

Phase 1: Data Preprocessing

```
1 pip install tensorflow nltk scikit-learn

→ Requirement already satisfied: tensorflow in /usr/local/lib/python3.11/dist-packages (2.18.0)
Requirement already satisfied: nltk in /usr/local/lib/python3.11/dist-packages (3.9.1)
Requirement already satisfied: scikit-learn in /usr/local/lib/python3.11/dist-packages (1.6.1)
Requirement already satisfied: absl-py>=1.0.0 in /usr/local/lib/python3.11/dist-packages (from tensorflow) (1.4.0)
Requirement already satisfied: astunparse>=1.6.0 in /usr/local/lib/python3.11/dist-packages (from tensorflow) (1.6.3)
Requirement already satisfied: flatbuffers>=24.3.25 in /usr/local/lib/python3.11/dist-packages (from tensorflow) (25.2.10)
Requirement already satisfied: gast!=0.5.0,!0.5.1,!0.5.2,>=0.2.1 in /usr/local/lib/python3.11/dist-packages (from tensorflow) (0.6.0)
Requirement already satisfied: google-pasta>=0.1.1 in /usr/local/lib/python3.11/dist-packages (from tensorflow) (0.2.0)
Requirement already satisfied: libclang>=13.0.0 in /usr/local/lib/python3.11/dist-packages (from tensorflow) (18.1.1)
Requirement already satisfied: opt-einsum>=2.3.2 in /usr/local/lib/python3.11/dist-packages (from tensorflow) (3.4.0)
Requirement already satisfied: packaging in /usr/local/lib/python3.11/dist-packages (from tensorflow) (24.2)
Requirement already satisfied: protobuf!=4.21.0,!4.21.1,!4.21.2,!4.21.3,!4.21.4,!4.21.5,<6.0.0dev,>=3.20.3 in /usr/local/lib/python
Requirement already satisfied: requests<3,>=2.21.0 in /usr/local/lib/python3.11/dist-packages (from tensorflow) (2.32.3)
Requirement already satisfied: setuptools in /usr/local/lib/python3.11/dist-packages (from tensorflow) (75.2.0)
Requirement already satisfied: six>=1.12.0 in /usr/local/lib/python3.11/dist-packages (from tensorflow) (1.17.0)
Requirement already satisfied: termcolor>=1.1.0 in /usr/local/lib/python3.11/dist-packages (from tensorflow) (3.0.1)
Requirement already satisfied: typing-extensions>=3.6.6 in /usr/local/lib/python3.11/dist-packages (from tensorflow) (4.13.1)
Requirement already satisfied: wrapt>=1.11.0 in /usr/local/lib/python3.11/dist-packages (from tensorflow) (1.17.2)
Requirement already satisfied: grpcio<2.0,>=1.24.3 in /usr/local/lib/python3.11/dist-packages (from tensorflow) (1.71.0)
Requirement already satisfied: tensorboard<2.19,>=2.18 in /usr/local/lib/python3.11/dist-packages (from tensorflow) (2.18.0)
Requirement already satisfied: keras>=3.5.0 in /usr/local/lib/python3.11/dist-packages (from tensorflow) (3.8.0)
Requirement already satisfied: numpy<2.1.0,>=1.26.0 in /usr/local/lib/python3.11/dist-packages (from tensorflow) (2.0.2)
Requirement already satisfied: h5py>=3.11.0 in /usr/local/lib/python3.11/dist-packages (from tensorflow) (3.13.0)
Requirement already satisfied: ml-dtypes<0.5.0,>=0.4.0 in /usr/local/lib/python3.11/dist-packages (from tensorflow) (0.4.1)
Requirement already satisfied: tensorflow-io-gcs-filesystem>=0.23.1 in /usr/local/lib/python3.11/dist-packages (from tensorflow) (0.37.1)
Requirement already satisfied: click in /usr/local/lib/python3.11/dist-packages (from nltk) (8.1.8)
Requirement already satisfied: joblib in /usr/local/lib/python3.11/dist-packages (from nltk) (1.4.2)
Requirement already satisfied: regex>=2021.8.3 in /usr/local/lib/python3.11/dist-packages (from nltk) (2024.11.6)
Requirement already satisfied: tqdm in /usr/local/lib/python3.11/dist-packages (from nltk) (4.67.1)
Requirement already satisfied: scipy>=1.6.0 in /usr/local/lib/python3.11/dist-packages (from scikit-learn) (1.14.1)
Requirement already satisfied: threadpoolctl>=3.1.0 in /usr/local/lib/python3.11/dist-packages (from scikit-learn) (3.6.0)
Requirement already satisfied: wheel<1.0,>=0.23.0 in /usr/local/lib/python3.11/dist-packages (from astunparse>=1.6.0->tensorflow) (0.45.
Requirement already satisfied: rich in /usr/local/lib/python3.11/dist-packages (from keras>=3.5.0->tensorflow) (13.9.4)
Requirement already satisfied: namex in /usr/local/lib/python3.11/dist-packages (from keras>=3.5.0->tensorflow) (0.0.8)
Requirement already satisfied: optree in /usr/local/lib/python3.11/dist-packages (from keras>=3.5.0->tensorflow) (0.14.1)
Requirement already satisfied: charset-normalizer<4,>=2 in /usr/local/lib/python3.11/dist-packages (from requests<3,>=2.21.0->tensorflow)
Requirement already satisfied: idna<4,>=2.5 in /usr/local/lib/python3.11/dist-packages (from requests<3,>=2.21.0->tensorflow) (3.10)
Requirement already satisfied: urllib3<3,>=1.21.1 in /usr/local/lib/python3.11/dist-packages (from requests<3,>=2.21.0->tensorflow) (2.3
Requirement already satisfied: certifi>=2017.4.17 in /usr/local/lib/python3.11/dist-packages (from requests<3,>=2.21.0->tensorflow) (202
Requirement already satisfied: markdown>=2.6.8 in /usr/local/lib/python3.11/dist-packages (from tensorboard<2.19,>=2.18->tensorflow) (3.
Requirement already satisfied: tensorboard-data-server<0.8.0,>=0.7.0 in /usr/local/lib/python3.11/dist-packages (from tensorboard<2.19,>
Requirement already satisfied: werkzeug>=1.0.1 in /usr/local/lib/python3.11/dist-packages (from tensorboard<2.19,>=2.18->tensorflow) (3.
Requirement already satisfied: MarkupSafe>=2.1.1 in /usr/local/lib/python3.11/dist-packages (from werkzeug>=1.0.1->tensorboard<2.19,>=2.
Requirement already satisfied: markdown-it-py>=2.2.0 in /usr/local/lib/python3.11/dist-packages (from rich->keras>=3.5.0->tensorflow) (3
Requirement already satisfied: pygments<3.0.0,>=2.13.0 in /usr/local/lib/python3.11/dist-packages (from rich->keras>=3.5.0->tensorflow)
Requirement already satisfied: mdurl~0.1 in /usr/local/lib/python3.11/dist-packages (from markdown-it-py>=2.2.0->rich->keras>=3.5.0->te
```

```
1 import pandas as pd
2
3 # Load datasets
4 insta_df = pd.read_csv("Instagram_sentiment.csv")
5 twitter_df = pd.read_csv("Twitter_sentiment.csv")
6
7 # Add source column
8 insta_df['platform'] = 'Instagram'
9 twitter_df['platform'] = 'Twitter'
10
11 # Combine datasets
12 combined_df = pd.concat([insta_df, twitter_df], ignore_index=True)
13
14 # Save to a new file
15 combined_df.to_csv("Combined_SocialMedia_Sentiment.csv", index=False)
```

Install necessary libraries

```
1 # -*- coding: utf-8 -*-
2 """
3 Fine-tuning_DistilBERT_Social_Media_Sentiment.ipynb
4
5 Automatically generated by Colaboratory.
6
7 Original file is located at
```

```
8     https://colab.research.google.com/notebooks/intro.ipynb
9
10 ## Fine-tuning a DistilBERT Model for Social Media Sentiment Analysis
11
12 This notebook demonstrates how to fine-tune a pre-trained DistilBERT model
13 from the Hugging Face library for sentiment analysis on a custom dataset
14 of social media posts (Instagram & Twitter).
15
16 **Dataset:** Combined_SocialMedia_Sentiment.csv
17 **Columns:** (See output of df.columns.tolist() below for actual columns)
18 **Goal:** Train a model to predict the sentiment score (1-5) based on the text.
19 """
20
21 # @title Install necessary libraries
22 # Install transformers, datasets (for handling data), evaluate (for metrics)
23 !pip install transformers datasets evaluate accelerate -q # -q for quiet installation
24
25 # @title Import Libraries
26 import pandas as pd
27 import numpy as np
28 import torch
29 from sklearn.model_selection import train_test_split
30 # Corrected import: Removed load_metric as it's deprecated
31 from datasets import Dataset, DatasetDict
32 import evaluate # Preferred way for metrics
33 from transformers import (
34     AutoTokenizer,
35     AutoModelForSequenceClassification,
36     TrainingArguments,
37     Trainer,
38     DataCollatorWithPadding,
39 )
40 import re # For basic text cleaning
41 import warnings
42
43 # Ignore warnings for cleaner output
44 warnings.filterwarnings("ignore")
45
46 # @title Configuration
47 # --- Parameters You Might Want to Tune ---
48 CSV_FILE_PATH = "/content/Combined_SocialMedia_Sentiment.csv" # <<< IMPORTANT: Upload your CSV file to Colab or change path
49 # --- VITAL: Updated based on user's df.columns output ---
50 TEXT_COLUMN = "comment"
51 # --- VITAL: Updated based on user's df.columns output ---
52 LABEL_COLUMN = "sentiment"
53 MODEL_CHECKPOINT = "distilbert-base-uncased" # A good balance of speed and performance
54 NUM_LABELS = 5 # We have 5 sentiment scores (1 to 5)
55 TEST_SIZE = 0.2 # Proportion of data to use for validation/testing
56 RANDOM_STATE = 42 # For reproducibility
57 BATCH_SIZE = 16 # Adjust based on GPU memory in Colab
58 LEARNING_RATE = 2e-5
59 NUM_EPOCHS = 3 # Number of training epochs
60 OUTPUT_DIR = "sentiment_model_finetuned" # Directory to save the model
61
62 # @title Load and Prepare Data
63 # Load the dataset
64 try:
65     df = pd.read_csv(CSV_FILE_PATH)
66     print(f"Successfully loaded {len(df)} rows from {CSV_FILE_PATH}")
67     # --- DEBUGGING: Print the actual columns found in the CSV ---
68     print("DataFrame Columns:", df.columns.tolist())
69     # --- Check this output and update TEXT_COLUMN/LABEL_COLUMN above if needed ---
70
71 except FileNotFoundError:
72     print(f"ERROR: File not found at {CSV_FILE_PATH}.")
73     print("Please upload the file to your Colab environment or update the path.")
74     # Exit or raise error if file not found
75     raise SystemExit("CSV file not found.")
76 except Exception as e:
77     print(f"An error occurred while loading the CSV: {e}")
78     raise SystemExit("Error loading CSV.")
79
80
81 # --- Data Cleaning & Preprocessing ---
82 print("\n--- Data Preprocessing ---")
83
84 # Handle missing text data
```

```

85 print(f"Initial rows: {len(df)}")
86 # Check if TEXT_COLUMN exists before trying to dropna
87 if TEXT_COLUMN not in df.columns:
88     print(f"\nERROR: The column specified by TEXT_COLUMN ('{TEXT_COLUMN}') was not found in the DataFrame.")
89     print(f"Available columns are: {df.columns.tolist()}")
90     print("Please update the TEXT_COLUMN variable in the 'Configuration' cell and rerun.")
91     raise KeyError(f"Column '{TEXT_COLUMN}' not found in DataFrame.")
92
93 # Handle potential non-numeric or NaN values in LABEL_COLUMN before conversion
94 if LABEL_COLUMN not in df.columns:
95     print(f"\nERROR: The column specified by LABEL_COLUMN ('{LABEL_COLUMN}') was not found in the DataFrame.")
96     print(f"Available columns are: {df.columns.tolist()}")
97     print("Please update the LABEL_COLUMN variable in the 'Configuration' cell and rerun.")
98     raise KeyError(f"Column '{LABEL_COLUMN}' not found in DataFrame.")
99
100 # Convert label column to numeric, coercing errors to NaN, then drop rows with NaN labels
101 df[LABEL_COLUMN] = pd.to_numeric(df[LABEL_COLUMN], errors='coerce')
102 initial_rows_before_label_dropna = len(df)
103 df = df.dropna(subset=[LABEL_COLUMN])
104 print(f"Rows after dropping NaN in label column '{LABEL_COLUMN}': {len(df)} (dropped {initial_rows_before_label_dropna - len(df)})")
105
106 # Drop rows with NaN in the text column
107 initial_rows_before_text_dropna = len(df)
108 df = df.dropna(subset=[TEXT_COLUMN])
109 print(f"Rows after dropping missing text in column '{TEXT_COLUMN}': {len(df)} (dropped {initial_rows_before_text_dropna - len(df)})")
110
111
112 # Basic text cleaning (optional, as BERT tokenizers handle much of this)
113 def clean_text(text):
114     if not isinstance(text, str):
115         # Attempt to convert to string, or return empty if fails
116         try:
117             text = str(text)
118         except:
119             return ""
120     text = re.sub(r"http\S+", "", text) # Remove URLs
121     text = re.sub(r"@[\w-]+", "", text) # Remove mentions
122     text = re.sub(r"#[\w-]+", "", text) # Remove hashtags (optional, might contain sentiment)
123     text = re.sub(r"\s+", " ", text).strip() # Remove extra whitespace
124     return text
125
126 df['cleaned_text'] = df[TEXT_COLUMN].apply(clean_text)
127
128 # Prepare labels: sentiment score is 1-5, models need 0-based labels (0-4)
129 # Ensure the label column is integer type before subtracting
130 df['labels'] = df[LABEL_COLUMN].astype(int) - 1
131 print(f"Label distribution (0-4):\n{df['labels'].value_counts().sort_index()}") # Sort index for clarity
132
133 # Check if all expected labels (0-4) are present after conversion
134 expected_labels = set(range(NUM_LABELS))
135 actual_labels = set(df['labels'].unique())
136 if not expected_labels.issubset(actual_labels):
137     print(f"\nWARNING: Some expected labels (0-{NUM_LABELS}-1) might be missing after processing.")
138     print(f"Expected: {expected_labels}, Found: {actual_labels}")
139     # Consider adjusting NUM_LABELS if the data truly doesn't have all 5 scores,
140     # but usually it's better to keep it as 5 if the original scale was 1-5.
141
142 # Keep only necessary columns
143 df_processed = df[['cleaned_text', 'labels']].rename(columns={'cleaned_text': 'text'})
144
145 # --- Split Data ---
146 print("\n--- Splitting Data ---")
147 # Check if there's enough data to stratify
148 if df_processed['labels'].nunique() > 1 and len(df_processed) > 10: # Basic check
149     try:
150         train_df, test_df = train_test_split(
151             df_processed,
152             test_size=TEST_SIZE,
153             random_state=RANDOM_STATE,
154             stratify=df_processed['labels'] # Ensure similar label distribution in splits
155         )
156     except ValueError as e:
157         print(f"Stratify failed: {e}. Splitting without stratify.")
158         train_df, test_df = train_test_split(
159             df_processed,
160             test_size=TEST_SIZE,
161             random_state=RANDOM_STATE

```

```

162         )
163 elif len(df_processed) > 0:
164     print("Not enough data or classes to stratify. Splitting without stratify.")
165     train_df, test_df = train_test_split(
166         df_processed,
167         test_size=TEST_SIZE,
168         random_state=RANDOM_STATE
169     )
170 else:
171     raise SystemExit("No data left after preprocessing to split.")
172
173
174 print(f"Training set size: {len(train_df)}")
175 print(f"Test set size: {len(test_df)}")
176
177 # Convert pandas DataFrames to Hugging Face Datasets
178 train_dataset = Dataset.from_pandas(train_df.reset_index(drop=True))
179 test_dataset = Dataset.from_pandas(test_df.reset_index(drop=True))
180
181 # Combine into a DatasetDict (standard format for Hugging Face)
182 dataset_dict = DatasetDict({
183     'train': train_dataset,
184     'test': test_dataset
185 })
186
187 print("\nDatasetDict created:")
188 print(dataset_dict)
189
190 # @title Tokenization
191 print("\n--- Tokenizing Data ---")
192 # Load tokenizer associated with the chosen model
193 tokenizer = AutoTokenizer.from_pretrained(MODEL_CHECKPOINT)
194
195 # Function to tokenize the text data
196 def tokenize_function(examples):
197     # `truncation=True` ensures sequences longer than model max length are cut.
198     # `padding=True` would pad here, but DataCollator is usually preferred for dynamic padding.
199     # Handle potential non-string data in the 'text' column just in case
200     texts = [str(t) if t is not None else "" for t in examples["text"]]
201     return tokenizer(texts, truncation=True, max_length=512) # Using max_length=512 for BERT-like models
202
203 # Apply tokenization to the entire dataset_dict
204 # `batched=True` processes multiple examples at once for speed
205 tokenized_datasets = dataset_dict.map(tokenize_function, batched=True)
206
207 # Remove the original text column as it's no longer needed after tokenization
208 # Also remove '__index_level_0__' if it exists from pandas conversion
209 columns_to_remove = ["text"]
210 if '__index_level_0__' in tokenized_datasets['train'].column_names:
211     columns_to_remove.append('__index_level_0__')
212
213 # Check if columns actually exist before removing
214 existing_columns_to_remove = [col for col in columns_to_remove if col in tokenized_datasets['train'].column_names]
215 if existing_columns_to_remove:
216     tokenized_datasets = tokenized_datasets.remove_columns(existing_columns_to_remove)
217
218
219 # Rename 'labels' column to 'label' if necessary (Trainer expects 'label')
220 if 'labels' in tokenized_datasets['train'].column_names and 'label' not in tokenized_datasets['train'].column_names:
221     tokenized_datasets = tokenized_datasets.rename_column("labels", "label")
222
223
224 # Set format to PyTorch tensors
225 tokenized_datasets.set_format("torch")
226
227 print("\nTokenized datasets prepared:")
228 print(tokenized_datasets)
229 # Print example only if dataset is not empty
230 if len(tokenized_datasets['train']) > 0:
231     print(f"\nExample tokenized input: {tokenized_datasets['train'][0]}")
232 else:
233     print("\nTraining dataset is empty after processing.")
234
235
236 # Data collator handles dynamic padding (pads sequences to the longest in a batch)
237 data_collator = DataCollatorWithPadding(tokenizer=tokenizer)
238 print("\nData Collator initialized.")

```

```

239
240 # @title Load Model and Define Metrics
241 print("\n--- Loading Model & Defining Metrics ---")
242
243 # Load the pre-trained model for sequence classification
244 # `num_labels` must match the number of unique sentiment scores (0-4)
245 model = AutoModelForSequenceClassification.from_pretrained(
246     MODEL_CHECKPOINT,
247     num_labels=NUM_LABELS,
248     # Add id2label and label2id for better inference later (optional but good practice)
249     id2label={i: f"LABEL_{i+1}" for i in range(NUM_LABELS)}, # Maps 0 -> LABEL_1, 1 -> LABEL_2 etc.
250     label2id={f"LABEL_{i+1}": i for i in range(NUM_LABELS)}
251 )
252 print(f"Loaded model: {MODEL_CHECKPOINT} with {NUM_LABELS} labels.")
253 # print(f"Model label mapping: {model.config.id2label}") # Uncomment to check mapping
254
255 # Define evaluation metrics using the 'evaluate' library
256 accuracy_metric = evaluate.load("accuracy")
257 f1_metric = evaluate.load("f1")
258
259 def compute_metrics(eval_pred):
260     """Computes accuracy and F1 score for evaluation predictions."""
261     logits, labels = eval_pred
262     # Get the predicted class by finding the index with the highest logit value
263     predictions = np.argmax(logits, axis=-1)
264
265     # Calculate accuracy
266     accuracy = accuracy_metric.compute(predictions=predictions, references=labels)
267
268     # Calculate F1 score (use 'weighted' for multi-class imbalance)
269     f1 = f1_metric.compute(predictions=predictions, references=labels, average="weighted")
270
271     return {
272         "accuracy": accuracy["accuracy"],
273         "f1": f1["f1"],
274     }
275
276 print("Metrics computation function defined.")
277
278 # @title Configure Training Arguments
279 print("\n--- Configuring Training ---")
280
281 # Define training arguments
282 # These control various aspects of the fine-tuning process
283 training_args = TrainingArguments(
284     output_dir=OUTPUT_DIR, # Directory to save model checkpoints and logs
285     num_train_epochs=NUM_EPOCHS, # Total number of training epochs
286     per_device_train_batch_size=BATCH_SIZE, # Batch size per GPU/CPU for training
287     per_device_eval_batch_size=BATCH_SIZE, # Batch size per GPU/CPU for evaluation
288     learning_rate=LEARNING_RATE, # Learning rate for the optimizer
289     weight_decay=0.01, # Strength of weight decay regularization
290     evaluation_strategy="epoch", # Evaluate model performance at the end of each epoch
291     save_strategy="epoch", # Save model checkpoint at the end of each epoch
292     logging_strategy="epoch", # Log training metrics at the end of each epoch
293     load_best_model_at_end=True, # Load the best model checkpoint (based on metric) at the end
294     metric_for_best_model="f1", # Metric used to identify the best model (use F1 for potentially imbalanced classes)
295     save_total_limit=2, # Only keep the best and the latest checkpoint
296     push_to_hub=False, # Set to True to push model to Hugging Face Hub (requires login)
297     report_to="none", # Disable reporting to integrations like wandb/tensorboard unless configured
298 )
299
300 print("Training Arguments configured:")
301 # print(training_args) # Can be verbose, print if needed
302
303 # @title Initialize Trainer
304 print("\n--- Initializing Trainer ---")
305
306 # Check if train dataset is empty before initializing Trainer
307 if len(tokenized_datasets['train']) == 0:
308     raise SystemExit("Training dataset is empty. Cannot initialize Trainer. Please check data loading and preprocessing steps.")
309
310
311 # Initialize the Trainer
312 # This class orchestrates the fine-tuning process
313 trainer = Trainer(
314     model=model, # The instantiated Transformers model to be trained
315     args=training_args, # Training arguments defined above

```

```

316     train_dataset=tokennized_datasets["train"], # Training dataset
317     eval_dataset=tokennized_datasets["test"], # Evaluation dataset (using test set here for simplicity, often a separate validation set
318     tokenizer=tokenizer,                      # Tokenizer for processing data (needed by Trainer)
319     data_collator=data_collator,             # Data collator to handle batch padding
320     compute_metrics=compute_metrics,         # Function to compute evaluation metrics
321 )
322
323 print("Trainer initialized.")
324
325 # @title Start Fine-tuning
326 print("\n--- Starting Fine-tuning ---")
327 # Start the training process
328 try:
329     train_result = trainer.train()
330     print("\n--- Fine-tuning Completed ---")
331
332     # Save training metrics
333     metrics = train_result.metrics
334     trainer.log_metrics("train", metrics)
335     trainer.save_metrics("train", metrics)
336
337     # Save the fine-tuned model and tokenizer
338     trainer.save_model(OUTPUT_DIR) # Saves the best model due to load_best_model_at_end=True
339     tokenizer.save_pretrained(OUTPUT_DIR)
340     print(f"\nBest fine-tuned model and tokenizer saved to {OUTPUT_DIR}")
341
342 except Exception as e:
343     print(f"\nAn error occurred during training: {e}")
344     # Potentially add more specific error handling (e.g., CUDA out of memory)
345     raise SystemExit("Training failed.")
346
347
348 # @title Evaluate the Fine-tuned Model
349 print("\n--- Evaluating Model on Test Set ---")
350 # Evaluate the performance on the test set (loads the best model automatically)
351 eval_results = trainer.evaluate()
352
353 print("\n--- Evaluation Results ---")
354 print(f"Accuracy: {eval_results.get('eval_accuracy'):.4f}")
355 print(f"F1 Score (Weighted): {eval_results.get('eval_f1'):.4f}")
356 # print(eval_results) # Print full results if needed
357
358 # Save evaluation metrics
359 trainer.log_metrics("eval", eval_results)
360 trainer.save_metrics("eval", eval_results)
361
362 # @title Example Prediction (Optional)
363 print("\n--- Example Prediction ---")
364
365 # Load the fine-tuned model explicitly (not strictly necessary as trainer holds the best model)
366 # model_fine_tuned = AutoModelForSequenceClassification.from_pretrained(OUTPUT_DIR)
367 # tokenizer_fine_tuned = AutoTokenizer.from_pretrained(OUTPUT_DIR)
368
369 # Example text
370 text_to_classify = "This movie was absolutely fantastic, loved the acting!"
371 # text_to_classify = "The service was slow and the food was cold."
372 # text_to_classify = "It was an okay experience, nothing special."
373
374 # Use the trainer's model and tokenizer for prediction
375 current_model = trainer.model
376 current_tokenizer = trainer.tokenizer
377
378 # Tokenize the example text
379 inputs = current_tokenizer(text_to_classify, return_tensors="pt", truncation=True, padding=True, max_length=512)
380
381 # Move inputs to the same device as the model (important if using GPU)
382 if torch.cuda.is_available() and current_model.device.type == 'cuda':
383     inputs = {k: v.to(current_model.device) for k, v in inputs.items()}
384
385 # Perform prediction
386 with torch.no_grad(): # Disable gradient calculation for inference
387     logits = current_model(**inputs).logits
388
389 # Get the predicted class index
390 predicted_class_id = torch.argmax(logits, dim=-1).item()
391
392 # Map back to original sentiment score (0-4 -> 1-5)

```

```
393 predicted_sentiment_score = predicted_class_id + 1
394
395 print(f"Text: '{text_to_classify}'")
396 print(f"Predicted Sentiment Score (1-5): {predicted_sentiment_score}")
397
398 # You can add a mapping for better interpretation:
399 sentiment_map = {1: "Very Negative", 2: "Negative", 3: "Neutral", 4: "Positive", 5: "Very Positive"}
400 print(f"Predicted Sentiment Label: {sentiment_map.get(predicted_sentiment_score, 'Unknown')}")
401
402 # Example 2
403 text_to_classify_2 = "The service was slow and the food was cold."
404 inputs_2 = current_tokenizer(text_to_classify_2, return_tensors="pt", truncation=True, padding=True, max_length=512)
405 if torch.cuda.is_available() and current_model.device.type == 'cuda':
406     inputs_2 = {k: v.to(current_model.device) for k, v in inputs_2.items()}
407 with torch.no_grad():
408     logits_2 = current_model(**inputs_2).logits
409 predicted_class_id_2 = torch.argmax(logits_2, dim=-1).item()
410 predicted_sentiment_score_2 = predicted_class_id_2 + 1
411 print(f"\nText: '{text_to_classify_2}'")
412 print(f"Predicted Sentiment Score (1-5): {predicted_sentiment_score_2}")
413 print(f"Predicted Sentiment Label: {sentiment_map.get(predicted_sentiment_score_2, 'Unknown')}")
414
415
```

```

SuccessFully loaded 2000 rows from /content/Combined_SocialMedia_Sentiment.csv
DataFrame Columns: ['url', 'comment_user', 'comment_user_url', 'comment_date', 'comment', 'likes_number', 'replies_number', 'replies', '--- Data Preprocessing ---'
Initial rows: 2000
Rows after dropping NaN in label column 'sentiment': 2000 (dropped 0)
Rows after dropping missing text in column 'comment': 1000 (dropped 1000)
Label distribution (0-4):
labels
0    129
1     19
2     57
3     39
4    756
Name: count, dtype: int64

--- Splitting Data ---
Training set size: 800
Test set size: 200

DatasetDict created:
DatasetDict({
  train: Dataset({
    features: ['text', 'labels'],
    num_rows: 800
  })
  test: Dataset({
    features: ['text', 'labels'],
    num_rows: 200
  })
})

--- Tokenizing Data ---
tokenizer_config.json: 100%                                         48.0/48.0 [00:00<00:00, 2.95kB/s]
config.json: 100%                                         483/483 [00:00<00:00, 42.7kB/s]
vocab.txt: 100%                                         232k/232k [00:00<00:00, 5.60MB/s]
tokenizer.json: 100%                                         466k/466k [00:00<00:00, 3.22MB/s]
Map: 100%                                         800/800 [00:00<00:00, 5703.76 examples/s]
Map: 100%                                         200/200 [00:00<00:00, 2978.51 examples/s]

Tokenized datasets prepared:
DatasetDict({
  train: Dataset({
    features: ['label', 'input_ids', 'attention_mask'],
    num_rows: 800
  })
  test: Dataset({
    features: ['label', 'input_ids', 'attention_mask'],
    num_rows: 200
  })
})

Example tokenized input: {'label': tensor(4), 'input_ids': tensor([ 101, 11320, 6843, 11409, 3527, 102]), 'attention_mask': tensc
Data Collator initialized.

--- Loading Model & Defining Metrics ---
Xet Storage is enabled for this repo, but the 'hf_xet' package is not installed. Falling back to regular HTTP download. For better per
WARNING:huggingface_hub.file_download:Xet Storage is enabled for this repo, but the 'hf_xet' package is not installed. Falling back to
model.safetensors: 100%                                         268M/268M [00:02<00:00, 165MB/s]
Some weights of DistilBertForSequenceClassification were not initialized from the model checkpoint at distilbert-base-uncased and are
You should probably TRAIN this model on a down-stream task to be able to use it for predictions and inference.
Loaded model: distilbert-base-uncased with 5 labels.

Downloading builder script: 100%                                         4.20k/4.20k [00:00<00:00, 444kB/s]
Downloading builder script: 100%                                         6.79k/6.79k [00:00<00:00, 819kB/s]
Metrics computation function defined.

--- Configuring Training ---
Training Arguments configured:

--- Initializing Trainer ---
Trainer initialized.

--- Starting Fine-tuning ---
[150/150 11:00, Epoch 3/3]
Epoch Training Loss Validation Loss Accuracy F1

```

1	0.923500	0.809751	0.755000	0.649601
2	0.734600	0.777292	0.755000	0.649601
3	0.679500	0.759023	0.755000	0.649601

--- Fine-tuning Completed ---

***** train metrics *****

```
epoch          =      3.0
total_flos     =    40895GF
train_loss      =      0.7792
train_runtime   = 0:11:43.29
train_samples_per_second = 3.413
train_steps_per_second = 0.213
```

Best fine-tuned model and tokenizer saved to sentiment_model_finetuned

--- Evaluating Model on Test Set ---

[13/13 00:09]

Trainer.tokenizer is now deprecated. You should use Trainer.processing_class instead.

--- Evaluation Results ---

Accuracy: 0.7550
F1 Score (Weighted): 0.6496

***** eval metrics *****

```
epoch          =      3.0
eval_accuracy  =      0.755
eval_f1        =      0.6496
eval_loss       =      0.8098
eval_runtime   = 0:00:09.51
eval_samples_per_second = 21.03
eval_steps_per_second = 1.367
```

--- Example Prediction ---

Text: 'This movie was absolutely fantastic, loved the acting!'

Predicted Sentiment Score (1-5): 5

Predicted Sentiment Label: Very Positive

Text: 'The service was slow and the food was cold.'

Predicted Sentiment Score (1-5): 5

Predicted Sentiment Label: Very Positive

✓ Take 2

✓ Cell 1: Install Libraries

```
1 ## Fine-tuning DistilBERT on Amazon Food Reviews
2
3 # @title Cell 1: Install Libraries
4 !pip install transformers datasets evaluate accelerate torch -q # -q for quiet installation
5 print("Libraries installed.")
```

491.2/491.2 kB 17.1 MB/s eta 0:00:00
 84.0/84.0 kB 6.5 MB/s eta 0:00:00
 363.4/363.4 kB 3.4 MB/s eta 0:00:00
 13.8/13.8 kB 36.2 MB/s eta 0:00:00
 24.6/24.6 kB 39.6 MB/s eta 0:00:00
 883.7/883.7 kB 30.3 MB/s eta 0:00:00
 664.8/664.8 kB 2.6 MB/s eta 0:00:00
 211.5/211.5 kB 4.7 MB/s eta 0:00:00
 56.3/56.3 kB 11.5 MB/s eta 0:00:00
 127.9/127.9 kB 7.9 MB/s eta 0:00:00
 207.5/207.5 kB 4.6 MB/s eta 0:00:00
 21.1/21.1 kB 41.4 MB/s eta 0:00:00
 116.3/116.3 kB 5.9 MB/s eta 0:00:00
 183.9/183.9 kB 8.6 MB/s eta 0:00:00
 143.5/143.5 kB 9.3 MB/s eta 0:00:00
 194.8/194.8 kB 11.5 MB/s eta 0:00:00

ERROR: pip's dependency resolver does not currently take into account all the packages that are installed. This behaviour is the source of gcsfs 2025.3.2 requires fsspec==2025.3.2, but you have fsspec 2024.12.0 which is incompatible.
 Libraries installed.

✓ Cell 2: Import Libraries

```
1 # @title Cell 2: Import Libraries
2 import pandas as pd
3 import numpy as np
4 import torch
5 from sklearn.model_selection import train_test_split
6 from datasets import Dataset, DatasetDict
7 import evaluate # Preferred way for metrics
8 from transformers import (
9     AutoTokenizer,
10    AutoModelForSequenceClassification,
11    TrainingArguments,
12    Trainer,
13    DataCollatorWithPadding,
14 )
15 import re # For basic text cleaning
16 import warnings
17
18 # Ignore warnings for cleaner output
19 warnings.filterwarnings("ignore")
20 print("Libraries imported.")
```

↳ Libraries imported.

✓ Cell 3: Configuration

```
1 # title Cell 3: Configuration
2
3 # --- Data Settings ---
4 CSV_FILE_PATH = "//content/Reviews.csv"
5 SUBSET_SIZE = 50000          # <<< Set to None to load all data (might cause memory issues!), or set a number (e.g., 50000)
6 TEXT_COLUMN = "Text"
7 LABEL_COLUMN = "Score"
8
9 # --- Model Settings ---
10 MODEL_CHECKPOINT = "distilbert-base-uncased" # A good balance of speed and performance
11 NUM_LABELS = 5                      # We have 5 sentiment scores (1 to 5)
12
13 # --- Training Settings ---
```

```

14 TEST_SIZE = 0.2          # Proportion of data to use for validation/testing
15 RANDOM_STATE = 42        # For reproducibility
16 BATCH_SIZE = 16          # Adjust based on GPU memory in Colab (16 or 32 often work)
17 LEARNING_RATE = 2e-5     # Standard learning rate for fine-tuning BERT-like models
18 NUM_EPOCHS = 3            # Start with 1 epoch for large datasets
19 OUTPUT_DIR = "sentiment_model_amazon_csv_finetuned" # Directory to save the model
20
21 print("Configuration set.")
22 print(f"CSV file: {CSV_FILE_PATH}")
23 print(f"Loading subset size: {'All' if SUBSET_SIZE is None else SUBSET_SIZE}")
24 print(f"Text column: {TEXT_COLUMN}")
25 print(f"Label column: {LABEL_COLUMN}")
26 print(f"Model: {MODEL_CHECKPOINT}")

```

→ Configuration set.
 CSV file: //content/Reviews.csv
 Loading subset size: 50000
 Text column: Text
 Label column: Score
 Model: distilbert-base-uncased

Cell 4: Load Data from CSV

```

1 # title Cell 4: Load Data from CSV
2
3 try:
4     print(f"Loading data from CSV: {CSV_FILE_PATH}...")
5     # Use nrows to load only a subset if SUBSET_SIZE is specified
6     df = pd.read_csv(CSV_FILE_PATH, nrows=SUBSET_SIZE)
7     print(f"Successfully loaded {len(df)} rows.")
8     print("Sample rows:")
9     print(df.head())
10
11    # --- DEBUGGING: Print the actual columns found in the CSV ---
12    print("\nDataFrame Columns:", df.columns.tolist())
13    # --- Check this output and update TEXT_COLUMN/LABEL_COLUMN in Cell 3 if needed ---
14
15    print("\nColumn data types:")
16    print(df.info())
17
18    # Verify required columns exist
19    if TEXT_COLUMN not in df.columns:
20        raise ValueError(f"Text column '{TEXT_COLUMN}' not found in CSV.")
21    if LABEL_COLUMN not in df.columns:
22        raise ValueError(f"Label column '{LABEL_COLUMN}' not found in CSV.")
23
24 except FileNotFoundError:
25     print(f"\nERROR: File not found at {CSV_FILE_PATH}.")
26     print("Please upload the CSV file to your Colab environment or update the path in Cell 3.")
27     raise SystemExit("CSV file not found.")
28 except ValueError as ve:
29     print(f"\nERROR: {ve}")
30     print(f"Available columns are: {df.columns.tolist() if 'df' in locals() else 'Could not load DataFrame'}")
31     print("Please check the TEXT_COLUMN and LABEL_COLUMN variables in Cell 3.")
32     raise SystemExit("Required column not found.")
33 except Exception as e:
34     print(f"\nAn error occurred during CSV loading: {e}")
35     raise SystemExit("CSV loading failed.")

```

→ Loading data from CSV: //content/Reviews.csv...
 Successfully loaded 50000 rows.
 Sample rows:

	Id	ProductId	UserId	ProfileName	\
0	1	B001E4KFG0	A3SGXH7AUHU8GW	delmartian	
1	2	B00813GRG4	A1D87F6ZCVE5NK	dll pa	
2	3	B000LQOCHO	ABXLMWJIXXAIN	Natalia Corres	"Natalia Corres"
3	4	B000UA0QIQ	A395BORC6FGVXV	Karl	
4	5	B006K2ZZ7K	A1UQRSCLF8GW1T	Michael D. Bigham	"M. Wassir"

	HelpfulnessNumerator	HelpfulnessDenominator	Score	Time	\
0	1	1	5	1303862400	
1	0	0	1	1346976000	
2	1	1	4	1219017600	
3	3	3	2	1307923200	
4	0	0	5	1350777600	

Summary	Text
0 Good Quality Dog Food	I have bought several of the Vitality canned d...
1 Not as Advertised	Product arrived labeled as Jumbo Salted Peanut...
2 "Delight" says it all	This is a confection that has been around a fe...
3 Cough Medicine	If you are looking for the secret ingredient i...
4 Great taffy	Great taffy at a great price. There was a wid...

DataFrame Columns: ['Id', 'ProductId', 'UserId', 'ProfileName', 'HelpfulnessNumerator', 'HelpfulnessDenominator', 'Score', 'Time', 'Summary', 'Text']

Column data types:

```
<class 'pandas.core.frame.DataFrame'>
RangeIndex: 50000 entries, 0 to 49999
Data columns (total 10 columns):
 #   Column           Non-Null Count  Dtype  
--- 
 0   Id               50000 non-null   int64  
 1   ProductId        50000 non-null   object  
 2   UserId            50000 non-null   object  
 3   ProfileName      49995 non-null   object  
 4   HelpfulnessNumerator  50000 non-null   int64  
 5   HelpfulnessDenominator 50000 non-null   int64  
 6   Score             50000 non-null   int64  
 7   Time              50000 non-null   int64  
 8   Summary            49998 non-null   object  
 9   Text               50000 non-null   object  
dtypes: int64(5), object(5)
memory usage: 3.8+ MB
None
```

Cell 5: Preprocess Data

```
1 # title Cell 5: Preprocess Data
2
3 print("\n--- Data Preprocessing ---")
4 initial_rows = len(df)
5 print(f"Initial rows loaded: {initial_rows}")
6
7 # 1. Handle missing labels
8 # Convert label column to numeric, coercing errors to NaN, then drop rows with NaN labels
9 if not pd.api.types.is_numeric_dtype(df[LABEL_COLUMN]):
10    print(f"Converting label column '{LABEL_COLUMN}' to numeric...")
11    df[LABEL_COLUMN] = pd.to_numeric(df[LABEL_COLUMN], errors='coerce')
12
13 rows_before_label_dropna = len(df)
14 df = df.dropna(subset=[LABEL_COLUMN])
15 rows_after_label_dropna = len(df)
16 print(f"Rows after dropping NaN in label column '{LABEL_COLUMN}': {rows_after_label_dropna} (dropped {rows_before_label_dropna - rows_af
17
18 # 2. Handle missing text
19 rows_before_text_dropna = len(df)
20 df[TEXT_COLUMN] = df[TEXT_COLUMN].astype(str) # Ensure text column is string type
21 df = df.dropna(subset=[TEXT_COLUMN]) # Drop rows where text is explicitly NaN
22 df = df[df[TEXT_COLUMN].str.strip() != ''] # Drop rows where text is empty or only whitespace
23 rows_after_text_dropna = len(df)
24 print(f"Rows after dropping missing/empty text in column '{TEXT_COLUMN}': {rows_after_text_dropna} (dropped {rows_before_text_dropna - r
25
26
27 # 3. Basic text cleaning (optional)
28 def clean_text(text):
29     # Remove HTML tags (common in Amazon reviews)
30     text = re.sub(r'<.*?>', '', text)
31     # Remove URLs
32     text = re.sub(r"http\S+", "", text)
33     # Remove extra whitespace
34     text = re.sub(r"\s+", " ", text).strip()
35     return text
36
37 print("Applying basic text cleaning...")
38 df['cleaned_text'] = df[TEXT_COLUMN].apply(clean_text)
39
40 # 4. Prepare labels: Score is 1-5, models need 0-based labels (0-4)
41 # Ensure the label column is integer type before subtracting
42 df['labels'] = df[LABEL_COLUMN].astype(int) - 1
43 print(f"\nLabel distribution (0-4):")
44 print(df['labels'].value_counts().sort_index()) # Sort index for clarity
45
```

```

46 # Check if all expected labels (0-4) are present after conversion
47 expected_labels = set(range(NUM_LABELS))
48 actual_labels = set(df['labels'].unique())
49 if not expected_labels.issubset(actual_labels):
50     print(f"\nWARNING: Some expected labels (0-{NUM_LABELS-1}) might be missing after processing.")
51     print(f"Expected: {expected_labels}, Found: {actual_labels}")
52     # Consider adjusting NUM_LABELS if the data truly doesn't have all 5 scores,
53     # but usually it's better to keep it as 5 if the original scale was 1-5.
54
55 # 5. Keep only necessary columns and rename
56 df_processed = df[['cleaned_text', 'labels']].rename(columns={'cleaned_text': 'text'})
57 print("\nPreprocessing finished. Final DataFrame sample:")
58 print(df_processed.head())
59 print(f"\nTotal rows after preprocessing: {len(df_processed)}")

```

→ --- Data Preprocessing ---
Initial rows loaded: 50000
Rows after dropping NaN in label column 'Score': 50000 (dropped 0)
Rows after dropping missing/empty text in column 'Text': 50000 (dropped 0)
Applying basic text cleaning...

Label distribution (0-4):
labels
0 4721
1 2814
2 4047
3 7288
4 31130
Name: count, dtype: int64

Preprocessing finished. Final DataFrame sample:

	text	labels
0	I have bought several of the Vitality canned d...	4
1	Product arrived labeled as Jumbo Salted Peanut...	0
2	This is a confection that has been around a fe...	3
3	If you are looking for the secret ingredient i...	1
4	Great taffy at a great price. There was a wide...	4

Total rows after preprocessing: 50000

Cell 6: Split Data

```

1 # @title Cell 6: Split Data
2
3 print("\n--- Splitting Data ---")
4 # Check if there's enough data to split and stratify
5 if len(df_processed) < 10:
6     raise SystemExit(f"Not enough data ({len(df_processed)} rows) remaining after preprocessing to split.")
7
8 if df_processed['labels'].nunique() > 1:
9     try:
10         train_df, test_df = train_test_split(
11             df_processed,
12             test_size=TEST_SIZE,
13             random_state=RANDOM_STATE,
14             stratify=df_processed['labels'] # Ensure similar label distribution in splits
15         )
16         print("Data split using stratification.")
17     except ValueError as e:
18         print(f"Stratify failed: {e}. Splitting without stratify.")
19         train_df, test_df = train_test_split(
20             df_processed,
21             test_size=TEST_SIZE,
22             random_state=RANDOM_STATE
23         )
24 else:
25     print("Only one class label present. Splitting without stratify.")
26     train_df, test_df = train_test_split(
27         df_processed,
28         test_size=TEST_SIZE,
29         random_state=RANDOM_STATE
30     )
31
32 print(f"Training set size: {len(train_df)}")
33 print(f"Test set size: {len(test_df)}")
34
35 # Convert pandas DataFrames to Hugging Face Datasets

```

```

36 train_dataset = Dataset.from_pandas(train_df.reset_index(drop=True))
37 test_dataset = Dataset.from_pandas(test_df.reset_index(drop=True))
38
39 # Combine into a DatasetDict (standard format for Hugging Face)
40 dataset_dict = DatasetDict({
41     'train': train_dataset,
42     'test': test_dataset
43 })
44
45 print("\nDatasetDict created:")
46 print(dataset_dict)

→ --- Splitting Data ---
Data split using stratification.
Training set size: 40000
Test set size: 10000

DatasetDict created:
DatasetDict({
    train: Dataset({
        features: ['text', 'labels'],
        num_rows: 40000
    })
    test: Dataset({
        features: ['text', 'labels'],
        num_rows: 10000
    })
})

```

Cell 7: Tokenization

```

1 # @title Cell 7: Tokenization
2
3 print("\n--- Tokenizing Data ---")
4 # Load tokenizer associated with the chosen model
5 try:
6     tokenizer = AutoTokenizer.from_pretrained(MODEL_CHECKPOINT)
7     print(f"Tokenizer loaded for {MODEL_CHECKPOINT}")
8 except Exception as e:
9     print(f"Error loading tokenizer for {MODEL_CHECKPOINT}: {e}")
10    raise SystemExit("Tokenizer loading failed.")
11
12
13 # Function to tokenize the text data
14 def tokenize_function(examples):
15     # `truncation=True` ensures sequences longer than model max length are cut.
16     # `padding=False` (default) - padding will be handled by DataCollator.
17     # Handle potential non-string data in the 'text' column just in case
18     texts = [str(t) if t is not None else "" for t in examples["text"]]
19     return tokenizer(texts, truncation=True, max_length=512) # Using max_length=512 for BERT-like models
20
21 # Apply tokenization to the entire dataset_dict
22 print("Applying tokenization (this may take a while for large datasets)... ")
23 # `batched=True` processes multiple examples at once for speed
24 tokenized_datasets = dataset_dict.map(tokenize_function, batched=True)
25 print("Tokenization complete.")
26
27 # Remove the original text column as it's no longer needed after tokenization
28 # Also remove '__index_level_0__' if it exists from pandas conversion
29 columns_to_remove = ["text"]
30 if '__index_level_0__' in tokenized_datasets['train'].column_names:
31     columns_to_remove.append('__index_level_0__')
32
33 # Check if columns actually exist before removing
34 existing_columns_to_remove = [col for col in columns_to_remove if col in tokenized_datasets['train'].column_names]
35 if existing_columns_to_remove:
36     tokenized_datasets = tokenized_datasets.remove_columns(existing_columns_to_remove)
37     print(f"Removed columns: {existing_columns_to_remove}")
38
39
40 # Rename 'labels' column to 'label' if necessary (Trainer expects 'label')
41 if 'labels' in tokenized_datasets['train'].column_names and 'label' not in tokenized_datasets['train'].column_names:
42     tokenized_datasets = tokenized_datasets.rename_column("labels", "label")
43     print("Renamed 'labels' column to 'label'.")
44
45

```

```
46 # Set format to PyTorch tensors
47 tokenized_datasets.set_format("torch")
48
49 print("\nTokenized datasets prepared:")
50 print(tokenized_datasets)
51 # Print example only if dataset is not empty
52 if len(tokenized_datasets['train']) > 0:
53     print(f"\nExample tokenized input: {tokenized_datasets['train'][0]}")
54 else:
55     print("\nTraining dataset is empty after processing.")
56
57
58 # Data collator handles dynamic padding (pads sequences to the longest in a batch)
59 data_collator = DataCollatorWithPadding(tokenizer=tokenizer)
60 print("\nData Collator initialized.")
```


Data Collator initialized.

▼ Cell 8: Load Model & Define Metrics

```

1 # @title Cell 8: Load Model & Define Metrics
2
3 print("\n--- Loading Model & Defining Metrics ---")
4
5 # Load the pre-trained model for sequence classification
6 # `num_labels` must match the number of unique sentiment scores (0-4)
7 try:
8     model = AutoModelForSequenceClassification.from_pretrained(
9         MODEL_CHECKPOINT,
10        num_labels=NUM_LABELS,
11        # Add id2label and label2id for better inference later (optional but good practice)
12        id2label={i: f"SCORE_{i+1}" for i in range(NUM_LABELS)}, # Maps 0 -> SCORE_1, 1 -> SCORE_2 etc.
13        label2id={f"SCORE_{i+1}": i for i in range(NUM_LABELS)})
14    )
15    print(f"Loaded model: {MODEL_CHECKPOINT} with {NUM_LABELS} labels.")
16    # print(f"Model label mapping: {model.config.id2label}") # Uncomment to check mapping
17 except Exception as e:
18     print(f"Error loading model {MODEL_CHECKPOINT}: {e}")
19     raise SystemExit("Model loading failed.")
20
21
22 # Define evaluation metrics using the 'evaluate' library
23 try:
24     accuracy_metric = evaluate.load("accuracy")
25     f1_metric = evaluate.load("f1")
26     print("Loaded evaluation metrics: Accuracy, F1")
27 except Exception as e:
28     print(f"Error loading evaluation metrics: {e}")
29     raise SystemExit("Metrics loading failed.")
30
31
32 def compute_metrics(eval_pred):
33     """Computes accuracy and F1 score for evaluation predictions."""
34     logits, labels = eval_pred
35     # Get the predicted class by finding the index with the highest logit value
36     predictions = np.argmax(logits, axis=-1)
37
38     # Calculate accuracy
39     accuracy = accuracy_metric.compute(predictions=predictions, references=labels)
40
41     # Calculate F1 score (use 'weighted' for multi-class imbalance)
42     f1 = f1_metric.compute(predictions=predictions, references=labels, average="weighted")
43
44     return {
45         "accuracy": accuracy["accuracy"],
46         "f1": f1["f1"],
47     }
48
49 print("Metrics computation function defined.")

```

→ Some weights of DistilBertForSequenceClassification were not initialized from the model checkpoint at distilbert-base-uncased and are new. You should probably TRAIN this model on a down-stream task to be able to use it for predictions and inference.

```
--- Loading Model & Defining Metrics ---  
Loaded model: distilbert-base-uncased with 5 labels.  
Loaded evaluation metrics: Accuracy, F1  
Metrics computation function defined.
```

▼ Cell 9: Configure Training Arguments

```
1 # @title Cell 9: Configure Training Arguments
2
3 print("\n--- Configuring Training ---")
4
5 # Define training arguments
```

```

6 # These control various aspects of the fine-tuning process
7 training_args = TrainingArguments(
8     output_dir=OUTPUT_DIR,                      # Directory to save model checkpoints and logs
9     num_train_epochs=NUM_EPOCHS,                 # Total number of training epochs
10    per_device_train_batch_size=BATCH_SIZE,      # Batch size per GPU/CPU for training
11    per_device_eval_batch_size=BATCH_SIZE*2,      # Can often use larger batch size for evaluation
12    learning_rate=LEARNING_RATE,                # Learning rate for the optimizer
13    weight_decay=0.01,                          # Strength of weight decay regularization
14    evaluation_strategy="epoch",               # Evaluate model performance at the end of each epoch
15    save_strategy="epoch",                     # Save model checkpoint at the end of each epoch
16    logging_strategy="epoch",                  # Log training metrics at the end of each epoch
17    load_best_model_at_end=True,              # Load the best model checkpoint (based on metric) at the end
18    metric_for_best_model="f1",                # Metric used to identify the best model (use F1 for potentially imbalanced classes)
19    save_total_limit=2,                      # Only keep the best and the latest checkpoint
20    fp16=torch.cuda.is_available(),           # Use mixed precision training if GPU is available (faster, less memory)
21    push_to_hub=False,                       # Set to True to push model to Hugging Face Hub (requires login)
22    report_to="none",                        # Disable reporting to integrations like wandb/tensorboard unless configured
23    # Add gradient accumulation to simulate larger batch sizes if memory is limited
24    # gradient_accumulation_steps=2,          # Uncomment and adjust if needed (effective batch size = BATCH_SIZE * accumulation_steps)
25 )
26
27 print("Training Arguments configured.")
28 # print(training_args) # Can be verbose, print if needed

```



--- Configuring Training ---
Training Arguments configured.

Cell 10: Initialize Trainer

```

1 # @title Cell 10: Initialize Trainer
2
3 print("\n--- Initializing Trainer ---")
4
5 # Check if train dataset is empty before initializing Trainer
6 if len(tokenized_datasets['train']) == 0:
7     raise SystemExit("Training dataset is empty. Cannot initialize Trainer. Please check data loading and preprocessing steps.")
8
9
10 # Initialize the Trainer
11 # This class orchestrates the fine-tuning process
12 trainer = Trainer(
13     model=model,                                # The instantiated Transformers model to be trained
14     args=training_args,                         # Training arguments defined above
15     train_dataset=tokenized_datasets["train"],   # Training dataset
16     eval_dataset=tokenized_datasets["test"],       # Evaluation dataset
17     tokenizer=tokenizer,                         # Tokenizer (needed for padding/batching consistency)
18     data_collator=data_collator,                # Data collator to handle batch padding
19     compute_metrics=compute_metrics,             # Function to compute evaluation metrics
20 )
21
22 print("Trainer initialized.")

```



--- Initializing Trainer ---
Trainer initialized.

Cell 11: Start Fine-tuning

```

1 # @title Cell 11: Start Fine-tuning
2
3 print("\n--- Starting Fine-tuning ---")
4 print(f"Training for {NUM_EPOCHS} epochs...")
5 # Start the training process
6 try:
7     train_result = trainer.train()
8     print("\n--- Fine-tuning Completed ---")
9
10    # Save training metrics
11    metrics = train_result.metrics
12    trainer.log_metrics("train", metrics)
13    trainer.save_metrics("train", metrics)
14
15    # Save the fine-tuned model and tokenizer
16    trainer.save_model(OUTPUT_DIR) # Saves the best model due to load_best_model_at_end=True

```

```

17     tokenizer.save_pretrained(OUTPUT_DIR)
18     print(f"\nBest fine-tuned model and tokenizer saved to {OUTPUT_DIR}")
19
20 except Exception as e:
21     print(f"\nAn error occurred during training: {e}")
22     # Potentially add more specific error handling (e.g., CUDA out of memory)
23     raise SystemExit("Training failed.")

```

→ --- Starting Fine-tuning ---
Training for 3 epochs... [7500/7500 18:05, Epoch 3/3]

Epoch	Training Loss	Validation Loss	Accuracy	F1
1	0.726800	0.650675	0.743900	0.741853
2	0.551200	0.629174	0.763400	0.755655
3	0.440800	0.660513	0.766200	0.760181

--- Fine-tuning Completed ---

***** train metrics *****

```

epoch          =      3.0
total_flos    =  9030560GF
train_loss     =      0.573
train_runtime  = 0:18:05.65
train_samples_per_second =  110.532
train_steps_per_second  =      6.908

```

Best fine-tuned model and tokenizer saved to sentiment_model_amazon.csv.fineturned

Cell 12: Evaluate Model

```

1 # @title Cell 12: Evaluate Model
2
3 print("\n--- Evaluating Model on Test Set ---")
4 # Evaluate the performance on the test set (loads the best model automatically)
5 try:
6     # Ensure the OUTPUT_DIR exists where the model was saved by the Trainer
7     if not os.path.isdir(OUTPUT_DIR):
8         print(f"ERROR: Model output directory '{OUTPUT_DIR}' not found.")
9         print("Please ensure training (Cell 11) completed successfully and saved the model.")
10        raise SystemExit("Evaluation failed: Model directory not found.")
11
12    eval_results = trainer.evaluate() # Trainer should have loaded the best model
13
14    print("\n--- Evaluation Results ---")
15    print(f"Accuracy: {eval_results.get('eval_accuracy', 'N/A'):.4f}")
16    print(f"F1 Score (Weighted): {eval_results.get('eval_f1', 'N/A'):.4f}")
17    # print(eval_results) # Print full results if needed
18
19    # Save evaluation metrics
20    trainer.log_metrics("eval", eval_results)
21    trainer.save_metrics("eval", eval_results)
22    print("\nEvaluation metrics saved.")
23
24 except Exception as e:
25     print(f"\nAn error occurred during evaluation: {e}")
26     raise SystemExit("Evaluation failed.")

```

→ --- Evaluating Model on Test Set --- [313/313 00:25]

```

--- Evaluation Results ---
Accuracy: 0.7662
F1 Score (Weighted): 0.7602
***** eval metrics *****
epoch          =      3.0
eval_accuracy  =      0.7662
eval_f1        =      0.7602
eval_loss       =      0.6605
eval_runtime   = 0:00:25.54
eval_samples_per_second =  391.532
eval_steps_per_second  =      12.255

```

Evaluation metrics saved.