIPEDIA

# Screen Shot # 10



# WORD CLOUD OF DATA SET OUTPUT

Word clouds or tag clouds are graphical representations of word frequency that give greater prominence to words that appear more frequently in a source text. TWITTER DATA SET

# Screen Shot # 11



# VECTOR SPACE MODELING

Vector space models are to consider the relationship between data that are represented by vectors. It is popular in information retrieval systems but also useful for other purposes. Generally, this allows us to compare the similarity of two vectors from a geometric perspective.

# Screen Shot # 12



# VECTOR SPACE MODELING

Vector space models are to consider the relationship between data that are represented by vectors. It is popular in information retrieval systems but also useful for other purposes. Generally, this allows us to compare the similarity of two vectors from a geometric perspective.

# Screen Shot # 13



# VECTOR SPACE MODELING OUTPUT

Vector space models are to consider the relationship between data that are represented by vectors. It is popular in information retrieval systems but also useful for other purposes. Generally, this allows us to compare the similarity of two vectors from a geometric perspective.
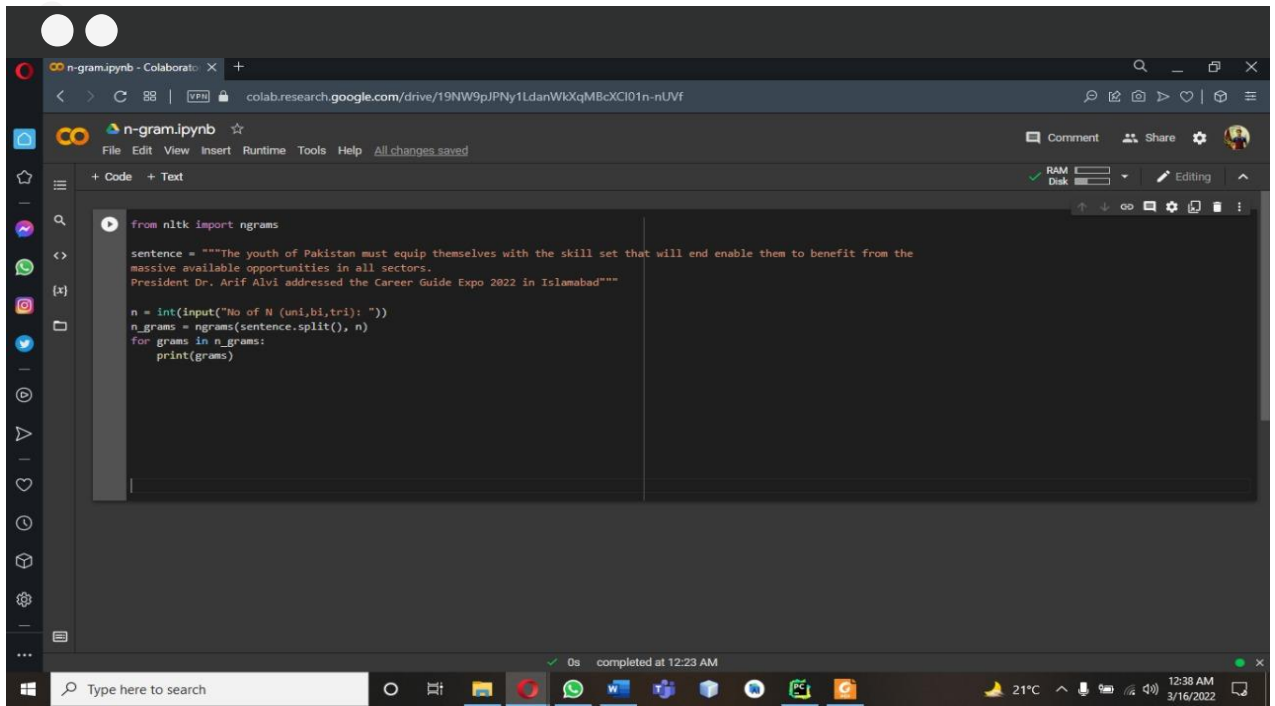
# Screen Shot # 14



# VECTOR SPACE MODELING OUTPUT

Vector space models are to consider the relationship between data that are represented by vectors. It is popular in information retrieval systems but also useful for other purposes. Generally, this allows us to compare the similarity of two vectors from a geometric perspective.

# Screen Shot # 15



# VSM DATA SET

Vector space models are to consider the relationship between data that are represented by vectors. It is popular in information retrieval systems but also useful for other purposes. Generally, this allows us to compare the similarity of two vectors from a geometric perspective.

# Screen Shot # 16



# N-GRAM

———

N-grams are continuous sequences of words or symbols or tokens in a document. In technical terms, they can be defined as the neighbouring sequences of items in a document. They come into play when we deal with text data in NLP(Natural Language Processing) tasks.

# Screen Shot # 17



# UNI-GRAM

---

unigram means taking only one word at a time, bigram means taking two words at a time and trigram means taking three words at a time.
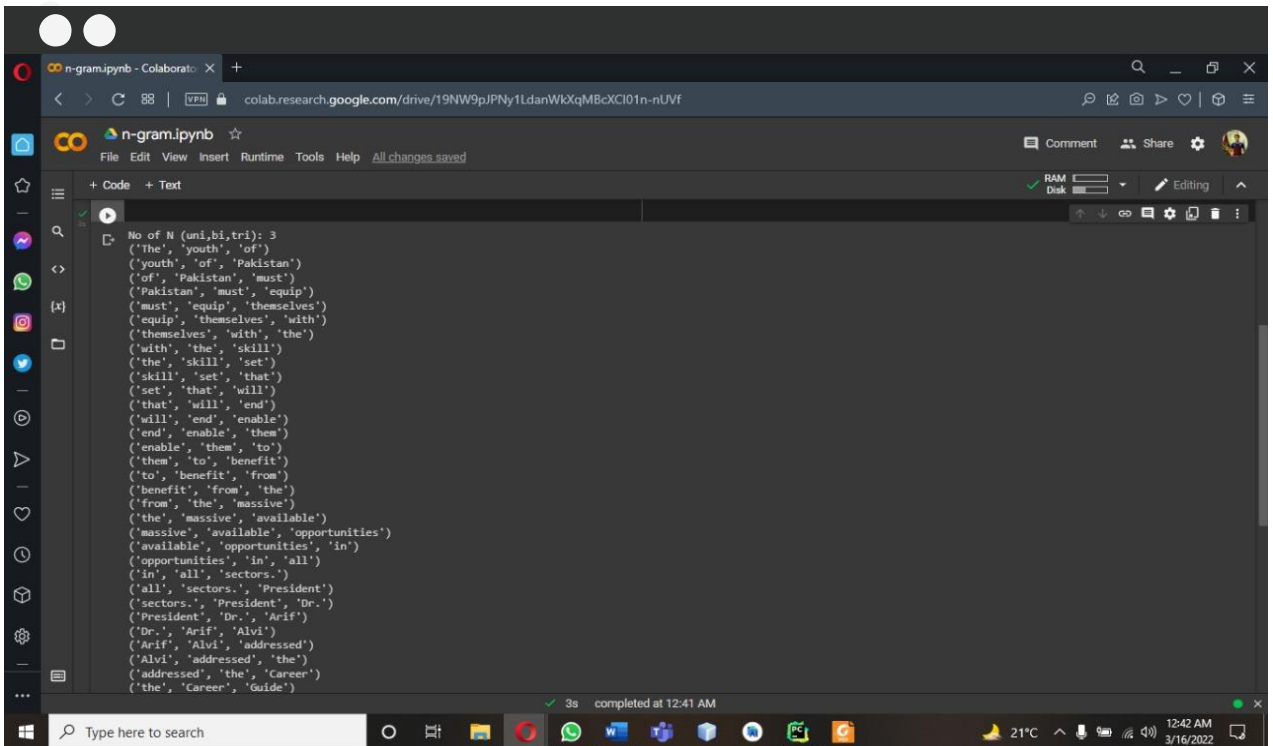
# Screen Shot # 18



# BI-GRAM

---

N-grams are continuous sequences of words or symbols or tokens in a document. In technical terms, they can be defined as the neighboring sequences of items in a document. They come into play when we deal with text data in NLP(Natural Language Processing) tasks.

# Screen Shot # 20



# TRI-GRAM

---

N-grams are continuous sequences of words or symbols or tokens in a document. In technical terms, they can be defined as the neighboring sequences of items in a document. They come into play when we deal with text data in NLP(Natural Language Processing) tasks.

# THE END

## REFERENCES

https://www.pngall.com/wp-content/uploads/4/Sherlock-Holmes-Silhouette-PNG-HD-Quality.png

https://www.askpython.com/python/examples/n-grams-python-nltk

https://colab.research.google.com/?utm_source=scs-index

https://twitter.com/PresOfPakistan

https://www.nltk.org