

```

1 import pandas as pd
2 import re
3 import torch
4 import numpy as np
5 import matplotlib.pyplot as plt
6
7 from sklearn.model_selection import train_test_split
8 from sklearn.metrics import accuracy_score, classification_report, confusion_matrix, ConfusionMatrixDisplay
9
10 from transformers import (
11     AutoTokenizer,
12     AutoModelForSequenceClassification,
13     Trainer,
14     TrainingArguments
15 )
16
17 from torch.utils.data import Dataset
18

```

```

1 df = pd.read_csv("BTC_tweets_daily_example.csv")
2 print(df.head())
3

```

```

Unnamed: 0      Date \
0      0  Fri Mar 23 00:40:32 +0000 2018
1      1  Fri Mar 23 00:40:34 +0000 2018
2      2  Fri Mar 23 00:40:35 +0000 2018
3      3  Fri Mar 23 00:40:36 +0000 2018
4      4  Fri Mar 23 00:40:36 +0000 2018

      Tweet      Screen_name \
0  RT @ALXTOKEN: Paul Krugman, Nobel Luddite. I h...  myresumerocket
1  @lopp @ Kevin_Pham @psycho_sage @naval But @Pr...  BitMocro
2  RT @tippereconomy: Another use case for #block...  hojachotopur
3      free coins https://t.co/DiuoePJdap  denies_distro
4  RT @payvxofficial: WE are happy to announce th...  aditzgraha

      Source \
0      []
1      [u'Bitcoin']
2  [u'blockchain', u'Tipper', u'TipperEconomy']
3      []
4      []

      Link      Sentiment \
0  <a href="http://twitter.com" rel="nofollow">Tw...  ['neutral']
1  <a href="http://twitter.com/download/android" ...  ['neutral']
2  <a href="http://twitter.com" rel="nofollow">Tw...  ['positive']
3  <a href="http://twitter.com" rel="nofollow">Tw...  ['positive']
4  <a href="http://twitter.com/download/android" ...  ['positive']

sent_score  New_Sentiment_Score  New_Sentiment_State
0      0.0      0.000000      0.0
1      0.0      0.000000      0.0
2      1.0      0.136364      1.0
3      1.0      0.400000      1.0
4      1.0      0.468182      1.0

```

```

1 def clean_text(text):
2     text = str(text).lower()
3     text = re.sub(r"http\S+", "", text)
4     text = re.sub(r"@w+", "", text)
5     text = re.sub(r"#w+", "", text)
6     text = re.sub(r"^[a-z\s]", "", text)
7     return text
8

```

```

1 df["Tweet"] = df["Tweet"].astype(str)
2 df["Tweet"] = df["Tweet"].apply(clean_text)
3
4 # Remove rows without labels
5 df = df.dropna(subset=["New_Sentiment_State"])
6

```

```

1 label_mapping = {
2     label: idx for idx, label in enumerate(sorted(df["New_Sentiment_State"].unique()))
3 }
4
5 df["label"] = df["New_Sentiment_State"].map(label_mapping)
6

```

```

7 num_labels = len(label_mapping)
8 print("Label mapping:", label_mapping)
9

```

```
Label mapping: {np.float64(-1.0): 0, np.float64(0.0): 1, np.float64(1.0): 2}
```

```

1 X_train, X_test, y_train, y_test = train_test_split(
2     df["Tweet"],
3     df["label"],
4     test_size=0.2,
5     random_state=42,
6     stratify=df["label"]
7 )
8

```

```

1 MODEL_NAME = "distilroberta-base"
2
3 tokenizer = AutoTokenizer.from_pretrained(MODEL_NAME)
4
5 device = torch.device("cuda" if torch.cuda.is_available() else "cpu")
6
7 model = AutoModelForSequenceClassification.from_pretrained(
8     MODEL_NAME,
9     num_labels=num_labels
10 ).to(device)

```

/usr/local/lib/python3.12/dist-packages/huggingface\_hub/utils/\_auth.py:94: UserWarning:  
The secret `HF\_TOKEN` does not exist in your Colab secrets.

To authenticate with the Hugging Face Hub, create a token in your settings tab (<https://huggingface.co/settings/tokens>), set  
You will be able to reuse this secret in all of your notebooks.

Please note that authentication is recommended but still optional to access public models or datasets.

```
warnings.warn(
```

Loading weights: 100%

101/101 [00:00<00:00, 314.86it/s, Materializing param=roberta.encoder.layer.5.output.dense.weight]

**RobertaForSequenceClassification LOAD REPORT** from: distilroberta-base

Key	Status
lm_head.layer_norm.bias	UNEXPECTED
roberta.pooler.dense.weight	UNEXPECTED
lm_head.dense.weight	UNEXPECTED
lm_head.layer_norm.weight	UNEXPECTED
roberta.pooler.dense.bias	UNEXPECTED
lm_head.dense.bias	UNEXPECTED
lm_head.bias	UNEXPECTED
classifier.dense.weight	MISSING
classifier.out_proj.weight	MISSING
classifier.out_proj.bias	MISSING
classifier.dense.bias	MISSING

Notes:

- **UNEXPECTED** :can be ignored when loading from different task/architecture; not ok if you expect identical arch.
- **MISSING** :those params were newly initialized because missing from the checkpoint. Consider training on your downstream task

```

1 class TweetDataset(Dataset):
2     def __init__(self, texts, labels):
3         self.encodings = tokenizer(
4             texts.tolist(),
5             truncation=True,
6             padding=True,
7             max_length=128
8         )
9         self.labels = labels.tolist()
10
11     def __getitem__(self, idx):
12         item = {k: torch.tensor(v[idx]) for k, v in self.encodings.items()}
13         item["labels"] = torch.tensor(self.labels[idx])
14         return item
15
16     def __len__(self):
17         return len(self.labels)
18

```

```

1 train_dataset = TweetDataset(X_train, y_train)
2 test_dataset = TweetDataset(X_test, y_test)
3

```

```

1 training_args = TrainingArguments(
2     output_dir="./results",
3     per_device_train_batch_size=32,
4     per_device_eval_batch_size=32,

```

```

5     num_train_epochs=1,          # 🔥 FAST
6     learning_rate=2e-5,
7     weight_decay=0.01,
8     logging_steps=100,
9     save_strategy="no",
10    eval_strategy="no",
11    fp16=True,                   # 🔥 MUCH FASTER ON GPU
12    report_to="none"
13 )

```

```

1 trainer = Trainer(
2     model=model,
3     args=training_args,
4     train_dataset=train_dataset
5 )

```

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

 [1272/1272 01:47, Epoch 1/1]

Step	Training Loss
100	0.824642
200	0.438055
300	0.342024
400	0.292431
500	0.248264
600	0.246954
700	0.171308
800	0.195838
900	0.157617
1000	0.153829
1100	0.161120
1200	0.164257

TrainOutput(global\_step=1272, training\_loss=0.2765104238342189, metrics={'train\_runtime': 109.0079, 'train\_samples\_per\_second': 373.193, 'train\_steps\_per\_second': 11.669, 'total\_flos': 831506221697598.0, 'train\_loss':

```

1 predictions = trainer.predict(test_dataset)
2 y_pred = np.argmax(predictions.predictions, axis=1)
3
4 print("Accuracy:", accuracy_score(y_test, y_pred))
5 print("\nClassification Report:\n")
6 print(classification_report(y_test, y_pred, target_names=[str(key) for key in label_mapping.keys()]))

```

Accuracy: 0.9584111690099302

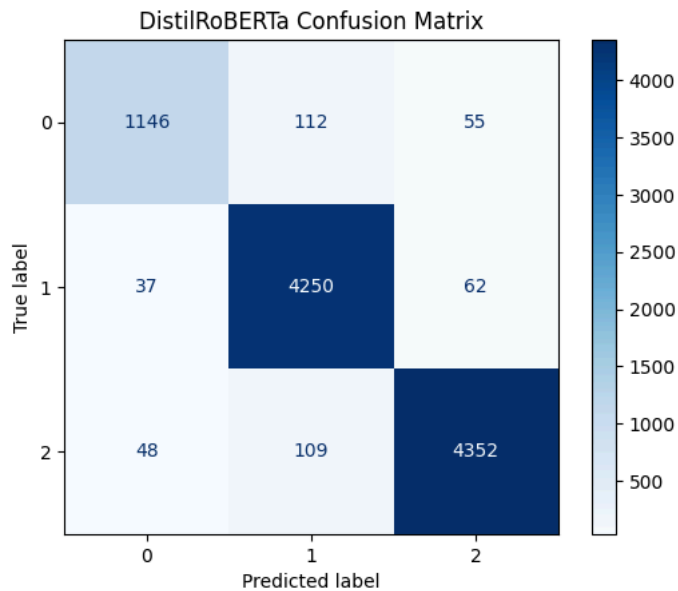
Classification Report:

	precision	recall	f1-score	support
-1.0	0.93	0.87	0.90	1313
0.0	0.95	0.98	0.96	4349
1.0	0.97	0.97	0.97	4509
accuracy			0.96	10171
macro avg	0.95	0.94	0.94	10171
weighted avg	0.96	0.96	0.96	10171

```

1 cm = confusion_matrix(y_test, y_pred)
2
3 disp = ConfusionMatrixDisplay(confusion_matrix=cm)
4 disp.plot(cmap="Blues")
5 plt.title("DistilRoBERTa Confusion Matrix")
6 plt.show()
7

```

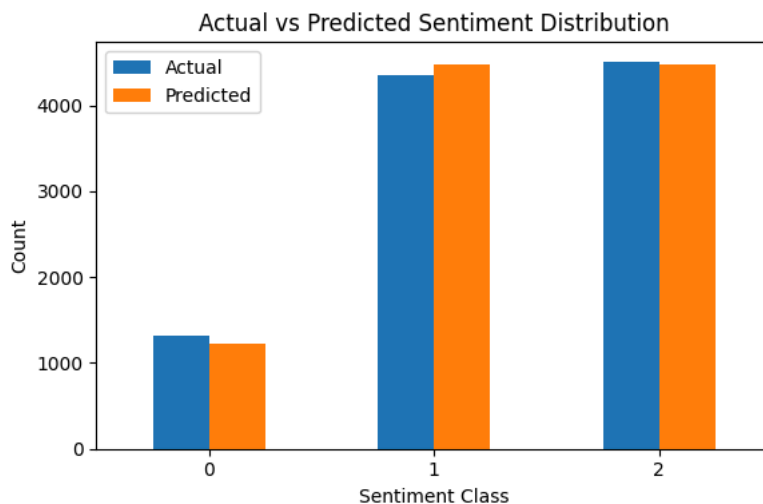


1 Start coding or [generate](#) with AI.

```

1 import pandas as pd
2 import matplotlib.pyplot as plt
3
4 # Convert numeric labels to Series
5 actual_counts = pd.Series(y_test).value_counts().sort_index()
6 pred_counts = pd.Series(y_pred).value_counts().sort_index()
7
8 # Align indexes so missing classes don't cause NaN
9 df_plot = pd.DataFrame({
10     "Actual": actual_counts,
11     "Predicted": pred_counts
12 }).fillna(0)
13
14 # Plot
15 df_plot.plot(
16     kind="bar",
17     figsize=(6,4)
18 )
19
20 plt.title("Actual vs Predicted Sentiment Distribution")
21 plt.xlabel("Sentiment Class")
22 plt.ylabel("Count")
23 plt.xticks(rotation=0)
24 plt.tight_layout()
25 plt.show()
26

```

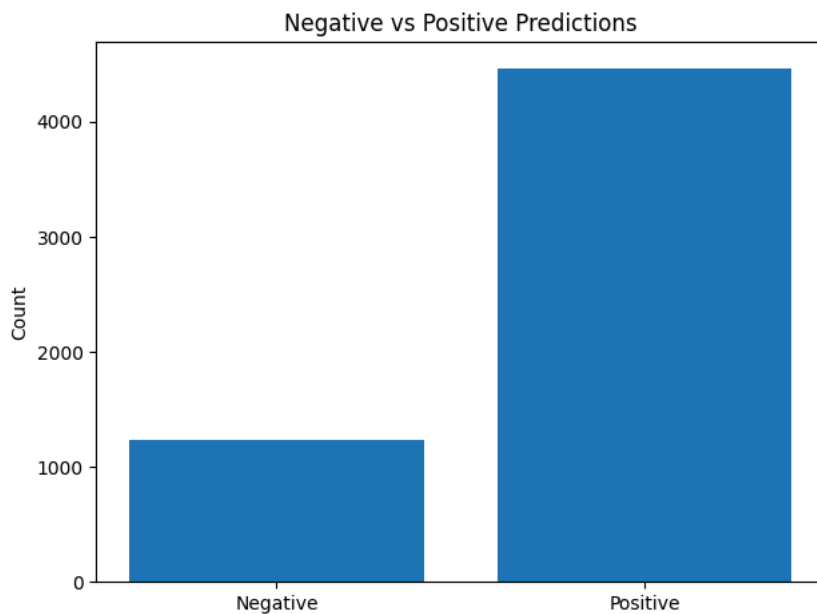


```

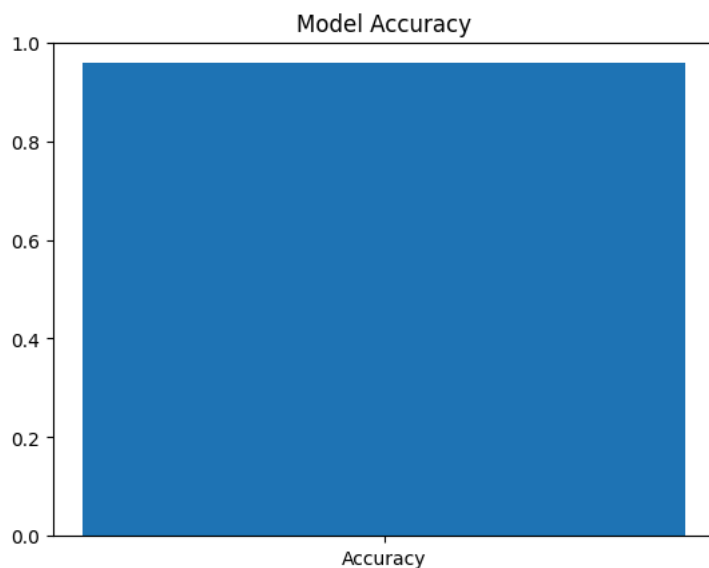
1 import matplotlib.pyplot as plt
2
3 neg_pos = {
4     "Negative": pred_counts.get(0, 0), # class 0

```

```
5     "Positive": pred_counts.get(2, 0)    # class 2
6 }
7
8 plt.bar(neg_pos.keys(), neg_pos.values())
9 plt.title("Negative vs Positive Predictions")
10 plt.ylabel("Count")
11 plt.tight_layout()
12 plt.show()
13
```



```
1 acc = accuracy_score(y_test, y_pred)
2
3 plt.bar(["Accuracy"], [acc])
4 plt.ylim(0, 1)
5 plt.title("Model Accuracy")
6 plt.show()
7
```



```
1 import torch
2 from transformers import AutoTokenizer, AutoModelForSequenceClassification
3
4 # Load DistilRoBERTa
5 MODEL_NAME = "distilroberta-base"
6
7 tokenizer = AutoTokenizer.from_pretrained(MODEL_NAME)
8 model = AutoModelForSequenceClassification.from_pretrained(
9     MODEL_NAME,
10     num_labels=3
11 )
12
13 model.eval()
14
```

```

15 # Input sentence
16 sentence = "I just love how Bitcoin destroyed my savings today"
17
18 # Tokenize
19 inputs = tokenizer(
20     sentence,
21     return_tensors="pt",
22     truncation=True,
23     padding=True
24 )
25
26 # Predict
27 with torch.no_grad():
28     outputs = model(**inputs)
29     probs = torch.softmax(outputs.logits, dim=1)
30
31 # Labels (model output space: 0,1,2)
32 labels = {
33     0: "Negative 🙄",
34     1: "Neutral 😐",
35     2: "Positive 😊"
36 }
37
38 # Max-probability decision
39 pred_id = torch.argmax(probs).item()
40
41 # Output
42 print("Sentence:", sentence)
43 print("Probabilities:")
44 print("Negative:", round(probs[0][0].item(), 3))
45 print("Neutral :", round(probs[0][1].item(), 3))
46 print("Positive:", round(probs[0][2].item(), 3))
47 print("\nFinal Prediction (max prob):", labels[pred_id])
48

```

Loading weights: 100%

101/101 [00:00&lt;00:00, 414.06it/s, Materializing param=roberta.encoder.layer.5.output.dense.weight]

**RobertaForSequenceClassification LOAD REPORT** from: distilroberta-base

Key	Status	
-----+-----+		
lm_head.layer_norm.bias	UNEXPECTED	
roberta.pooler.dense.weight	UNEXPECTED	
lm_head.dense.weight	UNEXPECTED	
lm_head.layer_norm.weight	UNEXPECTED	
roberta.pooler.dense.bias	UNEXPECTED	
lm_head.dense.bias	UNEXPECTED	
lm_head.bias	UNEXPECTED	
classifier.dense.weight	MISSING	
classifier.out_proj.weight	MISSING	
classifier.out_proj.bias	MISSING	
classifier.dense.bias	MISSING	

Notes:

- **UNEXPECTED** :can be ignored when loading from different task/architecture; not ok if you expect identical arch.
- **MISSING** :those params were newly initialized because missing from the checkpoint. Consider training on your downstream task

Sentence: I just love how Bitcoin destroyed my savings today

Probabilities:

Negative: 0.362

Neutral : 0.324

Positive: 0.315

Final Prediction (max prob): Negative 🙄