BUSINESS INTELLIGENCE AND
**DATA ANALYTICS**

**IQA** Institute of Analytics

Open in Colab

## Context Aware Word Embeddings - BERT

```
1 %pip install transformers torch
```

```
Requirement already satisfied: transformers in c:\users\me\appdata\local\programs\python\python310\lib\site-packages (4.51.3)
Requirement already satisfied: tqdm>=4.27 in c:\users\me\appdata\local\programs\python\python310\lib\site-packages (from transformers) (
Requirement already satisfied: safetensors>=0.4.3 in c:\users\me\appdata\local\programs\python\python310\lib\site-packages (from transfo
Requirement already satisfied: packaging>=20.0 in c:\users\me\appdata\roaming\python\python310\site-packages (from transformers) (24.2)
Requirement already satisfied: numpy>=1.17 in c:\users\me\appdata\local\programs\python\python310\lib\site-packages (from transformers)
Requirement already satisfied: tokenizers<0.22,>=0.21 in c:\users\me\appdata\local\programs\python\python310\lib\site-packages (from tra
Requirement already satisfied: filelock in c:\users\me\appdata\local\programs\python\python310\lib\site-packages (from transformers) (3.
Requirement already satisfied: regex!=2019.12.17 in c:\users\me\appdata\local\programs\python\python310\lib\site-packages (from transfor
Requirement already satisfied: huggingface-hub<1.0,>=0.30.0 in c:\users\me\appdata\local\programs\python\python310\lib\site-packages (fr
Requirement already satisfied: requests in c:\users\me\appdata\local\programs\python\python310\lib\site-packages (from transformers) (2.
Requirement already satisfied: pyyaml>=5.1 in c:\users\me\appdata\local\programs\python\python310\lib\site-packages (from transformers)
Requirement already satisfied: fsspec>=2023.5.0 in c:\users\me\appdata\local\programs\python\python310\lib\site-packages (from huggingfa
Requirement already satisfied: typing-extensions>=3.7.4.3 in c:\users\me\appdata\roaming\python\python310\site-packages (from huggingfac
Requirement already satisfied: colorama in c:\users\me\appdata\roaming\python\python310\site-packages (from tqdm>=4.27->transformers) (0
Requirement already satisfied: charset-normalizer<4,>=2 in c:\users\me\appdata\local\programs\python\python310\lib\site-packages (from r
Requirement already satisfied: idna<4,>=2.5 in c:\users\me\appdata\local\programs\python\python310\lib\site-packages (from requests->tra
Requirement already satisfied: urllib3<3,>=1.21.1 in c:\users\me\appdata\local\programs\python\python310\lib\site-packages (from request
Requirement already satisfied: certifi>=2017.4.17 in c:\users\me\appdata\local\programs\python\python310\lib\site-packages (from request
Note: you may need to restart the kernel to use updated packages.
WARNING: You are using pip version 21.2.3; however, version 25.0.1 is available.
You should consider upgrading via the 'c:\Users\me\AppData\Local\Programs\Python\Python310\python.exe -m pip install --upgrade pip' comm
```

## Display BERT Embeddings

```
 1 from transformers import BertTokenizer, BertModel
 2 import torch
 3
 4 # Load pretrained BERT
 5 tokenizer = BertTokenizer.from_pretrained('bert-base-uncased')
 6 model = BertModel.from_pretrained('bert-base-uncased')
 7
 8 # Sentence
 9 sentence = "He went to the bank to deposit money."
10
11 # Tokenize
12 inputs = tokenizer(sentence, return_tensors='pt')
13 tokens = tokenizer.convert_ids_to_tokens(inputs['input_ids'][0])
14
15 # Get outputs
16 with torch.no_grad():
17     outputs = model(**inputs)
18
19 # Get hidden states (embeddings)
20 embeddings = outputs.last_hidden_state.squeeze(0)  # shape: (seq_len, hidden_size)
21 # Print tokens and their embeddings
22 for token, embedding in zip(tokens, embeddings):
23     print(f"Token: {token}\n Embedding(First 10 Numbers):\n{embedding[:10]}\nShape: {embedding.shape}\n")
24
25
26 # # Find index of "bank"
27 # try:
28 #     idx = tokens.index("bank")
29 #     bank_embedding = embeddings[idx]
30 #     print(f"Embedding for 'bank':\n{bank_embedding}\n\nShape: {bank_embedding.shape}")
31 # except ValueError:
32 #     print("'bank' not found in tokenized input:", tokens)
33
```

```
Embedding for 'bank':
tensor([ 4.7019e-01, -1.9835e-01, -1.0122e-01, -1.3519e-01,  1.2612e+00,
        -9.6139e-03, -4.9014e-02,  1.0147e+00, -4.5361e-02,  1.7432e-01,
```

```
            1.2800e-01, -3.2356e-01, -1.3227e-01,  3.6582e-02, -7.8302e-01,
           -6.2770e-01,  5.2776e-01,  3.5693e-01,  1.3597e+00,  2.3784e-01,
           -3.0995e-01,  4.3136e-02,  3.2358e-01,  3.2144e-01,  3.3207e-01,
            4.5470e-01,  6.8660e-01,  5.2037e-01, -2.8076e-01, -5.2107e-01,
            5.3412e-01,  9.5313e-01,  3.6960e-01,  4.9074e-01,  1.0348e-01,
           -1.2543e-01,  1.8115e-01,  3.9604e-02, -1.1310e+00,  2.2161e-02,
           -4.4877e-01, -8.1382e-01, -6.2421e-01,  3.5284e-01, -2.4929e-01,
           -6.1539e-01,  1.9276e-01,  2.8171e-01, -7.0082e-01, -8.2422e-01,
           -3.0416e-01,  1.0278e+00,  4.3732e-01, -5.0054e-01,  1.1097e-01,
            4.7545e-01, -1.0476e+00, -4.6538e-01, -5.3300e-01, -2.1977e-01,
            7.0954e-01,  3.1443e-01,  5.0420e-01, -7.7659e-01,  2.3119e-01,
           -1.6568e-01,  4.8205e-01, -5.2181e-03, -1.1804e-01, -9.9157e-01,
           -2.4582e-01,  1.9177e+00, -8.5401e-01, -9.5023e-01, -4.6143e-01,
            5.0577e-01, -2.2459e-01,  2.9306e-01, -2.2287e-01, -9.2279e-03,
           -5.4351e-01, -7.2676e-01, -1.5370e-01,  3.8075e-01,  3.8973e-01,
           -6.8711e-01,  1.1629e-01,  1.9948e-01, -4.7043e-01,  1.2130e+00,
           -1.4237e-01,  2.6370e-01, -9.1200e-01, -2.1185e-01, -1.3493e+00,
           -2.0093e-02, -8.5553e-02,  7.4500e-02, -2.7861e-01,  5.7108e-01,
            6.0045e-01, -1.1580e+00,  2.2457e-01,  5.2492e-02,  7.7377e-02,
            2.1038e-01, -1.0968e+00, -1.8806e-01,  1.6193e-04,  3.0690e-01,
            1.5269e-01, -1.2186e+00,  4.4192e-01, -8.6209e-01, -3.1611e-01,
           -5.5617e-01,  2.4480e-02, -6.4080e-01, -8.4688e-01, -3.2271e-01,
            5.6128e-01,  6.9319e-01,  3.6391e-01,  1.4688e+00,  3.0880e-01,
            3.2401e-03,  5.1110e-02,  9.8396e-01, -5.8512e-01, -1.1058e-01,
            5.6782e-01,  5.1704e-01,  4.8607e-02,  3.8416e-02,  1.0723e-01,
           -5.8880e-01, -3.2970e-03, -4.8399e-01, -1.0776e+00,  3.2018e-01,
            8.1521e-01,  8.5122e-01,  5.0063e-01,  1.1447e+00, -1.8171e-01,
            4.7930e-01,  6.8499e-02, -9.1607e-01, -2.6628e-01, -1.0907e-01,
            1.1996e-01, -3.1764e-01, -3.9093e-01, -5.3363e-01,  2.1021e-01,
           -1.9328e-01, -4.6874e-01,  1.2553e-01, -5.8689e-01,  1.6500e-02,
            2.6675e-01,  4.8762e-03, -4.1349e-01, -1.0107e-01, -5.9281e-01,
           -3.6083e-01,  9.8077e-02,  4.3204e-01, -3.1348e-01,  2.6281e-01,
           -4.4020e-01,  9.3242e-02,  1.1335e+00, -3.1242e-01,  9.5865e-01,
            4.3642e-01, -9.1704e-01, -5.5271e-01,  5.3453e-01,  8.2077e-01,
           -3.4234e-02,  9.8388e-01,  5.4295e-01,  6.9819e-01,  1.0401e+00,
            5.3207e-01,  1.7786e-01,  1.1376e-01,  1.5848e-01, -2.3289e-01,
           -4.0887e-01, -4.7423e-01, -3.7506e-01,  3.0031e-02, -3.0034e-01,
           -3.0897e-01, -3.0131e-01, -2.1677e-01, -8.5815e-01, -3.3625e-01,
            1.7334e-01, -1.0664e+00,  2.9628e-01,  1.5440e-01, -8.2941e-02,
            1.9345e-01,  4.2606e-01,  3.3704e-01,  1.4883e+00, -1.4554e-01,
           -7.5085e-01,  5.5163e-01, -9.2985e-01,  1.2159e-01, -2.5228e-01,
            4.6588e-01, -1.2934e-02, -8.3803e-01,  3.8504e-02,  4.2223e-01,
            1.9769e-01, -4.9039e-02,  7.5554e-02,  7.1144e-01,  1.0876e-01,
            6.0351e-01,  3.1833e-01, -9.8136e-02,  8.2948e-02, -1.1635e-01,
            1.1131e-01,  1.5795e-01,  1.0441e+00, -7.2768e-02, -2.5971e-01,
           -4.1198e-01, -1.6680e+00, -7.8462e-01, -4.9429e-01, -2.1480e-02,
           -4.3348e-01, -3.0725e-01, -2.0169e-01, -3.1562e-01,  5.5264e-01,
            8.7739e-01, -2.5813e-01, -4.6994e-01,  2.2665e-01, -6.8825e-01,
           -1.0124e-01, -8.5311e-01, -4.8196e-01, -8.7491e-01,  8.3805e-01,
           -2.4941e-01, -4.3674e-02,  1.3940e-02, -7.0811e-01,  1.0864e+00,
            9.5805e-01, -1.4053e-01,  1.4219e-01,  5.0391e-01, -5.8840e-01,
           -6.1770e-01, -4.2367e-02,  2.9746e-01,  3.9472e-01,  6.5948e-01,
            6.4870e-01,  2.8621e-01, -5.4719e-01,  1.5666e+00,  7.9899e-01,
           -4.2446e-01,  1.0406e-01,  4.4579e-01,  3.1292e-01,  3.9760e-01,
```

## ⌄ Use BERT to Create Context-Aware Word Embeddings

```python
1  from transformers import BertTokenizer, BertModel
2  import torch
3  import torch.nn.functional as F
4
5  # Load pretrained BERT tokenizer and model
6  tokenizer = BertTokenizer.from_pretrained('bert-base-uncased')
7  model = BertModel.from_pretrained('bert-base-uncased')
8
9  # Function to extract contextual embedding for a word (handles subwords)
10 def get_token_embedding(sentence, target_word):
11     # Tokenize the sentence and get embeddings
12     inputs = tokenizer(sentence, return_tensors='pt')
13     outputs = model(**inputs)
14
15     # Get tokens and embeddings
16     tokens = tokenizer.convert_ids_to_tokens(inputs['input_ids'][0])
17     embeddings = outputs.last_hidden_state.squeeze(0)
18
19     # Tokenize the target word the same way BERT does
20     target_tokens = tokenizer.tokenize(target_word)
21
22     # Search for the position of the target word (handling subwords)
23     matches = []
```

```
24     for i in range(len(tokens) - len(target_tokens) + 1):
25         if tokens[i:i + len(target_tokens)] == target_tokens:
26             matches = list(range(i, i + len(target_tokens)))
27             break
28
29     if not matches:
30         raise ValueError(f"'{target_word}' not found in tokens: {tokens}")
31
32     # Average the embeddings over all subword tokens
33     return embeddings[matches].mean(dim=0)
34
35 # Contextual sentences
36 sentence_fruit = "He ate a fresh apple and enjoyed the fruit."
37 sentence_company = "Apple released a new product in the computer market."
38 sentence_orange = "An orange is a juicy fruit."
39 sentence_microsoft = "Microsoft computer was running the latest software."
40
41 # Get embeddings
42 apple_fruit = get_token_embedding(sentence_fruit, "apple")
43 apple_company = get_token_embedding(sentence_company, "apple")
44 orange = get_token_embedding(sentence_orange, "orange")
45 microsoft = get_token_embedding(sentence_microsoft, "Microsoft")
46
47 # Cosine similarity comparisons
48 sim_fruit = F.cosine_similarity(apple_fruit, orange, dim=0)
49 sim_company = F.cosine_similarity(apple_company, microsoft, dim=0)
50
51 # Results
52 print(f"Similarity between 'apple' (fruit) and 'orange': {sim_fruit.item():.4f}")
53 print(f"Similarity between 'apple' (company) and 'Microsoft': {sim_company.item():.4f}")
```

```
Similarity between 'apple' (fruit) and 'orange': 0.5839
Similarity between 'apple' (company) and 'Microsoft': 0.8549
```

## ⌄ Pipelines

Basic Pipeline Usage

### 1. Text Classification (Sentiment Analysis)

```
1 from transformers import pipeline
2
3 # Create a sentiment analysis pipeline
4 classifier = pipeline("sentiment-analysis")
5
6 # Analyze single text
7 result = classifier("I love using Hugging Face!")
8 print(result)
9 # Output: [{'label': 'POSITIVE', 'score': 0.9998}]
10
11 # Analyze multiple texts
12 texts = [
13     "I hate this product",
14     "This is amazing!",
15     "It's okay, nothing special"
16 ]
17 results = classifier(texts)
18 for text, result in zip(texts, results):
19     print(f"Text: {text}")
20     print(f"Sentiment: {result['label']}, Score: {result['score']:.4f}\n")
```

### 2. Named Entity Recognition (NER)

```
1 # NER pipeline
2 ner = pipeline("ner", aggregation_strategy="simple")
3
4 text = "My name is John and I live in New York. I work at Google."
5 entities = ner(text)
6
7 for entity in entities:
8     print(f"Entity: {entity['word']}")
9     print(f"Label: {entity['entity_group']}")
```

```
10      print(f"Score: {entity['score']:.4f}")
11      print(f"Start: {entity['start']}, End: {entity['end']}\n")
```

### 3. Question Answering

```
1   # Question answering pipeline
2   qa = pipeline("question-answering")
3
4   context = """
5   Hugging Face is a company that develops tools for building applications using machine learning.
6   They are especially known for their work in natural language processing. The company was founded in 2016
7   and is headquartered in New York.
8   """
9
10  questions = [
11      "When was Hugging Face founded?",
12      "Where is Hugging Face headquartered?",
13      "What is Hugging Face known for?"
14  ]
15
16  for question in questions:
17      result = qa(question=question, context=context)
18      print(f"Question: {question}")
19      print(f"Answer: {result['answer']}")
20      print(f"Score: {result['score']:.4f}\n")
```

### 4. Text Generation

```
1  # Text generation pipeline
2  generator = pipeline("text-generation", model="gpt2")
3
4  # Generate text with custom parameters
5  prompts = [
6      "The future of artificial intelligence is",
7      "In a world where robots exist,"
8  ]
9
10 for prompt in prompts:
11     generated = generator(
12         prompt,
13         max_length=50,
14         num_return_sequences=2,
15         temperature=0.7,
16         do_sample=True,
17         pad_token_id=generator.tokenizer.eos_token_id
18     )
19
20     print(f"Prompt: {prompt}")
21     for i, gen in enumerate(generated):
22         print(f"Generation {i+1}: {gen['generated_text']}\n")
```

### 5. Text Summarization

```
1  # Summarization pipeline
2  summarizer = pipeline("summarization")
3
4  article = """
5  Machine learning is a subset of artificial intelligence that enables computers to learn and improve
6  from experience without being explicitly programmed. It focuses on the development of computer programs
7  that can access data and use it to learn for themselves. The process of learning begins with observations
8  or data, such as examples, direct experience, or instruction, in order to look for patterns in data and
9  make better decisions in the future based on the examples that we provide. The primary aim is to allow
10 the computers to learn automatically without human intervention or assistance and adjust actions accordingly.
11 """
12
13 summary = summarizer(article, max_length=50, min_length=25, do_sample=False)
14 print("Original length:", len(article.split()))
15 print("Summary:", summary[0]['summary_text'])
16 print("Summary length:", len(summary[0]['summary_text'].split()))
```

### 6. Translation

```
1 # Translation pipeline
2 translator = pipeline("translation", model="Helsinki-NLP/opus-mt-en-fr")
3
4 texts = [
5     "Hello, how are you today?",
6     "Machine learning is fascinating.",
7     "I would like to order a coffee."
8 ]
9
10 for text in texts:
11     translated = translator(text)
12     print(f"English: {text}")
13     print(f"French: {translated[0]['translation_text']}\n")
```

## ⌄   Use specific model e.g. BERT to Create Questions Answering Pipeline

```
1 # Import required libraries
2 from transformers import AutoTokenizer, AutoModelForQuestionAnswering
3 from transformers import pipeline
4 import torch
5
6 # Using pipeline (High-level API)
7 qa_pipeline = pipeline( "question-answering",
8 model="bert-large-uncased-whole-word-masking-finetuned-squad",
9 tokenizer="bert-large-uncased-whole-word-masking-finetuned-squad" )
10
11 # Example usage
12 context = """ BERT is a method of pre-training language representations,
13 meaning that it trains a general-purpose language understanding
14 model on a large text corpus (like Wikipedia),
15 and then uses that model for downstream NLP tasks like question answering. """
16
17 question = "What is BERT?"
18 result = qa_pipeline(question=question, context=context)
19 print(f"Answer: {result['answer']}")
20 print(f"Confidence: {result['score']:.4f}")
21
```

```
c:\Users\me\myenv310\lib\site-packages\tqdm\auto.py:21: TqdmWarning: IProgress not found. Please update jupyter and ipywidgets. See http
  from .autonotebook import tqdm as notebook_tqdm
c:\Users\me\myenv310\lib\site-packages\huggingface_hub\file_download.py:144: UserWarning: `huggingface_hub` cache-system uses symlinks b
To support symlinks on Windows, you either need to activate Developer Mode or to run Python as an administrator. In order to activate de
  warnings.warn(message)
Xet Storage is enabled for this repo, but the 'hf_xet' package is not installed. Falling back to regular HTTP download. For better perfo
Some weights of the model checkpoint at bert-large-uncased-whole-word-masking-finetuned-squad were not used when initializing BertForQue
- This IS expected if you are initializing BertForQuestionAnswering from the checkpoint of a model trained on another task or with anoth
- This IS NOT expected if you are initializing BertForQuestionAnswering from the checkpoint of a model that you expect to be exactly ide
Device set to use cpu
Answer: a method of pre-training language representations
Confidence: 0.6874
```

## ⌄  Fine Tuning Large Language Models

**Tutorial: Fine-tuning a Language Model and Deploying with Hugging Face Spaces - Sentiment Analysis IMDB Reviews**

### ⌄   Step 1: Install Required Libraries

```
1 ! pip install transformers datasets huggingface_hub gradio
```

### ⌄   Step 2: Load and Prepare the Dataset

We will use a small portion of the IMDb dataset for binary sentiment classification.

```
1 from datasets import load_dataset
2
3 # Load a small subset for quicker training
4 dataset = load_dataset("imdb", split="train[:2000]")
5 dataset = dataset.train_test_split(test_size=0.2)
```

## ˅ Step 3: Load the Tokenizer and Model

We use `distilbert-base-uncased`, a lightweight version of BERT.

```
1 from transformers import AutoTokenizer, AutoModelForSequenceClassification
2
3 checkpoint = "distilbert-base-uncased"
4 tokenizer = AutoTokenizer.from_pretrained(checkpoint)
5 model = AutoModelForSequenceClassification.from_pretrained(checkpoint, num_labels=2)
```

## ˅ Step 4: Tokenize the Dataset

Tokenization prepares the text for input to the model.

```
1 def tokenize_function(example):
2     return tokenizer(example["text"], padding="max_length", truncation=True)
3
4 tokenized_dataset = dataset.map(tokenize_function, batched=True)
5 tokenized_dataset.set_format("torch", columns=["input_ids", "attention_mask", "label"])
```

## ˅ Step 5: Define Training Arguments and Trainer

```
 1 from transformers import TrainingArguments, Trainer
 2 import numpy as np
 3 from sklearn.metrics import accuracy_score, precision_recall_fscore_support
 4
 5 training_args = TrainingArguments(
 6     output_dir="./results",
 7     evaluation_strategy="epoch",
 8     num_train_epochs=3,
 9     per_device_train_batch_size=8,
10     per_device_eval_batch_size=8,
11     save_total_limit=1,
12     logging_dir="./logs",
13 )
14
15 def compute_metrics(eval_pred):
16     logits, labels = eval_pred
17     predictions = np.argmax(logits, axis=1)
18     precision, recall, f1, _ = precision_recall_fscore_support(labels, predictions, average='binary')
19     acc = accuracy_score(labels, predictions)
20     return {"accuracy": acc, "f1": f1, "precision": precision, "recall": recall}
21
22 trainer = Trainer(
23     model=model,
24     args=training_args,
25     train_dataset=tokenized_dataset["train"],
26     eval_dataset=tokenized_dataset["test"],
27     compute_metrics=compute_metrics,
28 )
```

## ˅ Step 6: Train the Model

```
1 trainer.train()
```

## ˅ Step 7: Log in to Hugging Face Hub

Do this only when you're ready to push your model.

```
1 from huggingface_hub import notebook_login
2
```

```
3 notebook_login()
```

After running this cell, you'll be prompted to enter your Hugging Face access token. You can create one here:
https://huggingface.co/settings/tokens

## ⌄ Step 8: Push Model and Tokenizer to Hugging Face Hub

Replace `"your-username/model-name"` with your actual username and desired model name.

```
1 model_name = "your-username/distilbert-sentiment-imdb-small"
2
3 model.push_to_hub(model_name)
4 tokenizer.push_to_hub(model_name)
```

This makes your model available for download and use in a Hugging Face Space.

## ⌄ Step 9: (Optional) Test with Gradio Locally in Colab

This is useful for debugging before deploying.

```
 1 import gradio as gr
 2 from transformers import pipeline
 3
 4 classifier = pipeline("sentiment-analysis", model=model_name)
 5
 6 def predict_sentiment(text):
 7     result = classifier(text)[0]
 8     return f"Label: {result['label']}, Confidence: {round(result['score'], 3)}"
 9
10 interface = gr.Interface(
11     fn=predict_sentiment,
12     inputs="text",
13     outputs="text",
14     title="Sentiment Analysis",
15     description="Enter a movie review to classify as POSITIVE or NEGATIVE."
16 )
17
18 interface.launch()
```

## ⌄ Step 10: Create a Hugging Face Space for Deployment

1. Go to https://huggingface.co/spaces

2. Click "Create New Space"

3. Choose:

   - **SDK**: Gradio
   - **Visibility**: Public or Private
   - Name: e.g. `sentiment-analyzer-student`

Add these two files to your Space:

1. `app.py`

```
 1 import gradio as gr
 2 from transformers import pipeline
 3
 4 model_name = "your-username/distilbert-sentiment-imdb-small"
 5 classifier = pipeline("sentiment-analysis", model=model_name)
 6
 7 def predict_sentiment(text):
 8     result = classifier(text)[0]
 9     return f"Label: {result['label']}, Confidence: {round(result['score'], 3)}"
10
11 interface = gr.Interface(
12     fn=predict_sentiment,
13     inputs="text",
14     outputs="text",
15     title="Sentiment Analysis",
```

```
16     description="Enter a movie review to classify as POSITIVE or NEGATIVE."
17 )
18
19 interface.launch()
```

2. requirements.txt

```
transformers
torch
gradio
```

After uploading both files, Hugging Face will automatically build and deploy your Space.

## ⌄ Conclusion

This complete workflow demonstrates how to:

- Fine-tune a transformer model on a small dataset
- Save and share the model using Hugging Face Hub
- Deploy the model as a web app with Hugging Face Spaces and Gradio

This structure is optimized for educational use, minimal setup, and reproducibility. If you would like a Colab version or a GitHub template, I can generate those for you as well.

---

## ⌄ Full Fine Tuning Code - IMDB Reviews

```
1 from datasets import load_dataset
2 from transformers import AutoTokenizer, AutoModelForSequenceClassification, TrainingArguments, Trainer
3 import numpy as np
4 from sklearn.metrics import accuracy_score, precision_recall_fscore_support
5
6 # Step 1: Load dataset
7 dataset = load_dataset("imdb", split='train[:2000]')
8 dataset = dataset.train_test_split(test_size=0.2)
9
10 # Step 2: Load tokenizer and model
11 tokenizer = AutoTokenizer.from_pretrained("distilbert-base-uncased")
12 model = AutoModelForSequenceClassification.from_pretrained("distilbert-base-uncased", num_labels=2)
13
14 # Step 3: Tokenize
15 def tokenize(example):
16     return tokenizer(example["text"], truncation=True, padding="max_length")
17
18 tokenized_data = dataset.map(tokenize, batched=True)
19 tokenized_data.set_format("torch", columns=["input_ids", "attention_mask", "label"])
20
21 # Step 4: Define training arguments
22 training_args = TrainingArguments(
23     output_dir="./results",
24     evaluation_strategy="epoch",
25     per_device_train_batch_size=8,
26     per_device_eval_batch_size=8,
27     num_train_epochs=3,
28     weight_decay=0.01,
29 )
30
31 # Step 5: Evaluation function
32 def compute_metrics(eval_pred):
33     logits, labels = eval_pred
34     preds = np.argmax(logits, axis=1)
35     precision, recall, f1, _ = precision_recall_fscore_support(labels, preds, average='binary')
36     acc = accuracy_score(labels, preds)
37     return {"accuracy": acc, "f1": f1, "precision": precision, "recall": recall}
38
39 # Step 6: Train
40 trainer = Trainer(
41     model=model,
42     args=training_args,
```

```
43     train_dataset=tokenized_data["train"],
44     eval_dataset=tokenized_data["test"],
45     compute_metrics=compute_metrics,
46 )
47
48 trainer.train()
49
```

---

## ⌄ Full Fine Tuning Code - Sentiment Analysis Amazon Reviews

```
1 # Required installations (uncomment if not already installed)
2 # !pip install transformers datasets scikit-learn
3
4 from transformers import BertTokenizer, BertForSequenceClassification, Trainer, TrainingArguments
5 from datasets import load_dataset, Dataset
6 from sklearn.model_selection import train_test_split
7 from sklearn.metrics import accuracy_score
8 import pandas as pd
9
10 # Load your CSV file (replace with your actual path)
11 df = pd.read_csv("amazon_reviews.csv")  # Columns: 'title', 'content', 'label'
12
13 # Combine title and content for input
14 df["text"] = df["title"] + " " + df["content"]
15 df = df[["text", "label"]]
16
17 # Split into train and validation
18 train_df, val_df = train_test_split(df, test_size=0.2, random_state=42)
19
20 # Convert to Hugging Face Dataset format
21 train_dataset = Dataset.from_pandas(train_df)
22 val_dataset = Dataset.from_pandas(val_df)
23
24 # Load pre-trained BERT tokenizer
25 tokenizer = BertTokenizer.from_pretrained("bert-base-uncased")
26
27 # Tokenize the text
28 def tokenize_function(example):
29     return tokenizer(example["text"], padding="max_length", truncation=True, max_length=256)
30
31 train_dataset = train_dataset.map(tokenize_function, batched=True)
32 val_dataset = val_dataset.map(tokenize_function, batched=True)
33
34 # Set format for PyTorch
35 train_dataset.set_format(type="torch", columns=["input_ids", "attention_mask", "label"])
36 val_dataset.set_format(type="torch", columns=["input_ids", "attention_mask", "label"])
37
38 # Load BERT model for binary classification
39 model = BertForSequenceClassification.from_pretrained("bert-base-uncased", num_labels=2)
40
41 # Define evaluation metrics
42 def compute_metrics(eval_pred):
43     logits, labels = eval_pred
44     preds = logits.argmax(axis=1)
45     return {"accuracy": accuracy_score(labels, preds)}
46
47 # Define training arguments
48 training_args = TrainingArguments(
49     output_dir="./results",
50     evaluation_strategy="epoch",
51     save_strategy="epoch",
52     per_device_train_batch_size=8,
53     per_device_eval_batch_size=8,
54     num_train_epochs=3,
55     logging_dir="./logs",
56     logging_steps=10,
57 )
58
59 # Trainer for training and evaluation
60 trainer = Trainer(
61     model=model,
62     args=training_args,
63     train_dataset=train_dataset,
```

```
64        eval_dataset=val_dataset,
65        compute_metrics=compute_metrics,
66 )
67
68 # Start fine-tuning
69 trainer.train()
70 # Save the model
71 trainer.save_model("fine_tuned_bert_amazon_reviews")
```

## ∨  More Fine Tuning Examples

### 1. Text Classification: News Topic Classification

Task: Classify news articles into topics (e.g., business, sports, politics)

- **Dataset**: AG News (4-class classification)
- **Model**: `distilbert-base-uncased`
- **Why it's good**: Multiclass instead of binary; introduces students to topic classification.

```
from datasets import load_dataset
from transformers import AutoTokenizer, AutoModelForSequenceClassification, TrainingArguments, Trainer
import numpy as np
from sklearn.metrics import accuracy_score, precision_recall_fscore_support

# Load data
dataset = load_dataset("ag_news")
dataset = dataset["train"].train_test_split(test_size=0.2)

# Tokenizer and model
tokenizer = AutoTokenizer.from_pretrained("distilbert-base-uncased")
model = AutoModelForSequenceClassification.from_pretrained("distilbert-base-uncased", num_labels=4)

def tokenize_function(example):
    return tokenizer(example["text"], padding="max_length", truncation=True)

tokenized_data = dataset.map(tokenize_function, batched=True)
tokenized_data.set_format("torch", columns=["input_ids", "attention_mask", "label"])

# Training arguments
training_args = TrainingArguments(
    output_dir="./results",
    evaluation_strategy="epoch",
    per_device_train_batch_size=8,
    num_train_epochs=3,
)

# Evaluation
def compute_metrics(eval_pred):
    logits, labels = eval_pred
    preds = np.argmax(logits, axis=1)
    acc = accuracy_score(labels, preds)
    return {"accuracy": acc}

trainer = Trainer(
    model=model,
    args=training_args,
    train_dataset=tokenized_data["train"],
    eval_dataset=tokenized_data["test"],
    compute_metrics=compute_metrics,
)

trainer.train()
```

Then follow the same login + push + deploy steps.

## 2. Text Generation: Simple Story Completion (using GPT-2)

Task: Given a prompt, generate the next few sentences of a story.

- **Model**: `gpt2`
- **Dataset**: A small set of fairy tales or a pre-tokenized open dataset like `wikitext`

> GPT-based fine-tuning takes longer and needs GPU memory, so keep the dataset very small.

```python
from datasets import load_dataset
from transformers import GPT2Tokenizer, GPT2LMHeadModel, Trainer, TrainingArguments

# Load small dataset
dataset = load_dataset("wikitext", "wikitext-2-raw-v1", split="train[:1%]")

tokenizer = GPT2Tokenizer.from_pretrained("gpt2")
tokenizer.pad_token = tokenizer.eos_token  # GPT2 has no padding token

def tokenize_function(examples):
    return tokenizer(examples["text"], return_special_tokens_mask=True, truncation=True, padding="max_length", max_length=128)

tokenized_dataset = dataset.map(tokenize_function, batched=True)
tokenized_dataset.set_format("torch", columns=["input_ids"])

# Load model
model = GPT2LMHeadModel.from_pretrained("gpt2")

training_args = TrainingArguments(
    output_dir="./gpt2-finetuned",
    overwrite_output_dir=True,
    per_device_train_batch_size=2,
    num_train_epochs=1,
    save_total_limit=1,
    prediction_loss_only=True,
)

trainer = Trainer(
    model=model,
    args=training_args,
    train_dataset=tokenized_dataset,
)

trainer.train()
```

You can then build a Gradio app with a text input (prompt) and text output (generated continuation).

## 3. Named Entity Recognition (NER)

Task: Identify entities like person names, locations, etc.

- **Dataset**: `conll2003`
- **Model**: `bert-base-cased`

NER gives students exposure to **token-level** classification.

```python
from datasets import load_dataset
from transformers import AutoTokenizer, AutoModelForTokenClassification, TrainingArguments, Trainer
from transformers import DataCollatorForTokenClassification
import numpy as np

dataset = load_dataset("conll2003")
tokenizer = AutoTokenizer.from_pretrained("bert-base-cased")

# Tokenize
```

```python
def tokenize_and_align_labels(examples):
    tokenized_inputs = tokenizer(examples["tokens"], truncation=True, is_split_into_words=True, padding="max_length")
    labels = []
    for i, label in enumerate(examples["ner_tags"]):
        word_ids = tokenized_inputs.word_ids(batch_index=i)
        previous_word_idx = None
        label_ids = []
        for word_idx in word_ids:
            if word_idx is None:
                label_ids.append(-100)
            elif word_idx != previous_word_idx:
                label_ids.append(label[word_idx])
            else:
                label_ids.append(label[word_idx])
            previous_word_idx = word_idx
        labels.append(label_ids)
    tokenized_inputs["labels"] = labels
    return tokenized_inputs

tokenized_data = dataset.map(tokenize_and_align_labels, batched=True)
tokenized_data.set_format("torch")

model = AutoModelForTokenClassification.from_pretrained("bert-base-cased", num_labels=9)

training_args = TrainingArguments(
    output_dir="./ner-model",
    evaluation_strategy="epoch",
    per_device_train_batch_size=8,
    num_train_epochs=2,
)

trainer = Trainer(
    model=model,
    args=training_args,
    train_dataset=tokenized_data["train"].select(range(1000)),
    eval_dataset=tokenized_data["validation"].select(range(200)),
    data_collator=DataCollatorForTokenClassification(tokenizer),
)

trainer.train()
```