

# CISC 691: Assignment 03 By Vikramvelupandian Palani

## Fine-Tuning Documentation for Sentiment Analysis:

The approach, objectives, and details of the implementation for fine-tuning a reduced *DistilBert* model on the IMDB sentiment analysis task. It also covers model performance evaluation, challenges encountered, and the solutions applied.

## Motivation for Fine-Tuning in Sentiment Analysis:

***Leveraging Pre-trained Knowledge:*** Through fine-tuning you can utilize models that have been pre-trained on extensive general text datasets. Models such as BERT and GPT-3 along with Roberta represent popular choices. The knowledge obtained from pre-training is tailored for sentiment analysis through selective updating of model weights.

***Confidence in efficiency and Resource Savings:*** Fine-tuning proves to be an efficient solution since it delivers strong performance on specialized tasks while eliminating the need to utilize extensive labeled data and computational power.

***Adaptability:*** The approach facilitates quick development and implementation across numerous practical applications including product service enhancement through customer review analysis and brand reputation management via social media sentiment understanding and user experience improvement through feedback survey analysis.

Training models from scratch requires extensive resources and data while fine-tuning pre-trained models achieves better efficiency.

***Pre-Trained Representations:*** The model possesses existing capabilities in natural language comprehension which necessitates only small modifications for effective sentiment classification specialization.

***Reduced Training Time:*** By utilizing fine-tuning the training time reduces because the model has already learned from extensive data.

***Lower Computational Demand:*** Using fine-tuning reduces the necessity of vast datasets and high computational resources which creates accessibility for organizations with limited resources.

## Practical Applications of Sentiment Analysis

**Customer Feedback:** Evaluate customer satisfaction levels through analysis of product reviews and service interactions.

**Social Media Monitoring:** Sentiment analysis enables businesses to comprehend how the public feels about brands, products or political happenings.

**Market Research:** Tracking consumer behavior trends alongside sentiment analysis helps shape effective marketing strategies.

## **Define Objectives:**

### **Fine-Tuning Goals**

**Adapt a Pre-Trained Model:** We aim to fine-tune a compressed *distilbert-base-uncased* model for positive and negative sentiment classification using a balanced IMDB dataset subset. The clear objective will direct our work and maintain our progress throughout the project.

**Improve Task-Specific Performance:** Fine-tune the model parameters to achieve precise sentiment prediction for new text inputs.

### **Target Performance Metrics**

**Accuracy:** The percentage of correct predictions overall predictions.

**Weighted F1 Score:** The weighted F1 Score measures the harmonic meaning between precision and recall while considering the imbalance between classes. The dataset shows a major imbalance between classes which helps in providing an accurate performance evaluation. The model's effectiveness is assessed through metrics that are calculated throughout both training and evaluation phases.

### **Baseline Model and Target Improvements**

**Baseline Model:** We utilize the original pre-trained DistilBert model as our baseline. The implementation uses a reduced model version with 3 layers instead of the original 6 to enhance training efficiency through a smaller hidden size. The baseline model provides the essential foundation needed for all subsequent improvement efforts.

**Target Improvements:** The fine-tuning process seeks to match the full-size model's performance metrics while creating a smaller model with faster execution for real-time applications.

## **Plan Implementation Steps:**

### **Pre-Trained Model Selection**

**Model Choice:** For model implementation we will use *distilbert-base-uncased* from Hugging Face. The code achieves a smaller footprint by using fewer layers and decreasing hidden size to accelerate both fine-tuning and inference processes.

## **Dataset and Preprocessing Approach**

**Dataset Selection:** Researchers use the IMDB dataset which serves as a popular benchmark for sentiment analysis tasks.

**Data Balancing and Filtering:** The dataset undergoes filtering through the code to maintain equal numbers of positive and negative examples which helps prevent training bias.

**Tokenization:** Text data goes through tokenization by the *AutoTokenizer* which is tied to the base model. The process of tokenization uses padding and truncation methods to ensure input size uniformity.

## **Evaluation Metrics and Testing Strategy**

**Metrics Computation:** The implementation tracks performance by measuring both accuracy and the weighted F1 score through (*sklearn.metrics*).

**Evaluation Strategy:** The system evaluates model performance after each epoch and stores the model with the highest performance. The evaluation function operates independently to test sample sentences and verify model generalization across new inputs.

## **Code Documentation and Explanations:**

### **Logging Configuration**

- **Function:** `configure_logging`  
Sets up both console and file logging to capture detailed debugging and execution flow information.

### **Model Loading and Reduction**

- **Function:** `load_base_model`
  - ✓ Loads the pre-trained model and tokenizer.
  - ✓ Creates a reduced model configuration using `DistilBertConfig` (with fewer layers and a smaller hidden size).
  - ✓ Saves the reduced model locally for fine-tuning later.

### **Dataset Loading, Filtering, and Tokenization**

- **Function:** `load_sentiment_dataset`
  - ✓ Downloads the IMDB dataset from Hugging Face.
  - ✓ Filters and balances the dataset to have an equal number of positive and negative samples.
  - ✓ Tokenize the text data using the associated tokenizer and saves the processed dataset to disk.

## Training Loop and Evaluation

- **Function:** `train_model`
  - ✓ Prepares training and testing data using a custom `SentimentDataset` class.
  - ✓ Configures training arguments (including batch sizes, learning rate, number of epochs, and device settings).
  - ✓ Implements the training loop using the Hugging Face Trainer API, which also computes metrics at each epoch.
  - ✓ Saves the fine-tuned model and tokenizer after training completes.

## Prediction and Evaluation of Sample Sentences

- **Functions:** `predict_sentiment` and `evaluate_sample_sentences`
  - ✓ `predict_sentiment`: Accepts a single text input, tokenizes it, and outputs a sentiment prediction.
  - ✓ `evaluate_sample_sentences`: Reads sample sentences from a JSON file and prints the predicted versus ground-truth sentiments for quick qualitative analysis.

## Main Execution Flow

- **Function:** `main()`

Orchestrates the entire workflow:

  1. Loads or creates the model and dataset.
  2. Triggers the training process.
  3. Runs predictions on test texts.
  4. Evaluates additional sample sentences from an external JSON file.

## Analysis of Model Performance:

• **Metrics Collection:** Accuracy and F1 scores are recorded at each epoch's conclusion during the training process. These metrics provide insights into the classification capabilities of the model.

```
PS C:\Users\vikramp > & C:\Users\vikramp\AppData\Local\Programs\Python\Python312\python.exe "c:\Users\vikramp\OneDrive - School Health Corporation\Desktop\Assignment Files CISC 691\A03\checkpoint 63_Scores.py"

Metrics for checkpoint directory: C:\Users\vikramp\OneDrive - School Health Corporation\Desktop\Assignment Files CISC 691\A03\checkpoint-63
Step: 63 | Accuracy: 0.5 | F1: 0.3333333333333333

Metrics for checkpoint directory: C:\Users\vikramp\OneDrive - School Health Corporation\Desktop\Assignment Files CISC 691\A03\checkpoint-126
Step: 63 | Accuracy: 0.5 | F1: 0.3333333333333333
Step: 126 | Accuracy: 0.538 | F1: 0.4153645240156105

Metrics for checkpoint directory: C:\Users\vikramp\OneDrive - School Health Corporation\Desktop\Assignment Files CISC 691\A03\checkpoint-189
Step: 63 | Accuracy: 0.5 | F1: 0.3333333333333333
Step: 126 | Accuracy: 0.538 | F1: 0.4153645240156105
Step: 189 | Accuracy: 0.762 | F1: 0.7605431444910838
```

• **Evaluation Results:** The evaluation function outputs performance metrics after training which allows for a comparison between the fine-tuned model and baseline models.

• **Logging Insights:** Detailed logging enables users to monitor training progress and detect potential overfitting or underfitting problems. The below image shows that the code restarted from the last checkpoint.

```
File Edit Selection View Go Run Terminal Help Search
cisc691_simple_fine_tuning_VP - Copy.py CISC_691_Simple_fine_tuning_Sample_Sentences.py CISC_691_Simple_fine_tuning_Sample_Sentences_Restartingcode.py X
C:\Users > vikramp > OneDrive - School Health Corporation > Desktop > Assignment Files CISC 691 > A03 > CISC_691_Simple_fine_tuning_Sample_Sentences_Restartingcode.py > train_model
164 def train_model(base_model, tokenizer, train_dataset, test_dataset):
205     train_dataset=train data.
PROBLEMS OUTPUT DEBUG CONSOLE TERMINAL PORTS COMMENTS
PS C:\Users\vikramp > & C:\Users\vikramp\AppData\Local\Programs\Python\Python312\python.exe "c:\Users\vikramp\OneDrive - School Health Corporation\Desktop\Assignment Files CISC 691\A03\checkpoint 63_Scores.py"
[2025-03-04 22:40:54,196] INFO _main_:load base model:83 - Loading model...distilbert-base-uncased
[2025-03-04 22:40:54,197] DEBUG urllib3.connectionpool: new_conn:1049 - Starting new HTTPS connection (1): huggingface.co:443
[2025-03-04 22:40:54,350] DEBUG urllib3.connectionpool: make_request:544 - https://huggingface.co:443 "HEAD /distilbert-base-uncased/resolve/main/config.json HTTP/1.1" 200 0
Some weights of DistilBertForSequenceClassification were not initialized from the model checkpoint at distilbert-base-uncased and are newly initialized: ['classifier.bias', 'classifier.weight', 'pre_classifier.bias', 'pre_classifier.weight']
You should probably TRAIN this model on a down-stream task to be able to use it for predictions and inference.
[2025-03-04 22:40:54,500] DEBUG urllib3.connectionpool: make_request:544 - https://huggingface.co:443 "HEAD /distilbert-base-uncased/resolve/main/tokenizer_config.json HTTP/1.1" 200 0
[2025-03-04 22:40:54,833] INFO _main_:load base model:106 - Base model loaded and saved.
[2025-03-04 22:40:54,843] INFO _main_:train model:165 - Starting training...
C:\Users\vikramp\AppData\Local\Programs\Python\Python312\Lib\site-packages\transformers\training_args.py:1594: FutureWarning: 'evaluation_strategy' is deprecated and will be removed in version 4.46 of Transformers. Use 'eval_strategy' instead
warnings.warn(
[2025-03-04 22:40:55,546] INFO _main_:train model:209 - Running training... This may take a while.
0% | 0/378 [00:00<?, ?it/s]C
:Users\vikramp\AppData\Local\Programs\Python\Python312\Lib\site-packages\transformers\trainer.py:3662: FutureWarning: 'torch.cpu.amp.autocast(args...)' is deprecated
. Please use 'torch.amp.autocast('cpu', args...)' instead.
ctx_manager = torch.cpu.amp.autocast(cache_enabled=cache_enabled, dtype=self.amp_dtype)
50% | 190/378 [07:17:07:12, 2.30s/it]
```

```
83% | 315/378 [7:31:52:58:32, 170.85s/it]C:\Users\vikramp\AppData\Local\Programs\Python\Python312\Lib\
site-packages\transformers\trainer.py:3662: FutureWarning: 'torch.cpu.amp.autocast(args...)' is deprecated. Please use 'torch.amp.autocast('cpu', args...)' instead.
ctx_manager = torch.cpu.amp.autocast(cache_enabled=cache_enabled, dtype=self.amp_dtype)
96% | 364/378 [14:55:49:41:39:48, 427.76s/it]
100% | 378/378 [16:38:01:00:00, 398.79s/it]C:\Users\vikramp\AppData\Local\Programs\Python\Python312\Lib\site-packages\transformers\trainer.py:3662: FutureWarning: 'torch.cpu.amp.autocast(args...)' is deprecated. Please use 'torch.amp.autocast('cpu', args...)' instead.
ctx_manager = torch.cpu.amp.autocast(cache_enabled=cache_enabled, dtype=self.amp_dtype)
('eval_loss': 0.482778935432434, 'eval_accuracy': 0.812, 'eval_f1': 0.8119879672299026, 'eval_runtime': 187.0257, 'eval_samples_per_second': 2.673, 'eval_steps_per_second': 0.086, 'epoch': 6.0)
{'train_runtime': 68070.5843, 'train_samples_per_second': 0.2, 'train_steps_per_second': 0.086, 'train_loss': 0.12282288412412614, 'epoch': 6.0}
100% | 378/378 [16:41:10:00:00, 158.92s/it]
Resumed global step: 378
Training logs: [{'epoch': 1.0, 'eval_accuracy': 0.5, 'eval_f1': 0.3333333333333333, 'eval_loss': 0.6913493871688843, 'eval_runtime': 144.3949, 'eval_samples_per_second': 3.463, 'eval_steps_per_second': 0.111, 'step': 63}, {'epoch': 1.5973015873015874, 'grad_norm': 0.048607063728333, 'learning_rate': 3.677248677248677e-05, 'loss': 0.6956, 'step': 180}, {'epoch': 2.0, 'eval_accuracy': 0.538, 'eval_f1': 0.4153645240156105, 'eval_loss': 0.676928414272, 'eval_runtime': 142.9995, 'eval_samples_per_second': 3.497, 'eval_steps_per_second': 0.112, 'step': 126}, {'epoch': 3.0, 'eval_accuracy': 0.762, 'eval_f1': 0.7605431444910838, 'eval_loss': 0.537821688652839, 'eval_runtime': 280.2415, 'eval_samples_per_second': 1.784, 'eval_steps_per_second': 0.057, 'step': 189}, {'loss': 0.3756, 'grad_norm': 2.664833872602354, 'learning_rate': 2.3544973544973545e-05, 'epoch': 3.1746831746831744, 'step': 200}, {'eval_loss': 0.4759533405303955, 'eval_accuracy': 0.794, 'eval_f1': 0.7936359738578853, 'eval_runtime': 152.7363, 'eval_samples_per_second': 3.274, 'eval_steps_per_second': 0.105, 'epoch': 4.0, 'step': 252}, {'loss': 0.2769, 'grad_norm': 1.9277321100234985, 'learning_rate': 1.8317460317460318e-05, 'epoch': 4.761904761904762, 'step': 308}, {'eval_loss': 0.4823992848396381, 'eval_accuracy': 0.812, 'eval_f1': 0.8119518596760771, 'eval_runtime': 15.33118, 'eval_samples_per_second': 3.261, 'eval_steps_per_second': 0.104, 'epoch': 5.0, 'step': 315}, {'eval_loss': 0.482778935432434, 'eval_accuracy': 0.812, 'eval_f1': 0.8119879672299026, 'eval_runtime': 187.0257, 'eval_samples_per_second': 2.673, 'eval_steps_per_second': 0.086, 'epoch': 6.0, 'step': 378}, {'train_runtime': 68070.5843, 'train_samples_per_second': 0.2, 'train_steps_per_second': 0.086, 'total_flos': 332377694208000.0, 'train_loss': 0.12282288412412614}
Text: I did not enjoy the movie. | Sentiment: Negative
[2025-03-05 15:25:40,200] INFO _main_:evaluate sample sentences:247 - Evaluating sample sentences from file: C:\Users\vikramp\OneDrive - School Health Corporation\Desktop\Assignment Files CISC 691\A03\results2\sample_sentence
```

## Discussion of Challenges and Solutions:

### Challenges

**Data Imbalance and Noise:** The initial dataset could demonstrate class imbalance or include inaccurate labels. To address class imbalance in the dataset we created a balanced subset that included both positive and negative reviews.

**Model Size vs. Performance Trade-Off:** The model performance could be affected when its size is reduced through a decrease in the number of layers and hidden dimensions. This compromise between model size and performance enabled quicker training and inference times which is essential for effective deployment.

**Device Compatibility:** The code performs a strong verification of Apple Silicon (MPS) compatibility while switching to CPU usage when necessary. The system remains adaptable to multiple hardware types because of the detailed attention given to device compatibility.

**Tokenization and Sequence Length:** The correct tokenization of text inputs with uniform padding and truncation was essential to prevent runtime errors while sustaining model performance.

### Solutions

**Balanced Sampling:** The data loading function implements equal sampling for both positive and negative examples.

**Logging and Monitoring:** The project includes detailed logging as an essential element because it delivers valuable information about how the training process unfolds. This system feature helps identify performance and convergence problems which supports better model accuracy.

**Flexible Device Assignment:** Device checks represent a critical component of this project. The system directs the model to operate on optimal hardware resources which results in substantial efficiency improvements.

**Modular Code Design:** The project uses modular functions such as model loading, dataset processing, training, and prediction which simplifies the process of adjusting and troubleshooting different pipeline parts.

## Conclusion and Results:

The process of fine-tuning a sentiment analysis system with a smaller DistilBert model. Through the integration of pre-trained models with balanced datasets and strict evaluation metrics the

system design ensures practical and efficient sentiment classification performance. Modular code design and comprehensive logging resolve the issues of data handling, model scaling, and device compatibility by delivering a solution that remains strong and versatile across various situations. As a result of successfully training the dataset and running the sample\_sentences\_2 file it provides clear results after running the file using the trained model. The final step of this entire project is to use this file to test our pre-trained model efficiently. The image below is the proof the model can predict the sentiment correctly.

```

CISC 691 Simple_fine_tuning_Sample_Sentences_Restartingcode.py X CISC 691 checkpoints 63_126_189_Scores.py
C:\Users\vikramp> OneDrive - School Health Corporation > Desktop > Assignment Files CISC 691 > A03 > CISC 691_Simple_fine_tuning_Sample_Sentences_Restartingcode.py > main
7/65 def main():
PROBLEMS OUTPUT DEBUG CONSOLE TERMINAL PORTS COMMENTS
| 2/16 [00:16<01:56, 8.30s/it]
100%|
Evaluation Results: {'eval_loss': 0.4759633406389955, 'eval_accuracy': 0.794, 'eval_f1': 0.7936359738578853, 'eval_runtime': 265.3275, 'eval_samples_per_second': 1.884, 'epoch': 6.0}
[2025-03-05 16:23:51,613] INFO __main__:train_model:1226 - Saving fine-tuned model...
[2025-03-05 16:23:52,090] INFO __main__:train_model:230 - Training and saving complete.
Text: The movie was great! | Sentiment: Positive
Text: I did not enjoy the movie. | Sentiment: Negative
[2025-03-05 16:24:06,071] INFO __main__:evaluate_sample_sentences:247 - Evaluating sample sentences from file: C:\Users\vikramp\OneDrive - School Health Corporation\Desktop\Assignment Files CISC 691\A03\results2\sample_sentences_2.json
-----
Sentence: This was the worst movie I've ever seen. A complete waste of time.
Language: en
Predicted Sentiment: Negative
Ground Truth: -1
-----
Sentence: Absolutely awful! I regret watching this movie.
Language: en
Predicted Sentiment: Negative
Ground Truth: -1
-----
Sentence: The plot was weak, and the acting was even worse.
Language: en
Predicted Sentiment: Negative
Ground Truth: -0.8
-----
Sentence: It had a few good moments, but overall, it was disappointing.
Language: en
Predicted Sentiment: Positive
Ground Truth: -0.5
-----
Sentence: Not bad, but I wouldn't watch it again.
Language: en
Predicted Sentiment: Negative
Ground Truth: -0.2
-----
Sentence: Es war der schlimmste Film, den ich je gesehen habe.
Language: de
Predicted Sentiment: Negative
Ground Truth: -1
-----
Sentence: Der Film war eine große Enttäuschung.
Language: de
Predicted Sentiment: Negative
Ground Truth: -0.8
-----
Sentence: Ein durchschnittlicher Film mit einigen guten Momenten.
Language: de
Predicted Sentiment: Positive
Ground Truth: 0
-----

```

Below image gives you the results efficiency of the metrics for checkpoints

```

PS C:\Users\vikramp> & C:\Users\vikramp\AppData\Local\Programs\Python\Python312\python.exe "C:\Users\vikramp\OneDrive - School Health Corporation\Desktop\Assignment Files CISC 691\A03\CISC 691_checkpoints 63_126_189_Scores.py"

Metrics for checkpoint directory: C:\Users\vikramp\OneDrive - School Health Corporation\Desktop\Assignment Files CISC 691\A03\checkpoint-63
Step: 63 | Accuracy: 0.5 | F1: 0.3333333333333333

Metrics for checkpoint directory: C:\Users\vikramp\OneDrive - School Health Corporation\Desktop\Assignment Files CISC 691\A03\checkpoint-126
Step: 63 | Accuracy: 0.5 | F1: 0.3333333333333333
Step: 126 | Accuracy: 0.538 | F1: 0.4153645240156105

Metrics for checkpoint directory: C:\Users\vikramp\OneDrive - School Health Corporation\Desktop\Assignment Files CISC 691\A03\checkpoint-189
Step: 63 | Accuracy: 0.5 | F1: 0.3333333333333333
Step: 126 | Accuracy: 0.538 | F1: 0.4153645240156105
Step: 189 | Accuracy: 0.762 | F1: 0.7605431444910838

Metrics for checkpoint directory: C:\Users\vikramp\OneDrive - School Health Corporation\Desktop\Assignment Files CISC 691\A03\checkpoint-252
Step: 63 | Accuracy: 0.5 | F1: 0.3333333333333333
Step: 126 | Accuracy: 0.538 | F1: 0.4153645240156105
Step: 189 | Accuracy: 0.762 | F1: 0.7605431444910838
Step: 252 | Accuracy: 0.794 | F1: 0.7936359738578853

Metrics for checkpoint directory: C:\Users\vikramp\OneDrive - School Health Corporation\Desktop\Assignment Files CISC 691\A03\checkpoint-315
Step: 63 | Accuracy: 0.5 | F1: 0.3333333333333333
Step: 126 | Accuracy: 0.538 | F1: 0.4153645240156105
Step: 189 | Accuracy: 0.762 | F1: 0.7605431444910838
Step: 252 | Accuracy: 0.794 | F1: 0.7936359738578853
Step: 315 | Accuracy: 0.812 | F1: 0.8119518596760771

Metrics for checkpoint directory: C:\Users\vikramp\OneDrive - School Health Corporation\Desktop\Assignment Files CISC 691\A03\checkpoint-378
Step: 63 | Accuracy: 0.5 | F1: 0.3333333333333333
Step: 126 | Accuracy: 0.538 | F1: 0.4153645240156105
Step: 189 | Accuracy: 0.762 | F1: 0.7605431444910838
Step: 252 | Accuracy: 0.794 | F1: 0.7936359738578853
Step: 315 | Accuracy: 0.812 | F1: 0.8119518596760771
Step: 378 | Accuracy: 0.812 | F1: 0.811987967299026
trainer_state.json not found in C:\Users\vikramp\OneDrive - School Health Corporation\Desktop\Assignment Files CISC 691\A03\fine_tuned_model

```

The checkpoints show a clear trend of improvement during training:



- **Checkpoint-63:**  
The initial evaluation shows baseline performance with 50% accuracy and an F1 score of 0.333.
- **Checkpoint-126:**  
A modest improvement is seen with accuracy rising to about 53.8% and the F1 score to approximately 0.415.
- **Checkpoint-189:**  
A significant jump occurs here with accuracy increasing to around 76.2% and F1 nearly matching at 0.7605.
- **Checkpoint-252:**  
Further improvement is observed, with metrics reaching about 79.4% accuracy and an F1 score of 0.7936.
- **Checkpoint-315:**  
Performance improves slightly more to roughly 81.2% accuracy and an F1 score of 0.812.
- **Checkpoint-378:**  
The metrics remain stable at approximately 81.2% for both accuracy and F1 score.

Below image gives you the results efficiency of the metrics of directory

```
Metrics for directory: C:\Users\vikramp\OneDrive - School Health Corporation\Desktop\Assignment Files CISC 691\A03\checkpoint-63
Step: 63 | Accuracy: 0.5 | F1: 0.3333333333333333

Metrics for directory: C:\Users\vikramp\OneDrive - School Health Corporation\Desktop\Assignment Files CISC 691\A03\checkpoint-126
Step: 63 | Accuracy: 0.5 | F1: 0.3333333333333333
Step: 126 | Accuracy: 0.538 | F1: 0.4153645240156105

Metrics for directory: C:\Users\vikramp\OneDrive - School Health Corporation\Desktop\Assignment Files CISC 691\A03\checkpoint-189
Step: 63 | Accuracy: 0.5 | F1: 0.3333333333333333
Step: 126 | Accuracy: 0.538 | F1: 0.4153645240156105
Step: 189 | Accuracy: 0.762 | F1: 0.7605431444910838

Metrics for directory: C:\Users\vikramp\OneDrive - School Health Corporation\Desktop\Assignment Files CISC 691\A03\checkpoint-252
Step: 63 | Accuracy: 0.5 | F1: 0.3333333333333333
Step: 126 | Accuracy: 0.538 | F1: 0.4153645240156105
Step: 189 | Accuracy: 0.762 | F1: 0.7605431444910838
Step: 252 | Accuracy: 0.794 | F1: 0.7936359738578853

Metrics for directory: C:\Users\vikramp\OneDrive - School Health Corporation\Desktop\Assignment Files CISC 691\A03\checkpoint-315
Step: 63 | Accuracy: 0.5 | F1: 0.3333333333333333
Step: 126 | Accuracy: 0.538 | F1: 0.4153645240156105
Step: 189 | Accuracy: 0.762 | F1: 0.7605431444910838
Step: 252 | Accuracy: 0.794 | F1: 0.7936359738578853
Step: 315 | Accuracy: 0.812 | F1: 0.8119518596760771

Metrics for directory: C:\Users\vikramp\OneDrive - School Health Corporation\Desktop\Assignment Files CISC 691\A03\checkpoint-378
Step: 63 | Accuracy: 0.5 | F1: 0.3333333333333333
Step: 126 | Accuracy: 0.538 | F1: 0.4153645240156105
Step: 189 | Accuracy: 0.762 | F1: 0.7605431444910838
Step: 252 | Accuracy: 0.794 | F1: 0.7936359738578853
Step: 315 | Accuracy: 0.812 | F1: 0.8119518596760771
Step: 378 | Accuracy: 0.812 | F1: 0.8119879672299026

Metrics for directory: C:\Users\vikramp\OneDrive - School Health Corporation\Desktop\Assignment Files CISC 691\A03\fine_tuned model
```

Finally, obtained the final fine-tuned model results. Here's what happened:



1. **Resuming from Checkpoint 378:**

The log shows “Resumed global step: 378,” meaning that the training run picked up from the last saved checkpoint (checkpoint-378). No additional training steps were taken beyond that point.

2. **Evaluation Metrics:**

After resuming, the training script immediately proceeded to evaluation. The evaluation results show an accuracy of about 79.4% and an F1 score of roughly 0.7936. These metrics represent the performance of the model as of checkpoint 378, which is the final state of your fine-tuning.

3. **Saving and Inference:**

The script then saved the fine-tuned model and ran predictions on sample sentences. The predictions (e.g., "The movie was great!" → Positive) confirm that the model from checkpoint 378 is now being used for inference.

In summary, by resuming from checkpoint 378 and running the evaluation right after, have effectively obtained and validated the final fine-tuned model results.