```
# Install required libraries
!pip install transformers datasets scikit-learn --quiet
!pip install evaluate --quiet
\overline{2}
                                            — 491.4/491.4 kB 15.5 MB/s eta 0:00:00
                    116.3/116.3 kB 11.8 MB/s eta 0:00:00
                                           — 193,6/193,6 kB 19,3 MB/s eta 0:00:00
                                        ----- 143.5/143.5 kB 13.1 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 the fo
    torch 2.6.0+cu124 requires nvidia-cublas-cu12==12.4.5.8; platform system == "Linux" and platform machine == "x86 64", but you have nvidia-cublas-
    torch 2.6.0+cu124 requires nvidia-cuda-cupti-cu12==12.4.127; platform system == "Linux" and platform machine == "x86 64", but you have nvidia-cud
    torch 2.6.0+cu124 requires nvidia-cuda-nvrtc-cu12==12.4.127; platform system == "Linux" and platform machine == "x86 64", but you have nvidia-cud
    torch 2.6.0+cu124 requires nvidia-cuda-runtime-cu12==12.4.127; platform system == "Linux" and platform machine == "x86 64", but you have nvidia-c
    torch 2.6.0+cu124 requires nvidia-cudnn-cu12==9.1.0.70; platform system == "Linux" and platform machine == "x86 64", but you have nvidia-cudnn-cu
    torch 2.6.0+cu124 requires nvidia-cufft-cu12==11.2.1.3; platform system == "Linux" and platform machine == "x86 64", but you have nvidia-cufft-cu
    torch 2.6.0+cu124 requires nvidia-curand-cu12==10.3.5.147; platform system == "Linux" and platform machine == "x86 64", but you have nvidia-curan
    torch 2.6.0+cu124 requires nvidia-cusolver-cu12==11.6.1.9; platform system == "Linux" and platform machine == "x86 64", but you have nvidia-cusol
    torch 2.6.0+cu124 requires nvidia-cusparse-cu12==12.3.1.170; platform system == "Linux" and platform machine == "x86 64", but you have nvidia-cus
    torch 2.6.0+cu124 requires nvidia-nvjitlink-cu12==12.4.127; platform system == "Linux" and platform machine == "x86 64", but you have nvidia-nvji
    gcsfs 2025.3.2 requires fsspec==2025.3.2. but you have fsspec 2025.3.0 which is incompatible.
                             84.0/84.0 kB 3.9 MB/s eta 0:00:00
import torch
torch.cuda.is available()
→ True
# 🖺 Import required packages
import pandas as pd
import numpy as np
import matplotlib.pyplot as plt
import os
import torch
from datasets import Dataset
from sklearn.model_selection import train_test_split
from sklearn.metrics import accuracy score
from transformers import AutoTokenizer, AutoModelForSequenceClassification, Trainer, TrainingArguments
```

See Load the Amazon Fine Food Reviews dataset

import pandas as pd

```
print(" Loading dataset...")
df = pd.read csv('/content/amazon reviews.csv')
# % Keep only 'Text' and 'Score' columns
df = df[['Text', 'Score']].rename(columns={'Text': 'text', 'Score': 'score'})
# <code>Map Sentiment (0: Negative, 1: Neutral, 2: Positive)</code>
def map sentiment(score):
    if score <= 2:
         return 0 # Negative
    elif score == 3:
         return 1 # Neutral
    else:
         return 2 # Positive
df['label'] = df['score'].apply(map_sentiment)
print(f"▼ Dataset loaded! Total samples: {len(df)}")
df.head()
     Loading dataset...

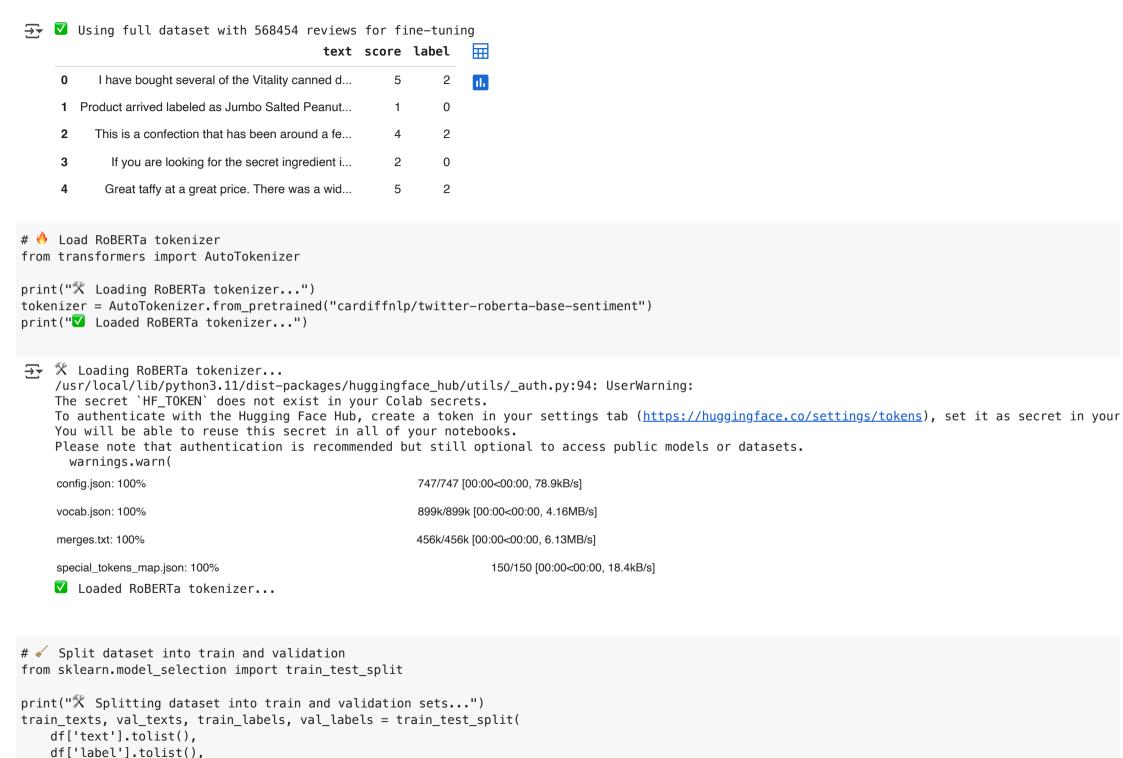
☑ Dataset loaded! Total samples: 568454

                                                                   \blacksquare
                                             text score label
           I have bought several of the Vitality canned d...
      0
                                                                   ıl.
      1 Product arrived labeled as Jumbo Salted Peanut...
                                                               0
      2
           This is a confection that has been around a fe...
                                                        4
                                                               2
                                                               0
      3
             If you are looking for the secret ingredient i...
                                                        2
      4
            Great taffy at a great price. There was a wid...
                                                        5
                                                               2
```

Duse full dataset without sampling

df.head()

print(f"♥ Using full dataset with {len(df)} reviews for fine-tuning")



```
test size=0.2,
    random state=42
print(f"☑ Training samples: {len(train texts)}")
print(f"♥ Validation samples: {len(val texts)}")
    X Splitting dataset into train and validation sets...
    ▼ Training samples: 454763

✓ Validation samples: 113691

# Nefine tokenization function
def tokenize(batch):
    return tokenizer(batch['text'], padding='max length', truncation=True, max length=128)
# % Format as Huggingface Datasets
from datasets import Dataset
train dataset = Dataset.from dict({'text': train texts, 'label': train labels})
val dataset = Dataset.from dict({'text': val texts, 'label': val labels})
# 	✓ Tokenize the datasets
train dataset = train dataset.map(tokenize, batched=True)
val dataset = val dataset.map(tokenize, batched=True)
print("
Tokenization and dataset formatting complete.")
    Map: 100%
                                                  454763/454763 [02:31<00:00, 2561.14 examples/s]
     Map: 100%
                                                  113691/113691 [00:43<00:00, 2169.28 examples/s]

▼ Tokenization and dataset formatting complete.

# 🤚 Load RoBFRTa model
from transformers import AutoModelForSequenceClassification
print("% Loading RoBERTa model...")
model = AutoModelForSequenceClassification.from pretrained(
    "cardiffnlp/twitter-roberta-base-sentiment",
    num_labels=3 # 3 classes: Negative, Neutral, Positive
```

X Loading RoBERTa model...
Xet Storage is enabled for this repo, but the 'hf_xet' package is not installed. Falling back to regular HTTP download. For better performance, i WARNING:huggingface_hub.file_download:Xet Storage is enabled for this repo, but the 'hf_xet' package is not installed. Falling back to regular HT pytorch_model.bin: 100%

499M/499M [00:01<00:00, 310MB/s]

```
from transformers import TrainingArguments
training_args = TrainingArguments(
    output dir="./results",
    per_device_train_batch_size=32,
    per device eval batch size=32,
    num_train_epochs=1,
    warmup steps=50,
    weight_decay=0.01,
    logging dir="./logs",
    save steps=10000,
    logging_steps=500,
    report to="none"
# Define compute_metrics function
from sklearn.metrics import accuracy score
import torch
def compute metrics(eval pred):
    logits, labels = eval_pred
    predictions = torch.argmax(torch.tensor(logits), axis=-1)
    return {"accuracy": accuracy score(labels, predictions)}
from transformers import Trainer
trainer = Trainer(
    model=model,
    args=training args,
    train_dataset=train_dataset,
    eval dataset=val dataset,
    compute_metrics=compute_metrics
```

```
# ✓ Mount Google Drive
from google.colab import drive
drive.mount('/content/drive')
# 🔥 Start fine-tuning RoBERTa
trainer.train()
# # 
Evaluate validation accuracy after training
results = trainer.evaluate()
print(f"▼ Final Validation Accuracy: {results['eval accuracy'] * 100:.2f}%")
# Define save path in Google Drive
drive model dir = "/content/drive/MyDrive/roberta finetuned/"
os.makedirs(drive model dir, exist ok=True)
# № Save model and tokenizer directly to Google Drive
model.save pretrained(drive model dir)
tokenizer.save_pretrained(drive_model_dir)
# h Save validation accuracy
accuracy = results['eval accuracy']
with open(os.path.join(drive model dir, "accuracy.txt"), "w") as f:
    f.write(str(accuracy))
print(f"▼ Validation Accuracy Saved to Drive: {accuracy:.4f}")
# C Zip and save to Google Drive
!zip -r "/content/drive/MyDrive/roberta_finetuned.zip" ./models/roberta_finetuned/
```

print("☑ Zipped model saved to Google Drive.")

-	_	_
	•	_
-	7	_

Step	Training	Loss
500	0.3	56100
1000	0.30	09600
1500	0.29	97500
2000	0.29	95300
2500	0.28	88200
3000	0.20	61400
3500	0.25	58800
4000	0.20	63900
4500	0.20	60100
5000	0.25	53800
5500	0.2	51800
6000	0.24	42900
6500	0.24	42500
7000	0.23	39400
7500	0.23	35100
8000	0.23	33100
8500	0.22	24400
9000	0.22	26200
9500	0.22	26900
10000	0.2	17900
10500	0.2	16900
11000	0.20	07000
11500	0.20	06900
12000	0.19	99500
12500	0.20	05400

```
13000
                  0.207100
     13500
                  0.205300
     14000
                  0.201100
                                         [1620/3553 05:25 < 06:28, 4.98 it/s]
                                        ) [3553/3553 11:51]
       Final Validation Accuracy: 92.92%
       Validation Accuracy Saved to Drive: 0.9292
             zip warning: name not matched: ./models/roberta finetuned/
    zip error: Nothing to do! (try: zip -r /content/drive/MyDrive/roberta_finetuned.zip . -i ./models/roberta_finetuned/)

✓ Zipped model saved to Google Drive.

results = trainer.evaluate()
print(f"▼ Final Validation Accuracy: {results['eval_accuracy']*100:.2f}%")
                                         [3553/3553 16:17]
model.save_pretrained("./models/roberta_finetuned/")
tokenizer.save_pretrained("./models/roberta_finetuned/")
with open('./models/roberta_finetuned/accuracy.txt', 'w') as f:
    f.write(str(results['eval accuracy']))
!zip -r roberta_finetuned.zip ./models/roberta_finetuned/
from google.colab import files
files.download('roberta_finetuned.zip')
Start coding or generate with AI.
```