

AMS691.02: Natural Language Processing – Fall 2024

Assignment 3

Vishesh Kumar
SBU ID: 115972814

Introduction

In this assignment, we utilize BERT features to classify DBPedia articles into categories. We implement a classifier that leverages pre-trained BERT embeddings and fine-tune it for our classification task. The dataset is pre-processed, and we use PyTorch along with the HuggingFace Transformers library for implementation.

Setup

We begin by installing the necessary libraries and importing modules.

```
1 !pip install tqdm boto3 requests regex sentencepiece sacremoses
2 !pip install transformers
```

Listing 1: Installing Libraries

```
1 import collections
2 import json
3 import torch
4 import torch.nn as nn
5 import tqdm
6 from torch.utils.data import Dataset, DataLoader
7 from transformers import AutoTokenizer, AutoModel
8 import numpy as np
9 import pandas as pd
10 import random
```

Listing 2: Importing Modules

Data Preparation

We define a custom dataset class for DBPedia and construct data loaders.

```
1 SPLITS = ['train', 'dev', 'test']
2
3 class DBPediaDataset(Dataset):
4     '''DBPedia dataset.'''
5     def __init__(self, path):
6         with open(path) as fin:
7             self._data = [json.loads(l) for l in fin]
8             self._n_classes = len(set([datum['label'] for datum in self._data
9                                     ]))
10
11     def __getitem__(self, index):
12         return self._data[index]
13
14     def __len__(self):
15         return len(self._data)
16
17     @property
18     def n_classes(self):
19         return self._n_classes
20
21     @staticmethod
22     def collate_fn(tokenizer, device, batch):
23         '''Collate function for batching.'''
24         labels = torch.tensor([datum['label'] for datum in batch]).long().
25             to(device)
26         sentences = tokenizer(
27             [datum['sentence'] for datum in batch],
28             return_tensors='pt',
29             padding=True)
30         for key in sentences:
31             sentences[key] = sentences[key].to(device)
32         return labels, sentences
33
34 def construct_datasets(prefix, batch_size, tokenizer, device):
35     '''Constructs datasets and data loaders.'''
36     datasets = collections.defaultdict()
37     dataloaders = collections.defaultdict()
38     for split in SPLITS:
39         datasets[split] = DBPediaDataset(f'{prefix}{split}.json')
40         dataloaders[split] = DataLoader(
41             datasets[split],
42             batch_size=batch_size,
43             shuffle=(split == 'train'),
44             collate_fn=lambda x:DBPediaDataset.collate_fn(tokenizer,
```

```
43         device, x))  
    return datasets, dataloaders
```

Listing 3: DBPedia Dataset Class

Problem 1.1: Implementing the Classifier

We implement a simple classifier with one hidden layer.

```
1 class Classifier(nn.Module):  
2     def __init__(self, input_size, hidden_size, num_classes):  
3         super(Classifier, self).__init__()  
4         self.fc1 = nn.Linear(input_size, hidden_size)  
5         self.relu = nn.ReLU()  
6         self.fc2 = nn.Linear(hidden_size, num_classes)  
7  
8     def forward(self, x):  
9         x = self.fc1(x)  
10        x = self.relu(x)  
11        x = self.fc2(x)  
12        return x
```

Listing 4: Classifier Implementation

Training and Evaluation

We set hyperparameters and initialize the tokenizer and BERT model.

```
1 batch_size = 32  
2 classifier_hidden_size = 32  
3  
4 tokenizer = AutoTokenizer.from_pretrained('bert-base-cased')  
5 bert_model = AutoModel.from_pretrained('bert-base-cased')  
6 if torch.cuda.is_available():  
7     bert_model = bert_model.cuda()  
8  
9 datasets, dataloaders = construct_datasets(  
10     prefix='dbpedia_',  
11     batch_size=batch_size,  
12     tokenizer=tokenizer,  
13     device=bert_model.device)
```

Listing 5: Hyperparameters and Initialization

We define functions for training and evaluation.

```
1 import random
2 import numpy as np
3 from sklearn.metrics import accuracy_score
4
5 def set_seed(seed):
6     random.seed(seed)
7     np.random.seed(seed)
8     torch.manual_seed(seed)
9     if torch.cuda.is_available():
10         torch.cuda.manual_seed_all(seed)
```

Listing 6: Training and Evaluation Functions

We train the classifier using different random seeds and evaluate its performance.

```
1 dev_accuracies = []
2 test_accuracies = []
3 seed_values = [42, 123, 999, 2021, 7]
4 best_dev_accuracy = 0.0
5 best_model_state = None
6
7 for seed in seed_values:
8     set_seed(seed)
9     classifier = Classifier(
10         input_size=bert_model.config.hidden_size,
11         hidden_size=classifier_hidden_size,
12         num_classes=datasets['train'].n_classes).to(bert_model.device)
13     optimizer = torch.optim.Adam(classifier.parameters(), lr=5e-4)
14     loss_func = nn.CrossEntropyLoss()
15     pbar = tqdm.tqdm(dataloaders['train'])
16     for labels, sentences in pbar:
17         with torch.no_grad():
18             unpooled_features = bert_model(**sentences)['last_hidden_state']
19             cls_embeddings = unpooled_features[:, 0, :]
20             logits = classifier(cls_embeddings)
21             loss = loss_func(logits, labels)
22             optimizer.zero_grad()
23             loss.backward()
24             optimizer.step()
```

```
25     pbar.set_description(f"Seed: {seed} | Loss: {loss.item():.4f}")
26     # Evaluation on development set
27     classifier.eval()
28     all_preds = []
29     all_labels = []
30     with torch.no_grad():
31         for labels, sentences in dataloaders['dev']:
32             unpooled_features = bert_model(**sentences)['last_hidden_state']
33             cls_embeddings = unpooled_features[:, 0, :]
34             logits = classifier(cls_embeddings)
35             preds = torch.argmax(logits, dim=1)
36             all_preds.extend(preds.cpu().numpy())
37             all_labels.extend(labels.cpu().numpy())
38     dev_accuracy = accuracy_score(all_labels, all_preds)
39     dev accuracies.append(dev_accuracy)
40     print(f"Seed: {seed} | Dev Accuracy: {dev_accuracy:.4f}")
41     if dev_accuracy > best_dev_accuracy:
42         best_dev_accuracy = dev_accuracy
43         best_model_state = classifier.state_dict()
44         best_seed = seed
```

Listing 7: Training Loop

Results

We compute the mean and standard deviation of development accuracies and evaluate the best model on the test set.

```
1 mean_dev_accuracy = np.mean(dev accuracies)
2 std_dev_accuracy = np.std(dev accuracies)
3 print(f"\nMean Dev Accuracy: {mean_dev_accuracy:.4f}")
4 print(f"Std Dev Accuracy: {std_dev_accuracy:.4f}")
5
6 # Evaluate on test set
7 classifier.load_state_dict(best_model_state)
8 classifier.eval()
9 all_preds = []
10 all_labels = []
11 with torch.no_grad():
12     for labels, sentences in dataloaders['test']:
13         unpooled_features = bert_model(**sentences)['last_hidden_state']
14         cls_embeddings = unpooled_features[:, 0, :]
```

```
15         logits = classifier(cls_embeddings)
16         preds = torch.argmax(logits, dim=1)
17         all_preds.extend(preds.cpu().numpy())
18         all_labels.extend(labels.cpu().numpy())
19 test_accuracy = accuracy_score(all_labels, all_preds)
20 print(f"\nBest Model Seed: {best_seed} | Test Accuracy: {test_accuracy:.4f}
      }")
```

Listing 8: Results

Outputs

Seed: 42 | Dev Accuracy: 0.9620
Seed: 123 | Dev Accuracy: 0.9620
Seed: 999 | Dev Accuracy: 0.9690
Seed: 2021 | Dev Accuracy: 0.9610
Seed: 7 | Dev Accuracy: 0.9080

Mean Dev Accuracy: 0.9524
Std Dev Accuracy: 0.0224

Best Model Seed: 999 | Test Accuracy: 0.9740

Problem 1.2: Mean and Max Pooling

We modify the classifier to use mean and max pooling over content tokens.

```
1 class Classifier(nn.Module):
2     def __init__(self, input_size, hidden_size, num_classes):
3         super(Classifier, self).__init__()
4         self.fc1 = nn.Linear(input_size * 2, hidden_size) # Input size
5                 doubled
6         self.relu = nn.ReLU()
7         self.fc2 = nn.Linear(hidden_size, num_classes)
8
9     def forward(self, x):
10         x = self.fc1(x)
11         x = self.relu(x)
12         x = self.fc2(x)
```

```
12         return x
```

Listing 9: Modified Classifier for Mean and Max Pooling

We adjust the input size since we concatenate mean-pooled and max-pooled vectors.

```
1 for labels, sentences in pbar:
2     with torch.no_grad():
3         outputs = bert_model(**sentences)
4         unpooled_features = outputs['last_hidden_state']
5         attention_mask = sentences['attention_mask']
6         mask_expanded = attention_mask.unsqueeze(-1).expand(
7             unpooled_features.size()).float()
8         masked_embeddings = unpooled_features * mask_expanded
9         sum_embeddings = torch.sum(masked_embeddings, dim=1)
10        sum_mask = torch.clamp(mask_expanded.sum(dim=1), min=1e-9)
11        mean_pooled = sum_embeddings / sum_mask
12        masked_embeddings[mask_expanded == 0] = -1e9
13        max_pooled = torch.max(masked_embeddings, dim=1)[0]
14        pooled_features = torch.cat((mean_pooled, max_pooled), dim=1)
15        logits = classifier(pooled_features)
16        loss = loss_func(logits, labels)
17        optimizer.zero_grad()
18        loss.backward()
19        optimizer.step()
20        pbar.set_description(f"Seed: {seed} | Loss: {loss.item():.4f}")
```

Listing 10: Training with Mean and Max Pooling

Results

Seed: 42 | Dev Accuracy: 0.9360

Seed: 123 | Dev Accuracy: 0.9430

Seed: 999 | Dev Accuracy: 0.9140

Seed: 2021 | Dev Accuracy: 0.9310

Seed: 7 | Dev Accuracy: 0.9080

Mean Dev Accuracy: 0.9264

Std Dev Accuracy: 0.0133

Best Model Seed: 123 | Test Accuracy: 0.9320

Problem 1.3: Fine-tuning Last Two Layers of BERT

We fine-tune the last two layers of BERT along with the classifier.

```
1 classifier = Classifier(  
2     bert_model.config.hidden_size,  
3     classifier_hidden_size,  
4     datasets['train'].n_classes).to(bert_model.device)  
5  
6 for name, param in bert_model.named_parameters():  
7     if name.startswith('encoder.layer.10') or name.startswith('encoder.  
8         layer.11'):  
9         param.requires_grad = True  
10    else:  
11        param.requires_grad = False  
12  
13 params_to_optimize = list(classifier.parameters()) + [param for param in  
14     bert_model.parameters() if param.requires_grad]  
15 optimizer = torch.optim.Adam(params_to_optimize, lr=5e-5)
```

Listing 11: Fine-tuning Setup

Training

We train the model without the `torch.no_grad()` context since we are updating BERT parameters.

```
1 for labels, sentences in pbar:  
2     outputs = bert_model(**sentences)  
3     unpooled_features = outputs['last_hidden_state']  
4     cls_features = unpooled_features[:, 0, :]  
5     logits = classifier(cls_features)  
6     loss = loss_func(logits, labels)  
7     optimizer.zero_grad()  
8     loss.backward()  
9     optimizer.step()  
10    pbar.set_description(f"Loss: {loss.item():.4f}")
```

Listing 12: Training Loop with Fine-tuning

Results

Loss: 0.0640: 100%|| 313/313 [00:55<00:00, 5.63it/s]

Problem 1.4: Evaluation with Fine-tuning

We evaluate the fine-tuned model using different seeds.

Seed: 42 | Dev Accuracy: 0.9960

Seed: 123 | Dev Accuracy: 0.9940

Seed: 999 | Dev Accuracy: 0.9960

Seed: 2021 | Dev Accuracy: 0.9940

Seed: 7 | Dev Accuracy: 0.9950

Mean Dev Accuracy: 0.9950

Std Dev Accuracy: 0.0009

Best Model Seed: 42 | Test Accuracy: 0.9970

Problem 1.5: Replacing BERT with GPT-2

We replace BERT with GPT-2 and adjust the code accordingly.

```
1 from transformers import GPT2Tokenizer, GPT2Model
2
3 tokenizer = GPT2Tokenizer.from_pretrained('gpt2')
4 tokenizer.pad_token = tokenizer.eos_token
5
6 gpt2_model = GPT2Model.from_pretrained('gpt2')
7 if torch.cuda.is_available():
8     gpt2_model = gpt2_model.cuda()
```

Listing 13: Using GPT-2 Model

We adjust the feature extraction method since GPT-2 is autoregressive.

```
1 with torch.no_grad():
2     outputs = gpt2_model(**sentences)
```

```
3     hidden_states = outputs['last_hidden_state']
4     attention_mask = sentences['attention_mask'].unsqueeze(-1).expand(
        hidden_states.size()).float()
5     masked_hidden_states = hidden_states * attention_mask
6     sum_hidden_states = torch.sum(masked_hidden_states, dim=1)
7     sum_mask = attention_mask.sum(dim=1).clamp(min=1e-9)
8     features = sum_hidden_states / sum_mask
```

Listing 14: Feature Extraction with GPT-2

Results

Seed: 42 | Dev Accuracy: 0.8960
Seed: 123 | Dev Accuracy: 0.8300
Seed: 999 | Dev Accuracy: 0.8910
Seed: 2021 | Dev Accuracy: 0.8920
Seed: 7 | Dev Accuracy: 0.8650

Mean Dev Accuracy: 0.8748
Std Dev Accuracy: 0.0249

Best Model Seed: 42 | Test Accuracy: 0.8900

Conclusion

We observe that fine-tuning the last two layers of BERT significantly improves the model's performance. Replacing BERT with GPT-2 yields lower accuracy, likely due to the autoregressive nature of GPT-2 and its suitability for generation tasks rather than classification.