**Submitted To: Sir Jamal Abdul Ahad**

**Submitted By: Areeba Khan**

**Roll No:        10172**

**Date:              30<sup>th</sup> Oct, 2024**

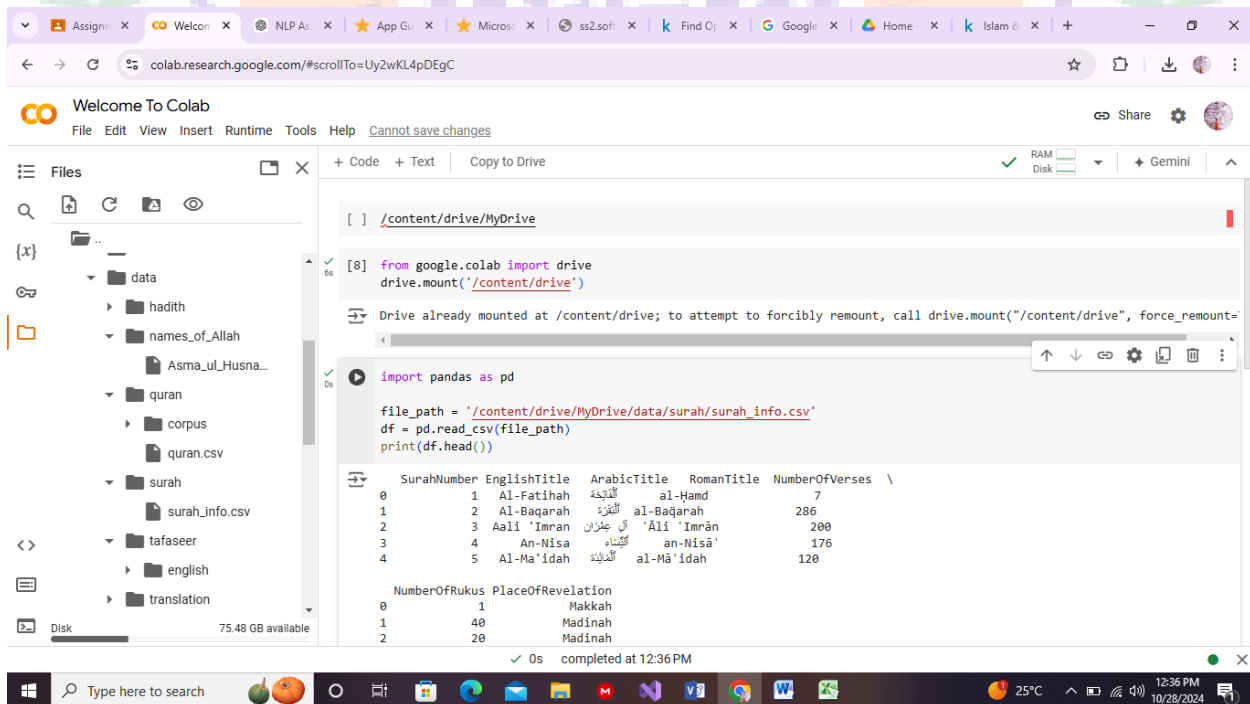**Subject:          Natural Language Processing**

## Question #1:

Choose any corpus of your choice of at least 200 MBs of any domain in NLP and perform

the following tasks:

• Text Preprocessing (Text Cleaning, Stemming / Lemmatization)

• Word Embedding (using an algorithm like Word2Vec, Glove, FastText)

• Encoding Techniques (Bag of Words, One – Hot)

• Parts of Speech tagging.

## Answer:

For this assignment, I chose the **Islam and AI Dataset** from **Kaggle,** as it offers a substantial corpus that combines discussions on religious perspectives and advancements in artificial intelligence, making it rich in context and language. This dataset will allow for comprehensive Natural Language Processing (NLP) analysis across different tasks. Firstly, text preprocessing, including text cleaning and lemmatization, will prepare the dataset by standardizing terms and removing extraneous characters, ensuring cleaner data for further analysis.

```
print(df.head()) # Print the first few rows to verify the data
```

```
   Arabic Name Name in English                     Name Meaning  \
0    الرَّحْمَٰن    AR-RAHMAAN                     The Beneficent
1    الرَّحِيم     AR RAHEEM                       The Merciful
2    الْمَلِك      AL MALIK                        The King
3    الْقُدُّوس    AL-QUDDUS     The Most Sacred / The Most Holy
4    السَّلَام     As-Salam                    The Giver of Peace

                                      Short Summary  \
0  He who wills goodness and mercy for all His cr...
1              He who acts with extreme kindness
2  The Sovereign Lord, The One with the complete ...
3  The one who is clear of any imperfection, weak...
4  The Most Perfect, The Source of Peace, The Giv...

                                       Long Summary  \
0  Allah is Ar-Rahmaan (in arabic: الرَّحْمَٰن...
1  The name Ar-Raheem (in Arabic: الرَّحِيم) c...
2  Allah is Al-Malik (in arabic: الْمَلِك), the ...
3  Allah سُبْحَٰنَهُ وَتَعَٰلَىٰ is Al-Quddus (i...
4  Allah is As-Salam (in Arabic: السَّلَام); He...

                                       Arabic Root  \
0  From the root ra-ha-mim (ر ح م), which has the...
1  From the root ra-ha-mim (ر ح م), which has the...
2  From the root mim-lam-kaf (م ل ك), which has t...
3  From the root qaf-dal-sin (ق د س), which has t...
4  From the root sin-lam-mim (س ل م), which has t...
```
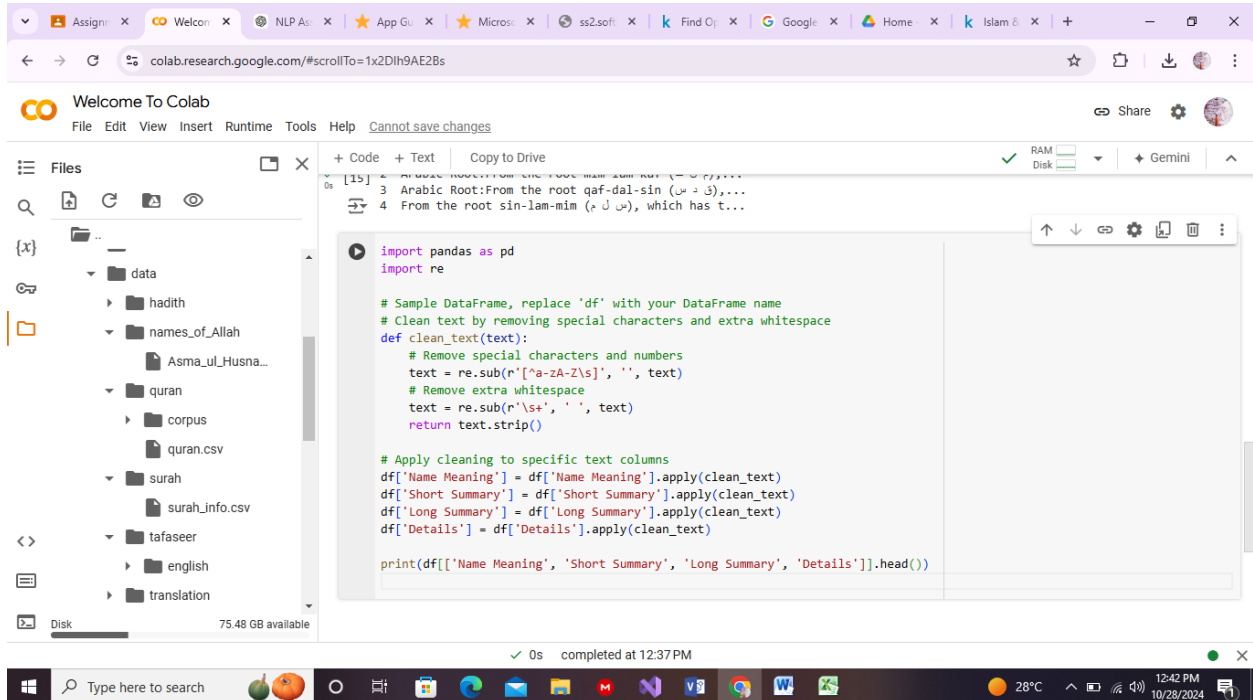
✓ 0s   completed at 12:37 PM

---

```
4  The Most Perfect, The Source of Peace, The Giv...

                                       Long Summary  \
0  Allah is Ar-Rahmaan (in arabic: الرَّحْمَٰن...
1  The name Ar-Raheem (in Arabic: الرَّحِيم) c...
2  Allah is Al-Malik (in arabic: الْمَلِك), the ...
3  Allah سُبْحَٰنَهُ وَتَعَٰلَىٰ is Al-Quddus (i...
4  Allah is As-Salam (in Arabic: السَّلَام); He...

                                       Arabic Root  \
0  From the root ra-ha-mim (ر ح م), which has the...
1  From the root ra-ha-mim (ر ح م), which has the...
2  From the root mim-lam-kaf (م ل ك), which has t...
3  From the root qaf-dal-sin (ق د س), which has t...
4  From the root sin-lam-mim (س ل م), which has t...

                                           Details
0  Arabic Root:From the root ra-ha-mim (ر ح م), w...
1  Arabic Root:From the root ra-ha-mim (ر ح م), w...
2  Arabic Root:From the root mim-lam-kaf (م ل ك),...
3  Arabic Root:From the root qaf-dal-sin (ق د س),...
4  Arabic Root:From the root sin-lam-mim (س ل م), which has t...
```

✓ 0s   completed at 12:37 PM

# Text Preprocessing (Text Cleaning, Stemming / Lemmatization):
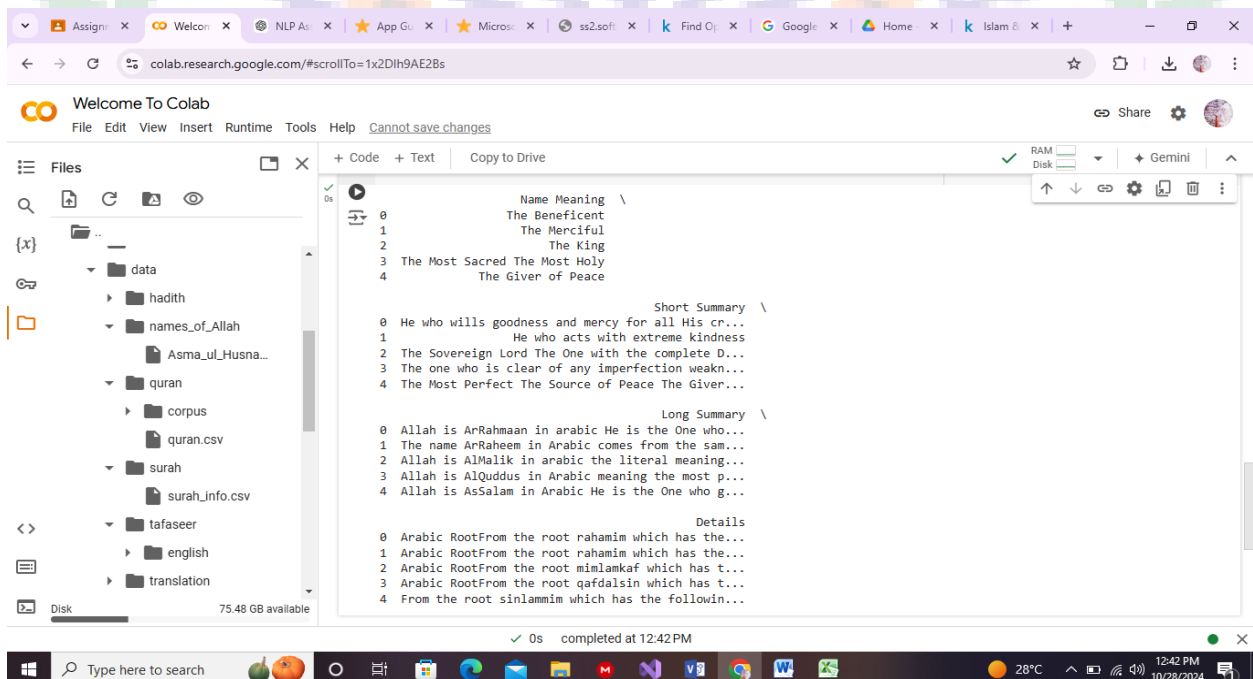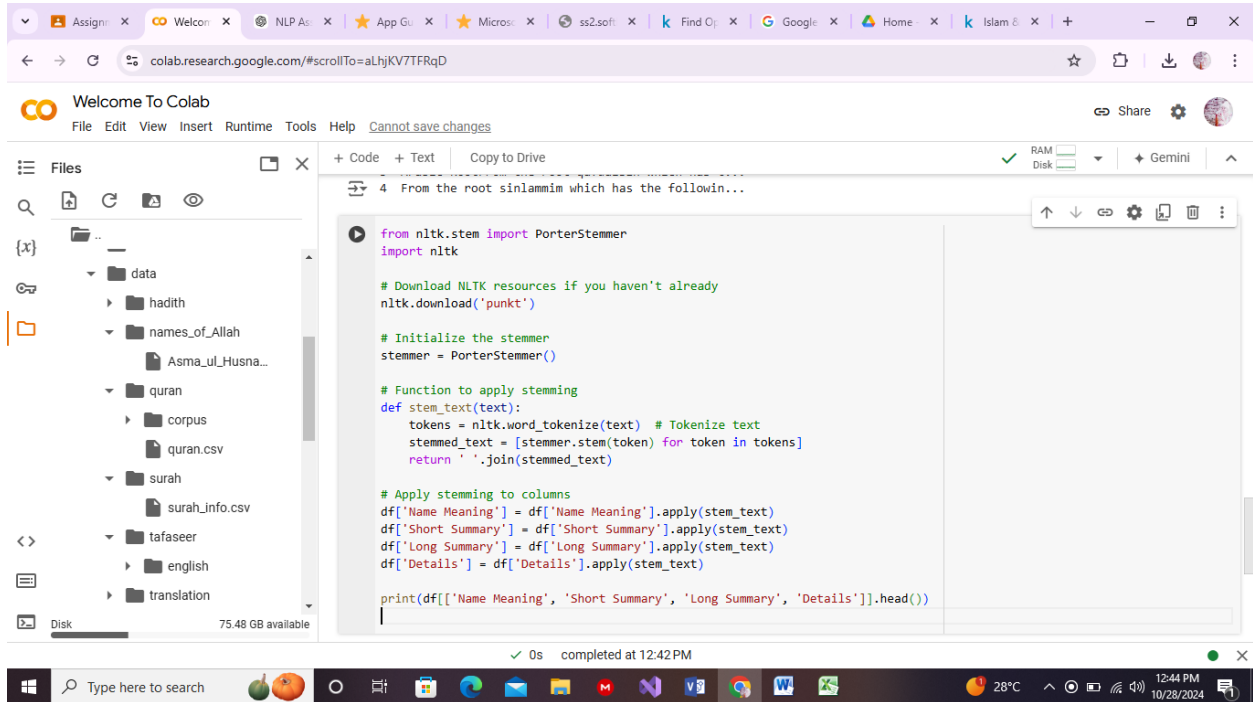
## Text cleaning:

## Code:



```python
import pandas as pd
import re

# Sample DataFrame, replace 'df' with your DataFrame name
# Clean text by removing special characters and extra whitespace
def clean_text(text):
    # Remove special characters and numbers
    text = re.sub(r'[^a-zA-Z\s]', '', text)
    # Remove extra whitespace
    text = re.sub(r'\s+', ' ', text)
    return text.strip()

# Apply cleaning to specific text columns
df['Name Meaning'] = df['Name Meaning'].apply(clean_text)
df['Short Summary'] = df['Short Summary'].apply(clean_text)
df['Long Summary'] = df['Long Summary'].apply(clean_text)
df['Details'] = df['Details'].apply(clean_text)

print(df[['Name Meaning', 'Short Summary', 'Long Summary', 'Details']].head())
```

## Output:

# Stemming:

## Code:



```python
from nltk.stem import PorterStemmer
import nltk

# Download NLTK resources if you haven't already
nltk.download('punkt')

# Initialize the stemmer
stemmer = PorterStemmer()

# Function to apply stemming
def stem_text(text):
    tokens = nltk.word_tokenize(text)  # Tokenize text
    stemmed_text = [stemmer.stem(token) for token in tokens]
    return ' '.join(stemmed_text)

# Apply stemming to columns
df['Name Meaning'] = df['Name Meaning'].apply(stem_text)
df['Short Summary'] = df['Short Summary'].apply(stem_text)
df['Long Summary'] = df['Long Summary'].apply(stem_text)
df['Details'] = df['Details'].apply(stem_text)

print(df[['Name Meaning', 'Short Summary', 'Long Summary', 'Details']].head())
```
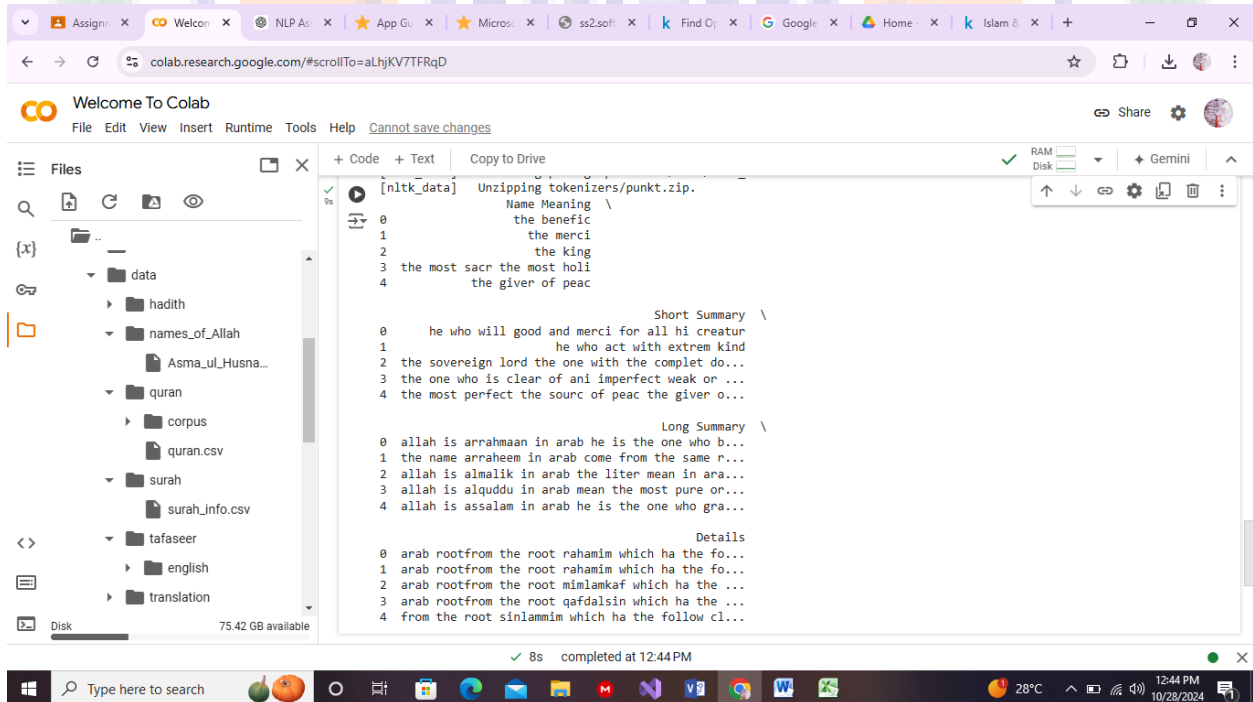
## Output:



```
[nltk_data]    Unzipping tokenizers/punkt.zip.
                    Name Meaning  \
0                     the benefic
1                      the merci
2                       the king
3       the most sacr the most holi
4                  the giver of peac

                               Short Summary  \
0       he who will good and merci for all hi creatur
1                     he who act with extrem kind
2       the sovereign lord the one with the complet do...
3       the one who is clear of ani imperfect weak or ...
4       the most perfect the sourc of peac the giver o...

                               Long Summary  \
0   allah is arrahmaan in arab he is the one who b...
1   the name arraheem in arab come from the same r...
2   allah is almalik in arab the liter mean in ara...
3   allah is alquddu in arab mean the most pure or...
4   allah is assalam in arab he is the one who gra...

                                      Details
0   arab rootfrom the root rahamim which ha the fo...
1   arab rootfrom the root rahamim which ha the fo...
2   arab rootfrom the root mimlamkaf which ha the ...
3   arab rootfrom the root qafdalsin which ha the ...
4   from the root sinlammim which ha the follow cl...
```

## Lemmatization:

## Code:



```python
from nltk.stem import WordNetLemmatizer

# Download necessary NLTK resources
nltk.download('wordnet')
nltk.download('omw-1.4')

# Initialize the lemmatizer
lemmatizer = WordNetLemmatizer()

# Function to apply lemmatization
def lemmatize_text(text):
    tokens = nltk.word_tokenize(text)
    lemmatized_text = [lemmatizer.lemmatize(token) for token in tokens]
    return ' '.join(lemmatized_text)

# Apply lemmatization to columns
df['Name Meaning'] = df['Name Meaning'].apply(lemmatize_text)
df['Short Summary'] = df['Short Summary'].apply(lemmatize_text)
df['Long Summary'] = df['Long Summary'].apply(lemmatize_text)
df['Details'] = df['Details'].apply(lemmatize_text)

print(df[['Name Meaning', 'Short Summary', 'Long Summary', 'Details']].head())
```

## Output:



```
[nltk_data] Downloading package wordnet to /root/nltk_data...
[nltk_data] Downloading package omw-1.4 to /root/nltk_data...
                          Name Meaning  \
0                         the benefic
1                         the merci
2                         the king
3    the most sacr the most holi
4               the giver of peac

                                    Short Summary  \
0     he who will good and merci for all hi creatur
1                       he who act with extrem kind
2    the sovereign lord the one with the complet do...
3    the one who is clear of ani imperfect weak or ...
4    the most perfect the sourc of peac the giver o...

                                     Long Summary  \
0    allah is arrahmaan in arab he is the one who b...
1    the name arraheem in arab come from the same r...
2    allah is almalik in arab the liter mean in ara...
3    allah is alquddu in arab mean the most pure or...
4    allah is assalam in arab he is the one who gra...

                                          Details
0    arab rootfrom the root rahamim which ha the fo...
1    arab rootfrom the root rahamim which ha the fo...
2    arab rootfrom the root mimlamkaf which ha the ...
3    arab rootfrom the root qafdalsin which ha the ...
4    from the root sinlammim which ha the follow cl...
```

# Word Embedding (using an algorithm like Word2Vec, Glove, FastText):

To implement word embedding on this dataset . I'll use the **Word2Vec** algorithm, which creates word vectors by examining co-occurrence within a certain context.

## Text Preprocessing:

We'll clean and tokenize the text to prepare it for Word2Vec.

## Word2Vec Training:

## Code:



```python
from gensim.models import Word2Vec

# Build the Word2Vec model
word2vec_model = Word2Vec(vector_size=100, window=5, min_count=1, workers=4)

# Build the vocabulary from the tokenized sentences
word2vec_model.build_vocab(sentences)

# Train the Word2Vec model
word2vec_model.train(sentences, total_examples=word2vec_model.corpus_count, epochs=10)

# Check if the word is in the vocabulary before accessing its vector
word_to_check = 'An-Nisa'  # Replace with the word you want to check

if word_to_check in word2vec_model.wv.key_to_index:
    print(word2vec_model.wv[word_to_check])
else:
    print(f"The word '{word_to_check}' is not in the model's vocabulary.")
```

```
[-8.7274825e-03  2.1301615e-03 -8.7354420e-04 -9.3190884e-03
 -9.4281426e-03 -1.4107180e-03  4.4324086e-03  3.7040710e-03
 -6.4986930e-03 -6.8730675e-03 -4.9994122e-03 -2.2868442e-03
 -7.2502876e-03 -9.6033178e-03 -2.7436293e-03 -8.3628409e-03
 -6.0388758e-03 -5.6709289e-03 -2.3441375e-03 -1.7069972e-03
 -8.9569986e-03 -7.3519943e-04  8.1525063e-03  7.6904297e-03
 -7.2061159e-03 -3.6668312e-03  3.1185520e-03 -9.5707225e-03
```

## Output:



```python
    word_to_check = 'An-Nisa'  # Replace with the word you want to check

    if word_to_check in word2vec_model.wv.key_to_index:
        print(word2vec_model.wv[word_to_check])
    else:
        print(f"The word '{word_to_check}' is not in the model's vocabulary.")
```

```
[-8.7274825e-03  2.1301615e-03 -8.7354420e-04 -9.3190884e-03
 -9.4281426e-03 -1.4107180e-03  4.4324086e-03  3.7040710e-03
 -6.4986930e-03 -6.8730675e-03 -4.9994122e-03 -2.2868442e-03
 -7.2502876e-03 -9.6033178e-03 -2.7436293e-03 -8.3628409e-03
 -6.0388758e-03 -5.6709289e-03 -2.3441375e-03 -1.7069972e-03
 -8.9569986e-03 -7.3519943e-04  8.1525063e-03  7.6904297e-03
 -7.2061159e-03 -3.6668312e-03  3.1185520e-03 -9.5707225e-03
  1.4764392e-03  6.5244664e-03  5.7464195e-03 -8.7630618e-03
 -4.5171441e-03 -8.1401607e-03  4.5956374e-05  9.2636338e-03
  5.9733056e-03  5.0673080e-03  5.0610625e-03 -3.2429171e-03
  9.5521836e-03 -7.3564244e-03 -7.2703874e-03 -2.2653891e-03
 -7.7856064e-04 -3.2161034e-03 -5.9258583e-04  7.4888230e-03
 -6.9751858e-04 -1.6249407e-03  2.7443992e-03 -8.3591007e-03
  7.8558037e-03  8.5361041e-03 -9.5840869e-03  2.4462664e-03
  9.9049713e-03 -7.6658037e-03 -6.9669187e-03 -7.7365171e-03
  8.3959233e-03 -6.8133592e-04  9.1444086e-03 -8.1582209e-03
  3.7430846e-03  2.6350426e-03  7.4271322e-04  2.3276759e-03
 -7.4690939e-03 -9.3583735e-03  2.3545765e-03  6.1484552e-03
  7.9856887e-03  5.7358947e-03 -7.7733636e-04  8.3061643e-03
 -9.3363142e-03  3.4061326e-03  2.6675343e-04  3.8572443e-03
  7.3857834e-03 -6.7251669e-03  5.5844807e-03 -9.5222248e-03
 -8.0445886e-04 -8.6887367e-03 -5.0986730e-03  9.2892265e-03
```

# Encoding Techniques (Bag of Words, One – Hot):

## Bag Of Words:

## Code:



## Output:

# One – Hot:

## Coding:



## Output:

**Parts of Speech tagging:**

**Code + Output:**



## Question #2:
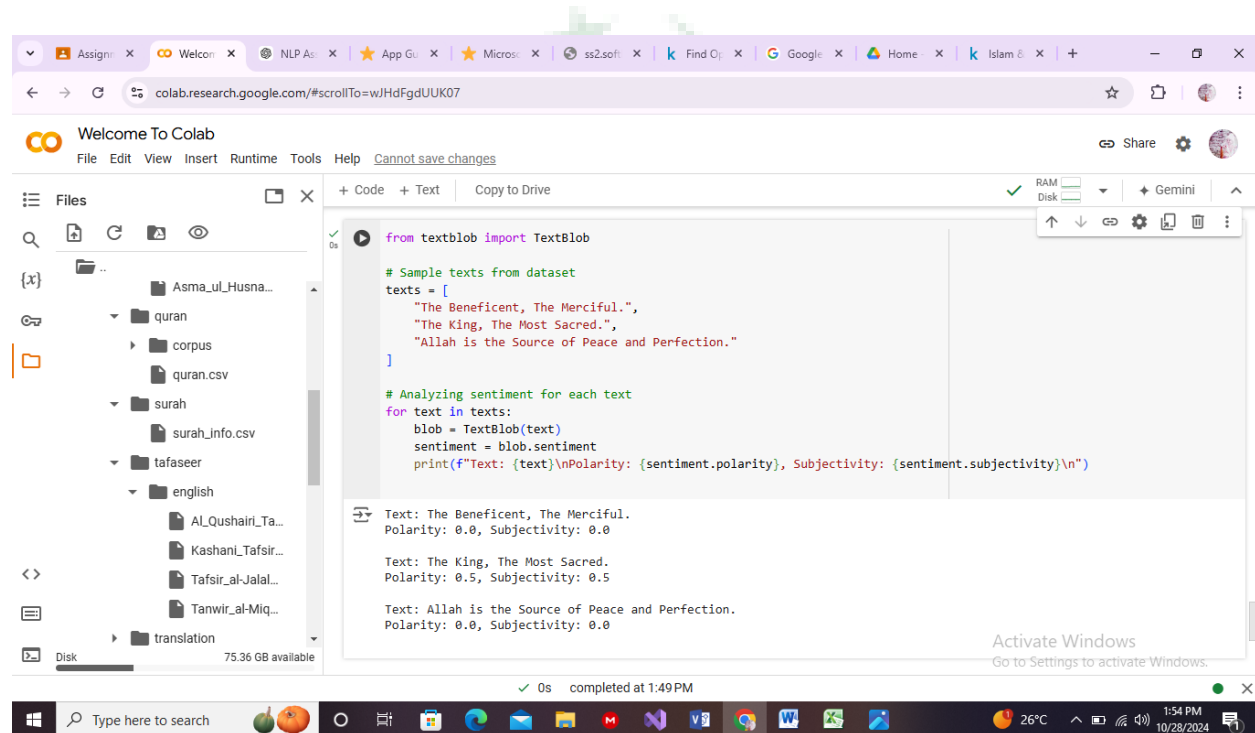
Perform any two of the basic NLP tasks listed below.

• Sentiment Analysis (using VADER, TextBlob)

• Named Entity Recognition (NER)

• Text Classification (Naive Bayes, Logistic Regression, SVM)

• Language Models (N-grams, Markov Chains)

• Topic Modeling (LDA, Latent Semantic Analysis)

**Answer:**

# 1. Sentiment Analysis (Using TextBlob):

Sentiment analysis will evaluate the emotional tone in the texts. Although religious and AI content might have more neutral sentiments, it's still insightful for polarity analysis.

## Code + Output:



# Topic Modeling (Using Latent Dirichlet Allocation - LDA):

Topic modeling will reveal the main themes within the dataset.

## Code:

## Output:

```
# Vectorize the text for LDA
vectorizer = CountVectorizer(stop_words='english')
text_vector = vectorizer.fit_transform(texts)

# Initialize LDA with 2 topics
lda_model = LatentDirichletAllocation(n_components=2, random_state=42)
lda_model.fit(text_vector)

# Display topics
words = vectorizer.get_feature_names_out()
for index, topic in enumerate(lda_model.components_):
    print(f"Topic #{index+1}:")
    print([words[i] for i in topic.argsort()[-5:]])
```

```
Topic #1:
['king', 'sovereign', 'peace', 'sacred', 'giver']
Topic #2:
['al', 'raheem', 'allah', 'ar', 'malik']
```