```
1  !pip install --upgrade transformers datasets accelerate
2  from google.colab import drive
3  drive.mount('/content/drive')
4  from sklearn.metrics import accuracy_score, precision_recall_fscore_support
5  import transformers
6  import pandas as pd
7  from datasets import Dataset
8  from transformers import BertTokenizer, BertForSequenceClassification, Trainer,
   TrainingArguments
9  from sklearn.model_selection import train_test_split
10 import torch
```

Show hidden output

```
1 def compute_metrics(pred):
2     labels = pred.label_ids
3     preds = pred.predictions.argmax(-1)
4     acc = accuracy_score(labels, preds)
5     precision, recall, f1, _ = precision_recall_fscore_support(labels, preds, average='binary')
6     return {
7         'accuracy': acc,
8         'precision': precision,
9         'recall': recall,
10        'f1': f1
11    }
```

Drive already mounted at /content/drive; to attempt to forcibly remount, call drive.mount("/content/drive", force_remount=True).

```
1  # Only use the input text and hate label
2  df = pd.read_csv('/content/drive/MyDrive/train_with_topics.csv')
3  df = df[['text', 'topic_label', 'hate_label']].dropna()
4  df['hate_label'] = df['hate_label'].astype(int)
5
6  # Split into train/test
7  from sklearn.model_selection import train_test_split
8  train_df, val_df = train_test_split(df, test_size=0.1, random_state=42)
9
10 # Convert to Hugging Face dataset
11 train_dataset = Dataset.from_pandas(train_df)
12 val_dataset = Dataset.from_pandas(val_df)
13 full_dataset = Dataset.from_pandas(df)
```

```
1  tokenizer = BertTokenizer.from_pretrained("bert-base-uncased")
2
3  def tokenize(example):
4      return tokenizer(example['text'], truncation=True, padding="max_length", max_length=128)
5
6  train_dataset = train_dataset.map(tokenize, batched=True)
7  val_dataset = val_dataset.map(tokenize, batched=True)
8  full_dataset = full_dataset.map(tokenize, batched=True)
9
10 train_dataset = train_dataset.rename_column("hate_label", "labels")
11 val_dataset = val_dataset.rename_column("hate_label", "labels")
12 full_dataset = full_dataset.rename_column("hate_label", "labels")
13
14 train_dataset.set_format("torch")
15 val_dataset.set_format("torch")
16 full_dataset.set_format("torch")
17
```

| Map: 100% | 13927/13927 [00:08<00:00, 1621.94 examples/s] |
| Map: 100% | 1548/1548 [00:00<00:00, 1679.21 examples/s] |
| Map: 100% | 15475/15475 [00:09<00:00, 1727.65 examples/s] |

```
1  # Load model
2  model = BertForSequenceClassification.from_pretrained("bert-base-uncased", num_labels=2)
3
4  # Define training arguments
5  training_args = TrainingArguments(
6      output_dir="/content/drive/MyDrive/hate_model_basic_full_dataset",
7      save_strategy="epoch",
8      per_device_train_batch_size=8,
9      per_device_eval_batch_size=8,
10     num_train_epochs=10,
11     logging_dir="/content/logs"
12 )
13
14 trainer = Trainer(
15     model=model,
16     args=training_args,
17     train_dataset=train_dataset,
18     eval_dataset=val_dataset,
19     tokenizer=tokenizer,
20     compute_metrics=compute_metrics
21 )
22
23 # Train and evaluate
24 trainer.train()
25 trainer.evaluate()
26
```

[13849/17410 12:47 < 03:17, 18.05 it/s, Epoch 7.95/10]

| Step | Training Loss |
| --- | --- |
| 500 | 0.461300 |
| 1000 | 0.450400 |
| 1500 | 0.423300 |
| 2000 | 0.389100 |
| 2500 | 0.354800 |
| 3000 | 0.357800 |
| 3500 | 0.342500 |
| 4000 | 0.248700 |
| 4500 | 0.260900 |
| 5000 | 0.264600 |
| 5500 | 0.209900 |
| 6000 | 0.188000 |
| 6500 | 0.159400 |
| 7000 | 0.193800 |
| 7500 | 0.100000 |
| 8000 | 0.117100 |
| 8500 | 0.120300 |
| 9000 | 0.095100 |
| 9500 | 0.068600 |
| 10000 | 0.060700 |
| 10500 | 0.063500 |
| 11000 | 0.044400 |
| 11500 | 0.058300 |
| 12000 | 0.055000 |
| 12500 | 0.040100 |
| 13000 | 0.033100 |
| 13500 | 0.032600 |

[17410/17410 16:06, Epoch 10/10]

| Step | Training Loss |
| --- | --- |
| 500 | 0.461300 |
| 1000 | 0.450400 |
| 1500 | 0.423300 |
| 2000 | 0.389100 |
| 2500 | 0.354800 |
| 3000 | 0.357800 |
| 3500 | 0.342500 |
| 4000 | 0.248700 |
| 4500 | 0.260900 |
| 5000 | 0.264600 |
| 5500 | 0.209900 |
| 6000 | 0.188000 |
| 6500 | 0.159400 |
| 7000 | 0.193800 |
| 7500 | 0.100000 |
| 8000 | 0.117100 |
| 8500 | 0.120300 |
| 9000 | 0.095100 |
| 9500 | 0.068600 |
| 10000 | 0.060700 |
| 10500 | 0.063500 |
| 11000 | 0.044400 |
| 11500 | 0.058300 |
| 12000 | 0.055000 |
| 12500 | 0.040100 |
| 13000 | 0.033100 |
| 13500 | 0.032600 |
| 14000 | 0.040500 |
| 14500 | 0.020700 |
| 15000 | 0.024100 |
| 15500 | 0.026300 |
| 16000 | 0.012700 |
| 16500 | 0.016100 |

```
        17000        0.020200
```

████████████████████ [194/194 00:03]
{'eval_loss': 1.321150779724121,
 'eval_accuracy': 0.8275193798449613,
 'eval_precision': 0.625,
 'eval_recall': 0.5539358600583091,
 'eval_f1': 0.5873261205564142,
 'eval_runtime': 3.1807,
 'eval_samples_per_second': 486.685,
 'eval_steps_per_second': 60.993,
 'epoch': 10.0}

```python
 1 # Running model in-domain on Twitter dataset
 2 # Load tokenizer
 3 tokenizer = BertTokenizer.from_pretrained("bert-base-uncased")
 4
 5 # Load test CSV
 6 df = pd.read_csv("/content/drive/MyDrive/train.csv")
 7
 8 # Batch size for prediction
 9 batch_size = 32
10
11 # Get the device
12 device = "cuda" if torch.cuda.is_available() else "cpu"
13 model.to(device)  # Ensure model is on the correct device
14
15 # Store predictions
16 all_preds = []
17 all_probs = []
18
19 # Iterate over the data in batches
20 for i in range(0, len(df), batch_size):
21     batch = df[i : i + batch_size]
22
23     test_encodings = tokenizer(
24         batch["text"].tolist(),
25         truncation=True,
26         padding=True,
27         max_length=128,
28         return_tensors="pt",
29     ).to(device)
30
31     model.eval()
32     with torch.no_grad():
33         outputs = model(**test_encodings)
34         probs = torch.softmax(outputs.logits, dim=1)
35         preds = torch.argmax(probs, dim=1)
36
37     all_preds.extend(preds.cpu().numpy())
38     all_probs.extend(probs.cpu().numpy().tolist())
39
40 # Add predictions to DataFrame
41 df["predicted_label"] = all_preds
42 df["predicted_prob"] = all_probs
43
44 # Calculate accuracy
45 correct = (df["predicted_label"] == df["hate_label"]).sum()
46 accuracy = correct / len(df)
47
48 # Print results
49 print(df[["text", "hate_label", "predicted_label"]].head())
50 print(f"Accuracy: {accuracy:.4f}")
51
```

```
                                       text  hate_label  \
0  The trans women reading this tweet right now i...           0
1  9) uhhhh i like being lgbt a lot. i feel proud...           0
2  @terryelaineh1 @UKLabour Why do 3.8 million #5...           0
3  I said it yesterday, I knew this is about to g...           0
4  White Small Little Invisible Clits Are A Disgr...           1

   predicted_label
0                0
1                0
2                0
3                0
4                1
Accuracy: 0.9899
```

```python
 1  # Running model out of domain on Reddit dataset
 2
 3  # Load tokenizer
 4  tokenizer = BertTokenizer.from_pretrained("bert-base-uncased")
 5
 6  # Load test CSV
 7  df = pd.read_csv("/content/drive/MyDrive/test_reddit.csv")
 8
 9  # Get the device
10  device = "cuda" if torch.cuda.is_available() else "cpu"
11  model.to(device)  # Ensure model is on the correct device
12
13  # Batch size for prediction
14  batch_size = 32
15
16  # Store predictions
17  all_preds = []
18  all_probs = []
19
20  # Iterate over the data in batches
21  for i in range(0, len(df), batch_size):
22      batch = df[i : i + batch_size]
23
24      test_encodings = tokenizer(
25          batch["text"].tolist(),
26          truncation=True,
27          padding=True,
28          max_length=128,
29          return_tensors="pt",
30      ).to(device)
```

```
31
32        model.eval()
33        with torch.no_grad():
34            outputs = model(**test_encodings)
35            probs = torch.softmax(outputs.logits, dim=1)
36            preds = torch.argmax(probs, dim=1)
37
38        all_preds.extend(preds.cpu().numpy())
39        all_probs.extend(probs.cpu().numpy().tolist())
40
41    # Add predictions to DataFrame
42    df["predicted_label"] = all_preds
43    df["predicted_prob"] = all_probs
44
45    # Calculate accuracy
46    correct = (df["predicted_label"] == df["hate_label"]).sum()
47    accuracy = correct / len(df)
48
49    # Print results
50    print(df[["text", "hate_label", "predicted_label"]].head())
51    print(f"Accuracy: {accuracy:.4f}")
52
```

```
                                              text  hate_label  \
0  For starters bend over the one in pink and kic...           1
1  Sounds like the kinda wholsesome life I'd die ...           0
2  Who the fuck is this insignificant simple mind...           0
3         Fuck off you insufferable retarded faggot.           1
4  Worthless whore, these tits with look nice wit...           0

   predicted_label
0                0
1                0
2                1
3                1
4                1
Accuracy: 0.7957
```

```
 1 # Running model out of domain on Youtube dataset
 2
 3 # Load tokenizer
 4 tokenizer = BertTokenizer.from_pretrained("bert-base-uncased")
 5
 6 # Load test CSV
 7 df = pd.read_csv("/content/drive/MyDrive/test_youtube.csv")
 8
 9 # Batch size for prediction
10 batch_size = 32
11
12 # Store predictions
13 all_preds = []
14 all_probs = []
15
16 # Iterate over the data in batches
17 for i in range(0, len(df), batch_size):
18     batch = df[i : i + batch_size]
19
20     test_encodings = tokenizer(
21         batch["text"].tolist(),
22         truncation=True,
23         padding=True,
24         max_length=128,
25         return_tensors="pt",
26     ).to(device)
27
28     model.eval()
29     with torch.no_grad():
30         outputs = model(**test_encodings)
31         probs = torch.softmax(outputs.logits, dim=1)
32         preds = torch.argmax(probs, dim=1)
33
34     all_preds.extend(preds.cpu().numpy())
35     all_probs.extend(probs.cpu().numpy().tolist())
36
37 # Add predictions to DataFrame
38 df["predicted_label"] = all_preds
39 df["predicted_prob"] = all_probs
40
41 # Calculate accuracy
42 correct = (df["predicted_label"] == df["hate_label"]).sum()
43 accuracy = correct / len(df)
44
45 # Print results
46 print(df[["text", "hate_label", "predicted_label"]].head())
47 print(f"Accuracy: {accuracy:.4f}")
48
```

```
                                              text  hate_label  \
0  Yes indeed. She sort of reminds me of the elde...           0
1  Question: These 4 broads who criticize America...           0
2  It is about time for all illegals to go back t...           0
3  OMG! The EGO's of these young, young, inexperi...           0
4  Joshua Lelo so you have seen all actors from e...           0

   predicted_label
0                0
1                0
2                0
3                0
4                0
Accuracy: 0.7754
```