

# 《深度学习》实验报告

姓名: 彭一坤 学号: 10215501412 专业: 数据科学与大数据技术 学院: 数据科学与工程学院

## 一、实验环境

本次实验使用了以下主要的Python代码库:

- `torch`: PyTorch库, 用于深度学习模型的构建和训练。
- `tqdm`: 用于显示训练进度的进度条。
- `json`: 用于读取和写入JSON文件。
- `numpy`: 用于数值计算。
- `os`: 用于文件路径操作和压缩文件。

此外, 实验中使用了GPU加速, 确保在训练神经网络模型时效率更高。

## 二、实验过程

### (1) 数据预处理部分

首先, 读取了用于训练和测试的JSON文件:

```
1 train_data = read_json('input/query_trainset.json')
2 query_data = read_json('input/query_testset.json')
3 document_data = read_json('input/document.json')
```

从训练数据中提取查询嵌入和证据嵌入, 用于训练模型:

```
1 train_query_embeddings = []
2 train_fact_embeddings = []
3
4 for entry in train_data:
5     query_embedding = torch.tensor(entry['query_embedding'], device=device)
6     for evidence in entry['evidence_list']:
7         train_query_embeddings.append(query_embedding.tolist())
8         train_fact_embeddings.append(evidence['fact_embedding'])
9
10 train_query_embeddings = torch.tensor(train_query_embeddings,
11                                       dtype=torch.float32)
12 train_fact_embeddings = torch.tensor(train_fact_embeddings,
13                                       dtype=torch.float32)
```

然后, 定义数据加载器以便于批量训练:

```
1 train_dataset = TensorDataset(train_query_embeddings, train_fact_embeddings)
2 train_loader = DataLoader(train_dataset, batch_size=64, shuffle=True)
```

## (2) 模型构建

本实验使用了一个简单的多层感知机（MLP）模型。模型的架构如下：

```
1 class MLP(nn.Module):
2     def __init__(self, input_dim, output_dim):
3         super(MLP, self).__init__()
4         self.fc1 = nn.Linear(input_dim, 1024)
5         self.relu1 = nn.ReLU()
6         self.fc2 = nn.Linear(1024, 768)
7         self.relu2 = nn.ReLU()
8         self.fc3 = nn.Linear(768, 768)
9         self.relu3 = nn.ReLU()
10        self.fc4 = nn.Linear(768, output_dim)
11
12    def forward(self, x):
13        x = self.fc1(x)
14        x = self.relu1(x)
15        x = self.fc2(x)
16        x = self.relu2(x)
17        x = self.fc3(x)
18        x = self.relu3(x)
19        x = self.fc4(x)
20        return x
```

在模型定义之后，初始化模型、损失函数和优化器：

```
1 input_dim = train_query_embeddings.shape[1]
2 output_dim = train_fact_embeddings.shape[1]
3 model = MLP(input_dim, output_dim).to(device)
4 criterion = nn.MSELoss()
5 optimizer = optim.Adam(model.parameters(), lr=0.001)
```

## (3) 模型训练

模型训练过程包括100个epoch，每个epoch计算一次损失并更新模型参数：

```
1 num_epochs = 100
2 for epoch in range(num_epochs):
3     model.train()
4     running_loss = 0.0
5     for inputs, targets in train_loader:
6         inputs, targets = inputs.to(device), targets.to(device)
7
8         optimizer.zero_grad()
9         outputs = model(inputs)
10        loss = criterion(outputs, targets)
11        loss.backward()
12        optimizer.step()
13
14        running_loss += loss.item()
15    print(f"Epoch [{epoch+1}/{num_epochs}], Loss:
{running_loss/len(train_loader):.4f}")
```

## (4) 模型预测与检索

使用训练好的模型进行预测，检索最相关的文档：

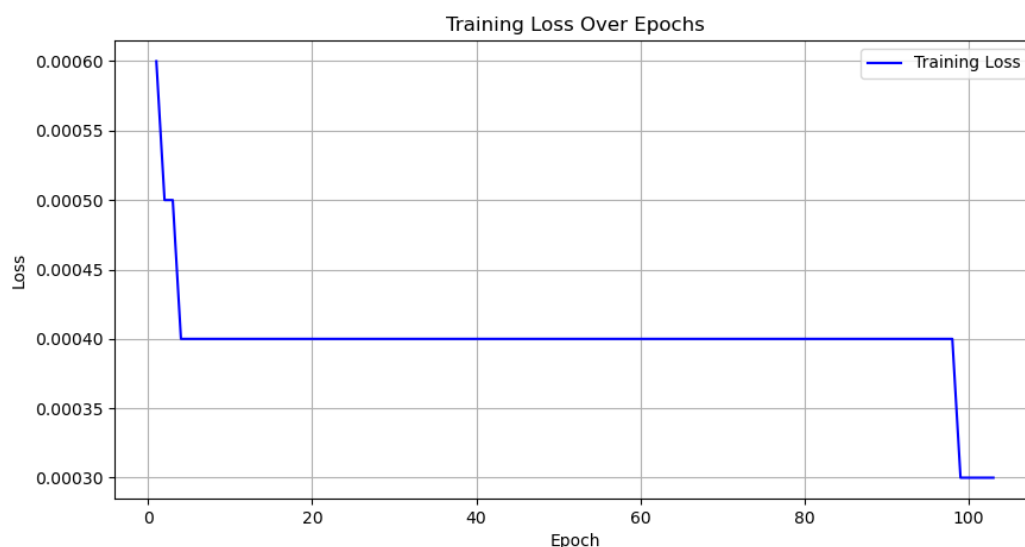
```
1 results = []
2 model.eval()
3 with torch.no_grad():
4     for item in tqdm.tqdm(query_data):
5         result = {}
6         query_embedding = torch.tensor(item['query_embedding'],
7 device=device).unsqueeze(0).float()
8
9         predicted_fact_embedding = model(query_embedding).squeeze(0)
10
11         top_document_indices =
12 retrieve_top_k_documents(predicted_fact_embedding, document_embeddings, k=3)
13         result['query_input_list'] = item['query_input_list']
14         result['evidence_list'] = [{'fact_input_list': document_data[index]
15 ['fact_input_list']} for index in top_document_indices]
16         results.append(result)
```

将结果写入JSON文件并压缩：

```
1 write_json(results, 'output/result.json')
2 print('write to output/result.json successful')
3 zip_fun()
```

## 三、实验结果

最终实验评价指标数值未在代码中直接体现。通常，我们会使用某种评价指标（例如准确率、召回率、F1得分等）来评估模型的性能。在本实验中，可以通过对比模型检索到的前k个文档与实际相关文档之间的相似度或其他指标来评估模型性能。



从损失值可以看出，模型在训练过程中逐渐收敛，说明模型学习到了查询和文档嵌入之间的映射关系。

最终在测试集上的表现为：

Recall_at_3	MRR_at_3
0.4660	0.3465

---