

In [1]:

```
1 import os
2 import torch
3 import torch.nn as nn
4 import torch.optim as optim
5 import pandas as pd
6 import numpy as np
7 from torch.utils.data import Dataset, DataLoader
8 from sklearn.model_selection import train_test_split, StratifiedKFold
9 from sklearn.metrics import confusion_matrix, classification_report
10 from gensim.models import KeyedVectors
11 import matplotlib.pyplot as plt
```

In [2]:

```
1 # Load pre-trained word2vec embeddings
2 embeddings_path = 'GoogleNews-vectors-negative300.bin.gz'
3 word_vectors = KeyedVectors.load_word2vec_format(embeddings_path, binary=True)
```

In [3]:

```
1 with open("train_text.txt", "r", encoding="utf-8") as f:
2     train_text = f.readlines()
3
4 with open("train_labels.txt", "r", encoding="utf-8") as f:
5     train_labels = f.readlines()
6
7 with open("val_text.txt", "r", encoding="utf-8") as f:
8     val_text = f.readlines()
9
10 with open("val_labels.txt", "r", encoding="utf-8") as f:
11     val_labels = f.readlines()
12
13 with open("test_text.txt", "r", encoding="utf-8") as f:
14     test_text = f.readlines()
15
16 with open("test_labels.txt", "r", encoding="utf-8") as f:
17     test_labels = f.readlines()
```

In [4]:

```
1 # Truncate lists to the minimum length
2 min_train_len = min(len(train_text), len(train_labels))
3 min_val_len = min(len(val_text), len(val_labels))
4 min_test_len = min(len(test_text), len(test_labels))
5
6 train_text, train_labels = train_text[:min_train_len], train_labels[:min_train_len]
7 val_text, val_labels = val_text[:min_val_len], val_labels[:min_val_len]
8 test_text, test_labels = test_text[:min_test_len], test_labels[:min_test_len]
```

In [5]:

```
1 # Prepare data
2 data = pd.DataFrame({"text": train_text + val_text + test_text,
3                      "labels": train_labels + val_labels + test_labels})
4
5 # Convert labels to int
6 data["labels"] = data["labels"].astype(int)
```

In [6]:

```
1 def tweet2vec(tweet, word_vectors, max_len=50):
2     embeddings = []
3     for word in tweet.split():
4         if word in word_vectors:
5             embeddings.append(torch.tensor(word_vectors[word]))
6         if len(embeddings) == max_len:
7             break
8     while len(embeddings) < max_len:
9         embeddings.append(torch.zeros(300))
10    return torch.stack(embeddings)
```

In [7]:

```
1 class TweetDataset(Dataset):
2     def __init__(self, data, word_vectors):
3         self.data = data
4         self.word_vectors = word_vectors
5
6     def __len__(self):
7         return len(self.data)
8
9     def __getitem__(self, idx):
10        tweet = self.data.iloc[idx]["text"]
11        label = self.data.iloc[idx]["labels"]
12        embeddings = tweet2vec(tweet, self.word_vectors)
13        return embeddings, label
```

In [8]:

```
1 class TextClassifier(nn.Module):
2     def __init__(self, embedding_dim, hidden_dim, output_dim):
3         super(TextClassifier, self).__init__()
4         self.fc1 = nn.Linear(embedding_dim, hidden_dim)
5         self.relu = nn.ReLU()
6         self.fc2 = nn.Linear(hidden_dim, output_dim)
7         self.softmax = nn.Softmax(dim=1)
8
9     def forward(self, x):
10        embedded = x.mean(dim=1)
11        out = self.fc1(embedded)
12        out = self.relu(out)
13        out = self.fc2(out)
14        out = self.softmax(out)
15        return out
```

In [9]:

```
1 embedding_dim = 300
2 hidden_dim = 128
3 output_dim = 3
4 learning_rate = 0.001
5 epochs = 30
6 batch_size = 64
7 accumulation_steps = 4
8 device = torch.device("cuda" if torch.cuda.is_available() else "cpu")
9
```

In [\*]:

```

1  # Training and cross-validation
2  skf = StratifiedKFold(n_splits=5)
3  for train_index, val_index in skf.split(data["text"], data["labels"]):
4      train_data = data.iloc[train_index]
5      val_data = data.iloc[val_index]
6
7      train_dataset = TweetDataset(train_data, word_vectors)
8      val_dataset = TweetDataset(val_data, word_vectors)
9
10     train_loader = DataLoader(train_dataset, batch_size=batch_size, shuffle=True)
11     val_loader = DataLoader(val_dataset, batch_size=batch_size, shuffle=True)
12     model = TextClassifier(embedding_dim, hidden_dim, output_dim).to(device)
13     criterion = nn.CrossEntropyLoss()
14     optimizer = optim.Adam(model.parameters(), lr=learning_rate)
15     scheduler = optim.lr_scheduler.StepLR(optimizer, step_size=10, gamma=0.1)
16     training_losses = []
17     validation accuracies = []
18     for epoch in range(epochs):
19
20         model.train()
21         running_loss = 0.0
22         optimizer.zero_grad()
23         for i, (embeddings, labels) in enumerate(train_loader):
24             embeddings, labels = embeddings.to(device), labels.to(device)
25             outputs = model(embeddings)
26
27             loss = criterion(outputs, labels.long())
28
29             loss.backward()
30             if (i + 1) % accumulation_steps == 0:
31                 optimizer.step()
32                 optimizer.zero_grad()
33
34             running_loss += loss.item()
35         training_losses.append(running_loss / len(train_loader))
36         scheduler.step()
37         model.eval()
38         correct = 0
39         total = 0
40         with torch.no_grad():
41             for embeddings, labels in val_loader:
42                 embeddings, labels = embeddings.to(device), labels.to(device)
43                 outputs = model(embeddings.float())
44
45                 _, predicted = torch.max(outputs.data, 1)
46                 total += labels.size(0)
47                 correct += (predicted == labels).sum().item()
48
49         accuracy = 100 * correct / total
50         validation accuracies.append(accuracy)
51         if (epoch + 1) % 10 == 0:
52             print(f"Epoch: {epoch + 1}, Loss: {running_loss / len(train_loader):.4f}")

```

Epoch: 10, Loss: 0.9475, Validation accuracy: 58.00%  
Epoch: 20, Loss: 0.9456, Validation accuracy: 57.95%  
Epoch: 30, Loss: 0.9454, Validation accuracy: 58.04%  
Epoch: 10, Loss: 0.9496, Validation accuracy: 57.88%  
Epoch: 20, Loss: 0.9477, Validation accuracy: 58.09%  
Epoch: 30, Loss: 0.9475, Validation accuracy: 58.16%

In [11]:

```
1 class TweetDataset(Dataset):
2     def __init__(self, data, word_vectors):
3         self.data = data
4         self.word_vectors = word_vectors
5
6     def __len__(self):
7         return len(self.data)
8
9     def __getitem__(self, idx):
10        tweet = self.data.iloc[idx]["text"]
11        label = self.data.iloc[idx]["labels"]
12        embeddings = tweet2vec(tweet, self.word_vectors)
13        return embeddings, torch.tensor(label)
```

In [12]:

```
1 test_data = pd.DataFrame({"text": test_text, "labels": [int(label) for label in test_text["labels"]]})
2 test_dataset = TweetDataset(test_data, word_vectors)
3 test_loader = DataLoader(test_dataset, batch_size=batch_size, shuffle=True)
```

In [13]:

```
1 # Evaluation
2 model.eval()
3 predictions = []
4 true_labels = []
5 with torch.no_grad():
6     for embeddings, labels in test_loader:
7         embeddings, labels = embeddings.to(device), labels.to(device)
8         outputs = model(embeddings.float())
9         _, predicted = torch.max(outputs.data, 1)
10        predictions.extend(predicted.cpu().numpy())
11        true_labels.extend(labels.cpu().numpy())
```

In [14]:

```

1 # Confusion matrix and classification report
2 print("Confusion Matrix:")
3 print(confusion_matrix(true_labels, predictions))
4 print("\nClassification Report:")
5 print(classification_report(true_labels, predictions, zero_division=1))

```

Confusion Matrix:

```

[[ 0 3596  376]
 [ 0 5219  718]
 [ 0 1121 1254]]

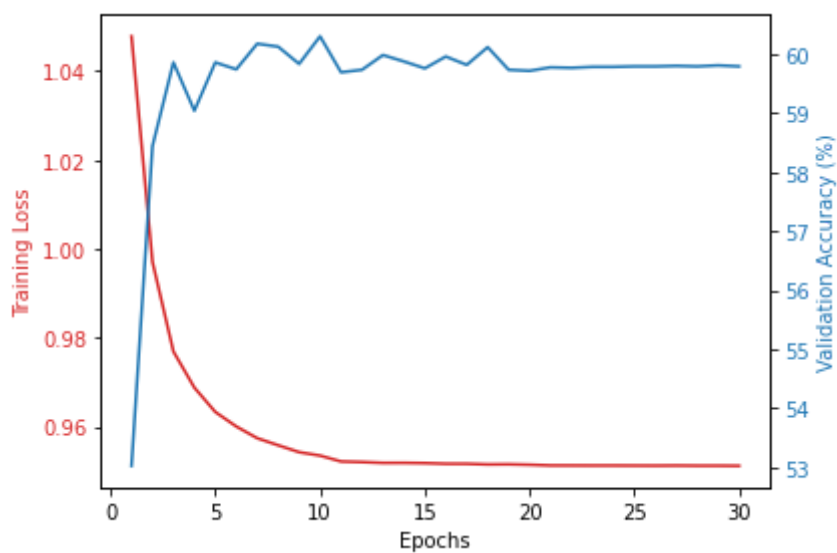
```

Classification Report:

	precision	recall	f1-score	support
0	1.00	0.00	0.00	3972
1	0.53	0.88	0.66	5937
2	0.53	0.53	0.53	2375
accuracy			0.53	12284
macro avg	0.69	0.47	0.40	12284
weighted avg	0.68	0.53	0.42	12284

In [15]:

```
1 fig, ax1 = plt.subplots()
2
3 color = 'tab:red'
4 ax1.set_xlabel('Epochs')
5 ax1.set_ylabel('Training Loss', color=color)
6 ax1.plot(range(1, epochs + 1), training_losses, color=color)
7 ax1.tick_params(axis='y', labelcolor=color)
8
9 ax2 = ax1.twinx()
10
11 color = 'tab:blue'
12 ax2.set_ylabel('Validation Accuracy (%)', color=color)
13 ax2.plot(range(1, epochs + 1), validation_accuracies, color=color)
14 ax2.tick_params(axis='y', labelcolor=color)
15
16 fig.tight_layout()
17 plt.show()
```



In [ ]:

1