

Data and Artificial Intelligence

Cyber Shujaa Program

Week 11 Assignment

Natural Language Processing

Student Name: Denis Kombe

Student ID: CS-DA01-25094

Table of Contents

Objectives	3
Introduction.....	3
Tasks Completed	4
Step 1: Importing and loading the pre-trained BERT tokenizer.	4
Step 2: Listing the 10 sentence pairs.	5
Step 3: Compare sentence pairs using BERT embeddings.	6
Step 4: Predicting semantic similarities using a defined threshold.	7
Step 5: Evaluating prediction accuracy	8
Step 6: Explaining the NLP concepts in this assignment	9
Link to Google Colab Notebook.....	10
Notebook Link to Code:.....	10
Conclusion	10

Objectives

1. Importing and loading the pre-trained BERT tokenizer.
2. Listing 10 sentences.
3. Compare sentence pairs using BERT embeddings.
4. Predicting semantic similarity using a defined threshold.
5. Evaluate prediction accuracy against ground truth.
6. Explaining the NLP concepts used in this assignment.

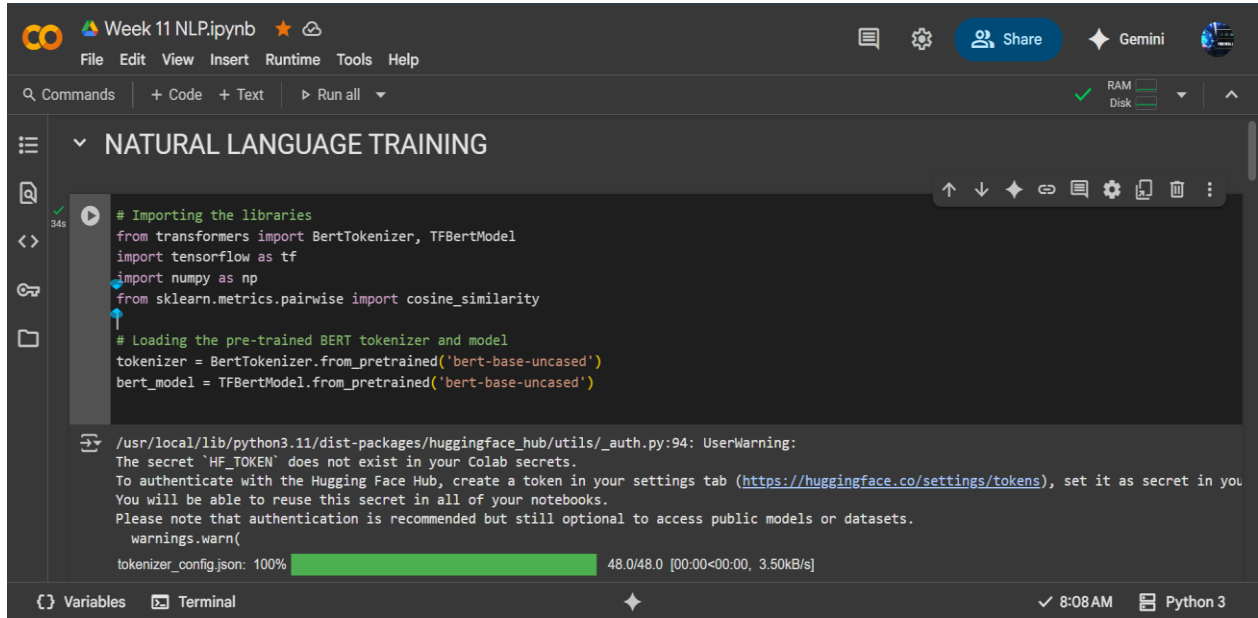
Introduction

This assignment focused on exploring sentence similarity using a pre-trained BERT model. The goal was to analyze how BERT understands the meaning of different sentences in context and how this can be used to measure how similar or different they are. By comparing 10 sentence pairs, I applied cosine similarity to the embeddings and used a set threshold to predict whether the sentences were similar or not. This assignment made me understand important NLP concepts and how they apply in real tasks.

Tasks Completed

Step 1: Importing and loading the pre-trained BERT tokenizer.

Started by importing the necessary libraries needed for this project and loading the pre-trained BERT tokenizer and model.



The screenshot shows a Jupyter Notebook titled "Week 11 NLP.ipynb" in a dark-themed interface. The notebook is open to a cell containing Python code for importing libraries and loading the BERT tokenizer and model. The code is as follows:

```
# Importing the libraries
from transformers import BertTokenizer, TFBertModel
import tensorflow as tf
import numpy as np
from sklearn.metrics.pairwise import cosine_similarity

# Loading the pre-trained BERT tokenizer and model
tokenizer = BertTokenizer.from_pretrained('bert-base-uncased')
bert_model = TFBertModel.from_pretrained('bert-base-uncased')
```

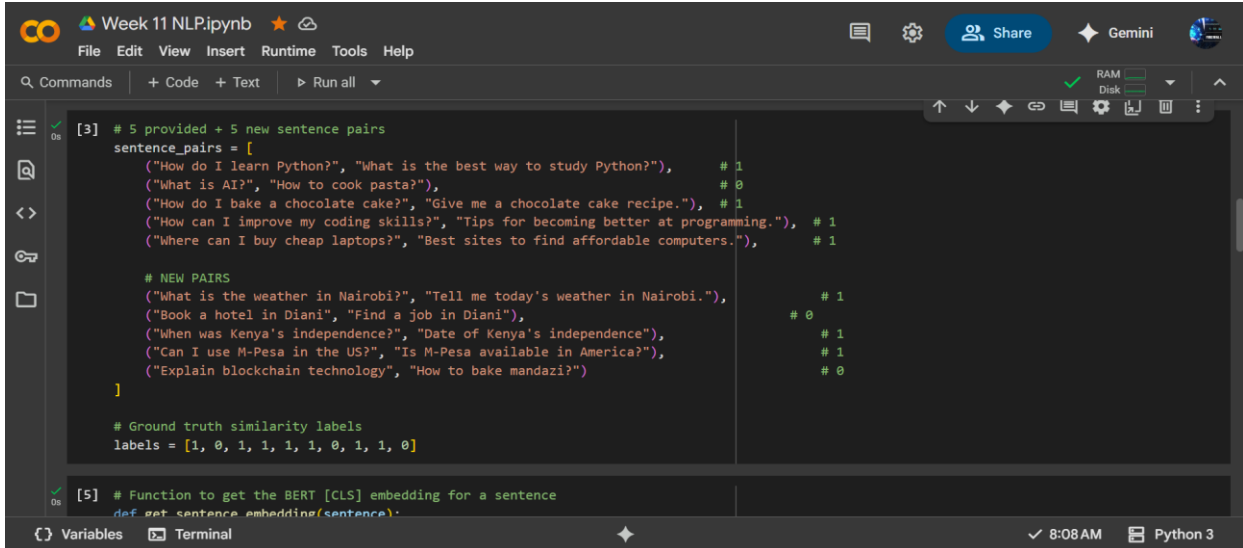
Below the code cell, there is a warning message from the Hugging Face Hub:

```
/usr/local/lib/python3.11/dist-packages/huggingface_hub/utils/_auth.py:94: UserWarning:
The secret `HF_TOKEN` does not exist in your Colab secrets.
To authenticate with the Hugging Face Hub, create a token in your settings tab (https://huggingface.co/settings/tokens), set it as secret in you
You will be able to reuse this secret in all of your notebooks.
Please note that authentication is recommended but still optional to access public models or datasets.
warnings.warn(
tokenizer_config.json: 100% 48.0/48.0 [00:00<00:00, 3.50kB/s]
```

The bottom of the interface shows the "Variables" and "Terminal" tabs, and the status bar indicates the time is 8:08 AM and the environment is Python 3.

Step 2: Listing the 10 sentence pairs.

I added 5 more sentences including and labeled each sentence based on the similarities they had.



```
Week 11 NLP.ipynb
File Edit View Insert Runtime Tools Help
Q Commands + Code + Text ▶ Run all
[3] # 5 provided + 5 new sentence pairs
sentence_pairs = [
    ("How do I learn Python?", "What is the best way to study Python?"), # 1
    ("What is AI?", "How to cook pasta?"), # 0
    ("How do I bake a chocolate cake?", "Give me a chocolate cake recipe."), # 1
    ("How can I improve my coding skills?", "Tips for becoming better at programming."), # 1
    ("Where can I buy cheap laptops?", "Best sites to find affordable computers."), # 1

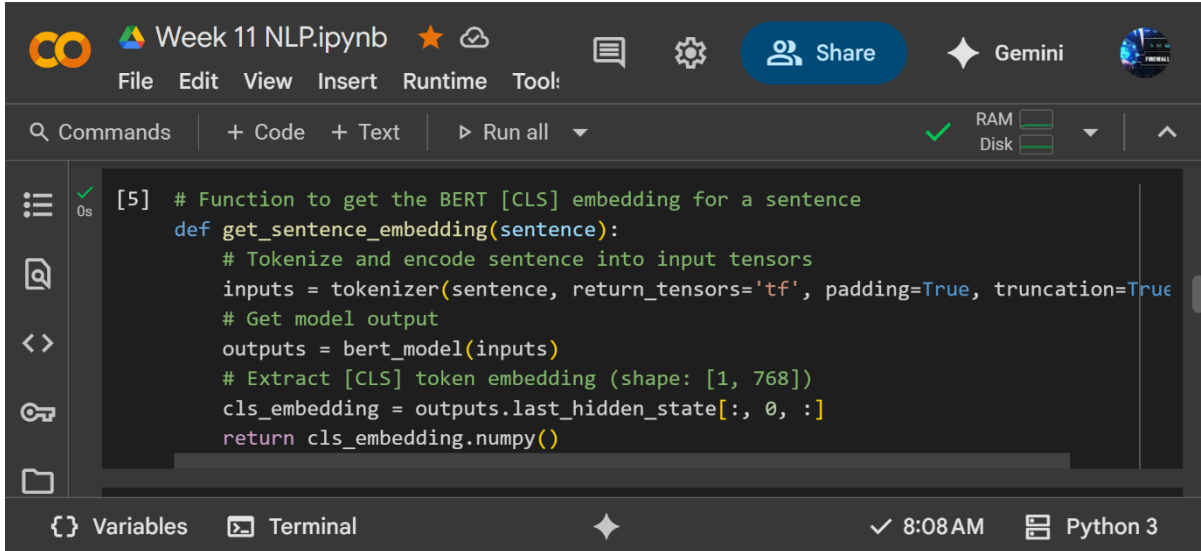
    # NEW PAIRS
    ("What is the weather in Nairobi?", "Tell me today's weather in Nairobi."), # 1
    ("Book a hotel in Diani", "Find a job in Diani"), # 0
    ("When was Kenya's independence?", "Date of Kenya's independence"), # 1
    ("Can I use M-Pesa in the US?", "Is M-Pesa available in America?"), # 1
    ("Explain blockchain technology", "How to bake mandazi?") # 0
]

# Ground truth similarity labels
labels = [1, 0, 1, 1, 1, 1, 0, 1, 1, 0]

[5] # Function to get the BERT [CLS] embedding for a sentence
def get_sentence_embedding(sentence):
```

Step 3: Compare sentence pairs using BERT embeddings.

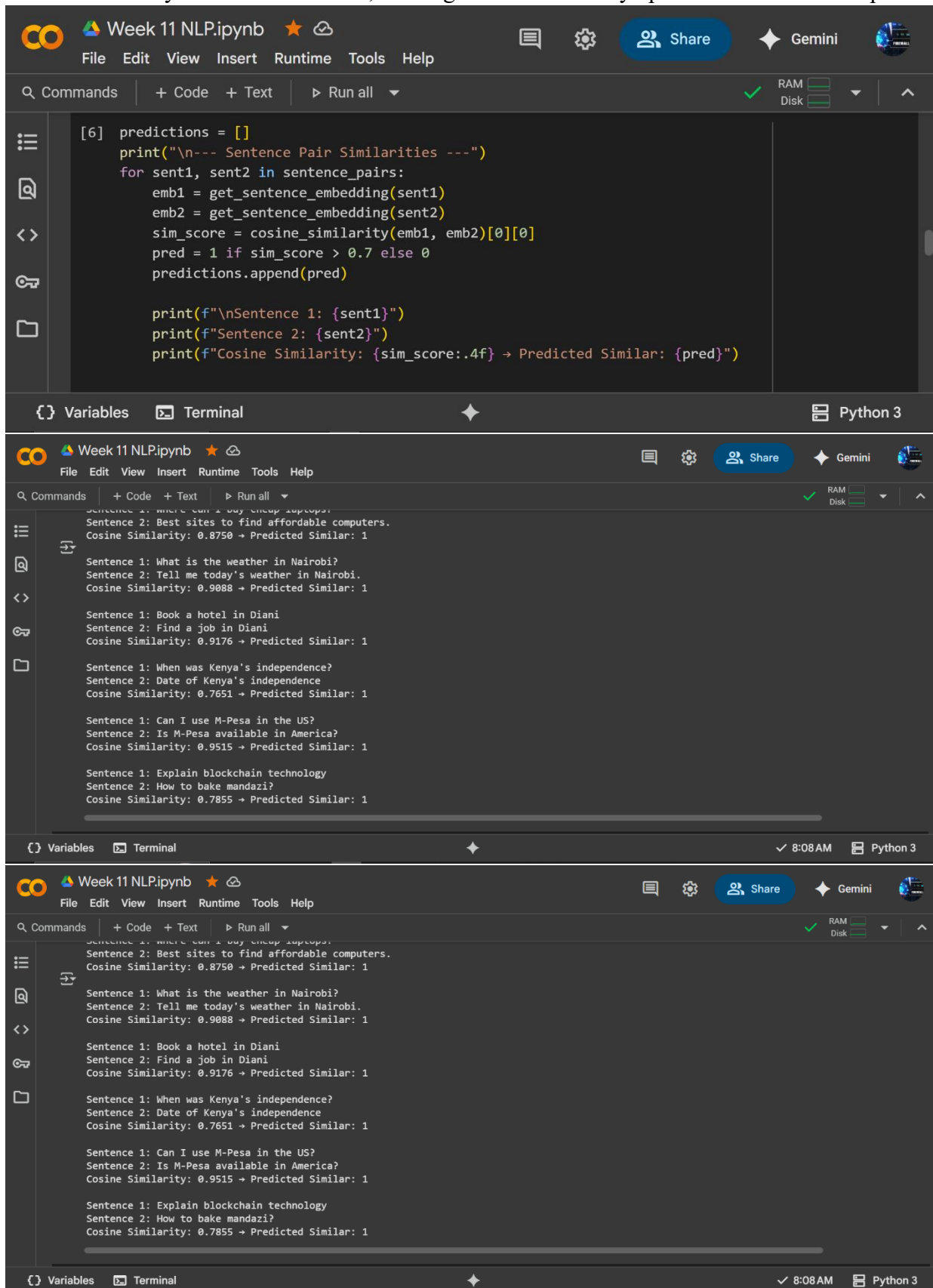
Applying a pre-trained BERT model for encoding a total of 10 sentence pair into contextual embeddings. Each pair is processed to extract sentence-level meaning using the CLS token.



```
[5] # Function to get the BERT [CLS] embedding for a sentence
def get_sentence_embedding(sentence):
    # Tokenize and encode sentence into input tensors
    inputs = tokenizer(sentence, return_tensors='tf', padding=True, truncation=True)
    # Get model output
    outputs = bert_model(inputs)
    # Extract [CLS] token embedding (shape: [1, 768])
    cls_embedding = outputs.last_hidden_state[:, 0, :]
    return cls_embedding.numpy()
```

Step 4: Predicting semantic similarities using a defined threshold.

To classify each sentence pair as either similar (label = 1) or not similar (label = 0) using a cosine similarity threshold of 0.7, and generate a binary prediction for each pair.



The following code is shown in the Jupyter Notebook:

```
[6] predictions = []
print("\n--- Sentence Pair Similarities ---")
for sent1, sent2 in sentence_pairs:
    emb1 = get_sentence_embedding(sent1)
    emb2 = get_sentence_embedding(sent2)
    sim_score = cosine_similarity(emb1, emb2)[0][0]
    pred = 1 if sim_score > 0.7 else 0
    predictions.append(pred)

print(f"\nSentence 1: {sent1}")
print(f"Sentence 2: {sent2}")
print(f"Cosine Similarity: {sim_score:.4f} → Predicted Similar: {pred}")
```

The output of the code is shown in the terminal, displaying the cosine similarity and predicted similarity for several sentence pairs:

```
Sentence 1: Where can I buy cheap laptops?
Sentence 2: Best sites to find affordable computers.
Cosine Similarity: 0.8750 → Predicted Similar: 1

Sentence 1: What is the weather in Nairobi?
Sentence 2: Tell me today's weather in Nairobi.
Cosine Similarity: 0.9088 → Predicted Similar: 1

Sentence 1: Book a hotel in Diani
Sentence 2: Find a job in Diani
Cosine Similarity: 0.9176 → Predicted Similar: 1

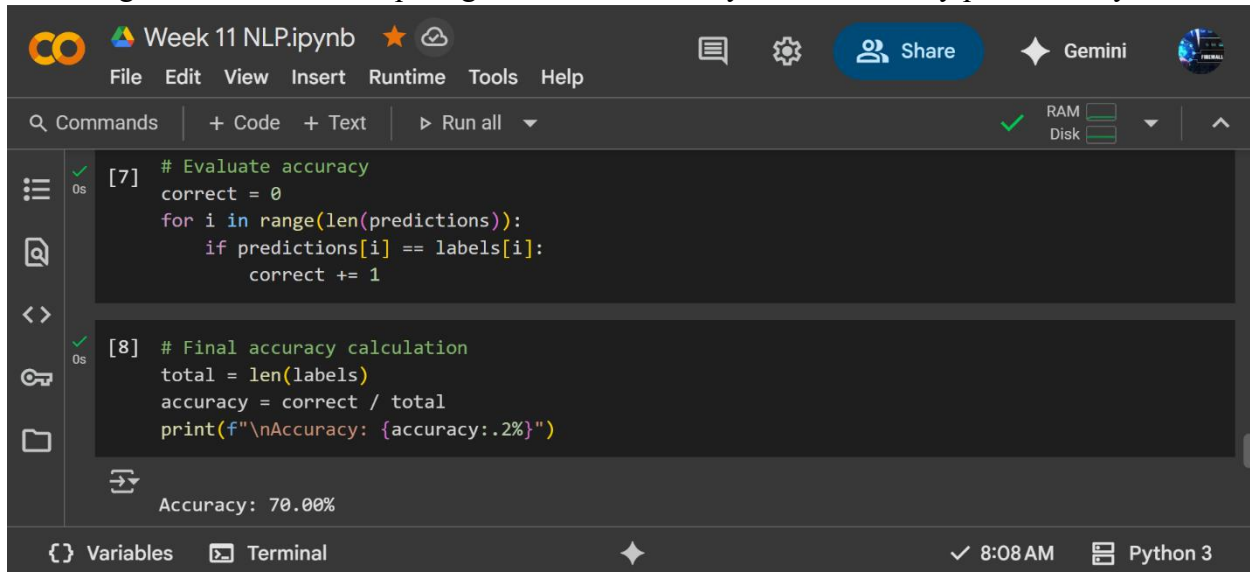
Sentence 1: When was Kenya's independence?
Sentence 2: Date of Kenya's independence
Cosine Similarity: 0.7651 → Predicted Similar: 1

Sentence 1: Can I use M-Pesa in the US?
Sentence 2: Is M-Pesa available in America?
Cosine Similarity: 0.9515 → Predicted Similar: 1

Sentence 1: Explain blockchain technology
Sentence 2: How to bake mandazi?
Cosine Similarity: 0.7855 → Predicted Similar: 1
```

Step 5: Evaluating prediction accuracy

Assigning similarity labels for all 10 sentence pairs, then comparing the model's predictions with the ground truth and computing the overall accuracy of the similarity prediction system.



```
Week 11 NLP.ipynb
```

File Edit View Insert Runtime Tools Help

Commands + Code + Text Run all

RAM Disk

```
[7] # Evaluate accuracy
correct = 0
for i in range(len(predictions)):
    if predictions[i] == labels[i]:
        correct += 1
```

```
[8] # Final accuracy calculation
total = len(labels)
accuracy = correct / total
print(f"\nAccuracy: {accuracy:.2%}")
```

Accuracy: 70.00%

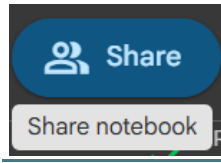
Variables Terminal 8:08 AM Python 3

Step 6: Explaining the NLP concepts in this assignment

1. How does BERT differ from traditional NLP approaches like Bag of Words or TF-IDF??
 - Traditional NLP approaches like Bag of Words or TF-IDF create fixed-size vectors based only on word frequency hence ignoring the meaning or order of words.
 - BERT understands context and word meaning in a sentence using its attention-based architecture. It generates contextual embeddings which differ based on sentence meaning.
2. What is the role of the encoder in the BERT model, and how is it used in this assignment?
 - The encoder in BERT learns the relationships between words in a sentence.
 - In this assignment, sentences are fed into BERT and the encoder outputs embeddings representing the meaning of each token. The CLS token's output from the encoder is used as a representation of the entire sentence.
3. What are contextual embeddings? How are they generated and used in this code?
 - Contextual embeddings are word vectors that change depending on surrounding words. For example, "bank" in "river bank" vs "financial bank" will have different vectors.
 - BERT generates these using self-attention and multiple transformer layers.
 - In this assignment, the CLS token embedding is used as a context-aware sentence representation to compare similarity between sentence pairs.
4. Why is the [CLS] token used for sentence similarity in this code?
 - To compare sentence meaning using cosine similarity.
5. What is cosine similarity, and why is it useful in comparing embeddings?
 - Cosine similarity measures the angle between two vectors and ranges from -1 (opposite) and 1 (identical).
 - It is useful because it helps quantify how similar two sentence embeddings are.

Link to Google Colab Notebook

Notebook Link to Code:



Conclusion

From this assignment, I have learned how powerful BERT is when it comes to understanding the actual meaning of words depending on how they are used in a sentence. Unlike traditional methods, BERT gives context to words, and that makes a big difference when measuring similarity. I also got to work with cosine similarity and understand how it helps us decide how close two sentences are in meaning. Overall, the experience gave me confidence in using pre-trained models and strengthened my understanding of contextual embeddings.