



INTERNATIONAL INSTITUTE OF
INFORMATION TECHNOLOGY

H Y D E R A B A D

INTRODUCTION TO NLP (CS7.401)

Assignment 4

Submitted by:

Bhanuj Gandhi

2022201068

April 27, 2023

Contents

1	ELMo	1
1.1	Introduction	1
1.2	Datasets	2
2	Implementation	3
2.1	Preprocessing	3
2.2	Build the Elmo model	3
2.3	PreTraining the Elmo model	4
2.4	Evaluate the Elmo model	4
3	Visualize the results	6
3.1	SST Dataset	6
3.2	Multi NLI Dataset	7
4	Conclusion	10

1 ELMo

1.1 Introduction

ELMo (Embeddings from Language Models) is a deep contextualized word representation model. Unlike traditional word embedding models, ELMo is capable of capturing the contextual meaning of words by taking into account the entire sentence in which they appear.

The ELMo architecture contains forward and backward language models trained through its Bi-LSTMs. The stacked Bi-LSTM-based language model has a layer of non-contextual word embeddings (like word2vec, or a character convolutional network as in the original ELMo) at the input. While a forward language model can be trained by the next word prediction task, for a backward language model - to put it simply - it's nothing but previous word prediction task. In either of the cases, the LSTM model makes the word prediction based on the context. By context, we mean all the words preceeding (forward LM) or succeeding (backward LM) the word to be predicted in the sentence.

By training these Bi-LSTMs using the language modeling objective, we essentially pretrain the weights of the ELMo architecture, making it suitable for application on downstream tasks. To break this down, when we pretrain the model, we make the model learn the nature of the language itself. The model does not really learn any useful task in this process - instead it captures the general dependencies in the language, which makes it smarter to be able to solve any task thrown at it later quickly. This task is what we refer to as the downstream task. On pretraining the model, we expect it to give good performance on a given NLP task with lesser training and data than it would take without pretraining.

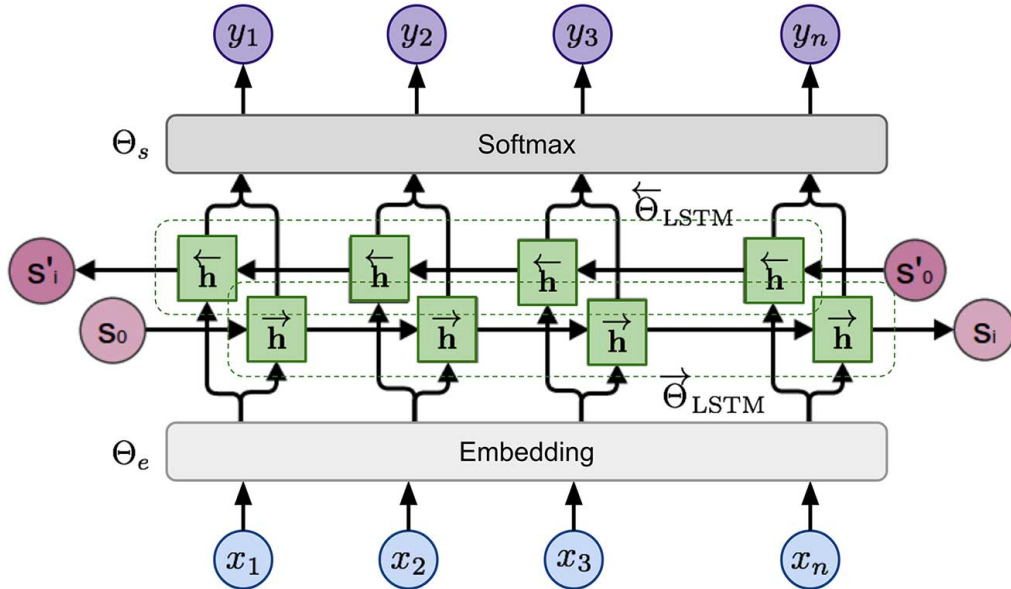


Figure 1: ELMo Architecture

1.2 Datasets

Datasets Used :

- **Stanford Sentiment Treebank:** It contains 8544 train, 1101 validation and 2210 test sentences. Each complete sentence is annotated with a float label that indicates its level of positive sentiment from 0.0 to 1.0. To convert the given task into a binary sentiment classification task, we need to assign a label of either 0 or 1 to the generated elmo embedding using MLP classifier.
- **Multi-Genre NLI Corpus :** MultiNLI is a comprehensive dataset comprising more than 400,000 sentence pairs (hypothesis and premise), each of which is categorized as either entailment, contradiction, or neutral. Similarly in this corpus, we have to classify each elmo embedding of a particular sentence into three categories.

2 Implementation

2.1 Preprocessing

All the preprocessing steps are explained in the code comments below

```
1
2 def preprocess_sentence(sent: str) -> str:
3     # Remove HTML
4     soup = BeautifulSoup(sent, "html.parser")
5     sent = soup.get_text(separator=" ")
6
7     # Remove whitespaces
8     sent = sent.strip()
9     sent = " ".join(sent.split())
10
11    # Lowercase
12    sent = sent.lower()
13
14    # Remove accent characters
15    sent = unidecode.unidecode(sent)
16
17    # Expand the contractions
18    sent = contractions.fix(sent)
19
20    # Remove punctuations and non-ASCII characters
21    sent = re.sub(r"[^\w\s]", "", sent)
22    sent = re.sub(r"[^\x00-\x7f]", "", sent)
23
24    # Remove stopwords and lemmatize
25    sent = " ".join(
26        [lemmatizer.lemmatize(word) for word in sent.split() if word not in
27         stop_words]
28    )
29    return sent
```

2.2 Build the Elmo model

I have used two Bi-LSTM stacked and passed the input. These layers are trained on the prediction of the next word.

```
1
2 class ELMo(nn.Module):
3     def __init__(self, vocab_size, embedding_size, hidden_size, dropout,
4         embeddings):
5         super(ELMo, self).__init__()
6         self.embedding = nn.Embedding.from_pretrained(embeddings)
7
8         self.layer_1 = nn.LSTM(
9             input_size=embedding_size,
10             hidden_size=hidden_size,
11             num_layers=1,
12             batch_first=True,
13             bidirectional=True,
14         )
15
16         self.layer_2 = nn.LSTM(
17             input_size=2 * hidden_size,
```

```

17         hidden_size=hidden_size,
18         num_layers=1,
19         batch_first=True,
20         bidirectional=True,
21     )
22
23     self.dropout = nn.Dropout(dropout)
24
25     self.linear = nn.Linear(hidden_size * 2, vocab_size)
26
27     def forward(self, X):
28         embeddings = self.embedding(X)
29
30         lstm1_output, _ = self.layer_1(embeddings)
31
32         lstm2_output, _ = self.layer_2(lstm1_output)
33         lstm2_output = self.dropout(lstm2_output)
34
35         output = self.linear(lstm2_output)
36         output = torch.transpose(output, 1, 2)
37         return output

```

2.3 PreTraining the Elmo model

```

1
2 class Sentiment_Classifier(nn.Module):
3     def __init__(self, embedding_size):
4         super(Sentiment_Classifier, self).__init__()
5
6         self.s1 = nn.Parameter(torch.ones(1))
7         self.s2 = nn.Parameter(torch.ones(1))
8         self.s3 = nn.Parameter(torch.ones(1))
9         self.alpha = nn.Parameter(torch.ones(1))
10
11         self.linear = nn.Linear(embedding_size, 2)
12
13         self.sigmoid = nn.Sigmoid()
14
15     def forward(self, sentence):
16         embeddings = elmo.embedding(sentence)
17         out_1, _ = elmo.lstm1(embeddings)
18         out_2, _ = elmo.lstm2(out_1)
19
20         s_sum = self.s1 + self.s2 + self.s3
21
22         output = self.alpha * (
23             self.s1 / s_sum * embeddings
24             + self.s2 / s_sum * out_1
25             + self.s3 / s_sum * out_2
26         ).to(torch.float32)
27
28         output = self.linear(output)
29         output = output.mean(dim=1)
30         output = self.sigmoid(output)
31
32         return output

```

2.4 Evaluate the Elmo model

```

1 sentiment_model = Sentiment_Classifier(100).to(device)
2 criterion = nn.CrossEntropyLoss().to(device)
3 optimizer = optim.Adam(sentiment_model.parameters(), lr=0.01)
4
5
6 prev_val = np.inf
7 for epoch in range(20):
8     total_loss = 0
9     total_loss_train = 0
10    for i, (sentence, label) in enumerate(tqdm(train_dataloader)):
11        sentiment_model.train()
12        sentence, label = sentence.to(device), label.to(device)
13
14        optimizer.zero_grad()
15        outputs = sentiment_model(sentence)
16
17        label = (label >= 0.5).long()
18        batch_loss = criterion(outputs, label)
19        total_loss_train += batch_loss.item()
20
21        batch_loss.backward()
22        optimizer.step()
23
24    print(f"Train Loss: {total_loss_train/len(train_dataloader)}")
25    with torch.no_grad():
26        sentiment_model.eval()
27        for inputs, targets in tqdm(test_dataloader):
28            model_input = inputs.to(device)
29            targets = targets.to(device)
30            out = sentiment_model(model_input)
31
32            targets = (targets >= 0.5).long()
33            loss = criterion(out, targets)
34            total_loss += loss.item()
35
36    if prev_val < total_loss:
37        break
38    else:
39        prev_val = total_loss
40
41    print(f"Validation Loss {total_loss/len(val_dataloader)}")

```

3 Visualize the results

3.1 SST Dataset

Classification Report:					
	precision	recall	f1-score	support	
0	0.70	0.78	0.74	1143	
1	0.73	0.64	0.68	1067	
accuracy			0.71	2210	
macro avg	0.72	0.71	0.71	2210	
weighted avg	0.71	0.71	0.71	2210	

Figure 2: Classification Report for SST

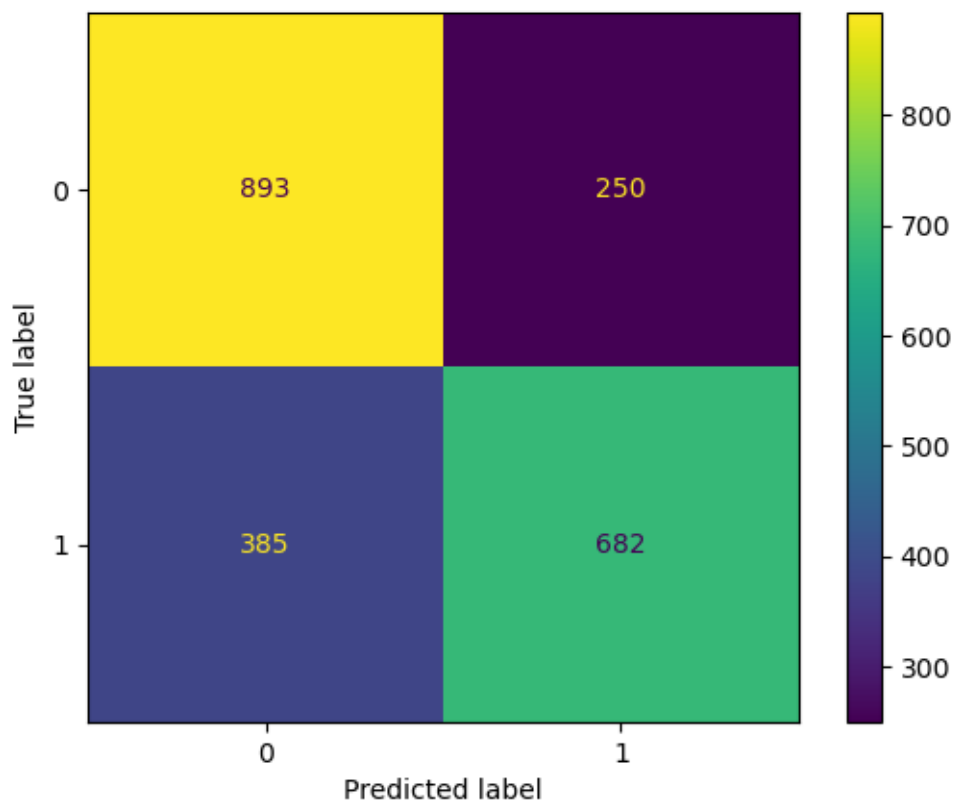


Figure 3: Confusion Matrix for SST

Positive Label - Positive Sentiment
 Negative Label - Negative Sentiment

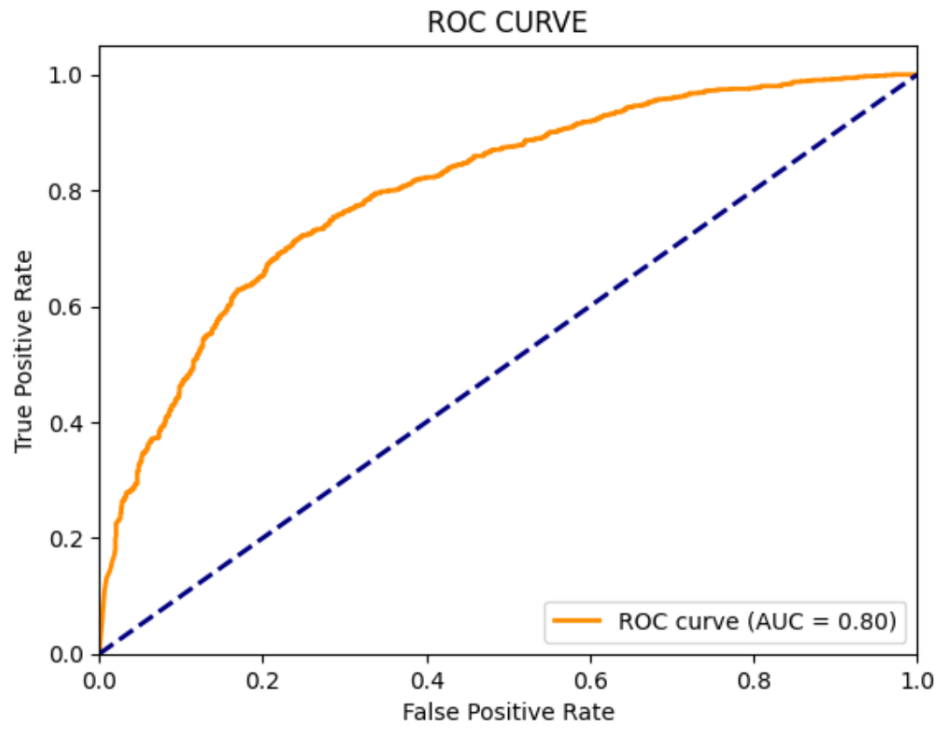


Figure 4: ROC Curve for SST Dataset

3.2 Multi NLI Dataset

Classification Report:					
	precision	recall	f1-score	support	
0	0.46	0.36	0.40	3479	
1	0.42	0.47	0.45	3123	
2	0.48	0.54	0.51	3213	
accuracy			0.46	9815	
macro avg	0.46	0.46	0.45	9815	
weighted avg	0.46	0.46	0.45	9815	

Figure 5: Classification Report for Multi NLI

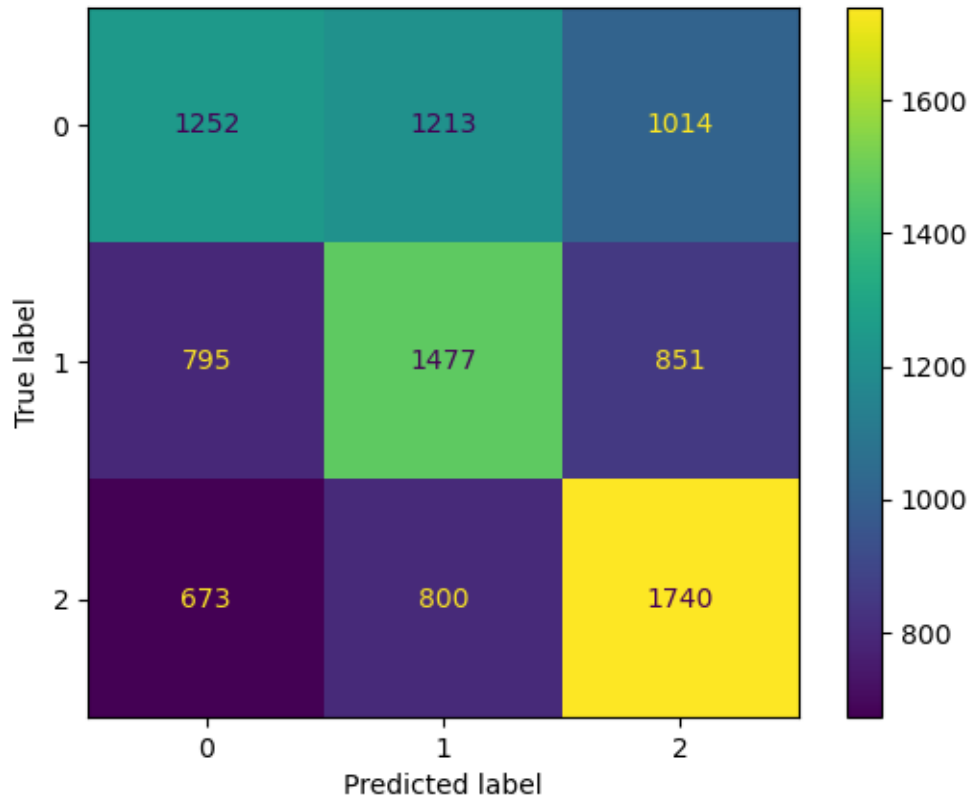


Figure 6: Confusion Matrix for Multi NLI

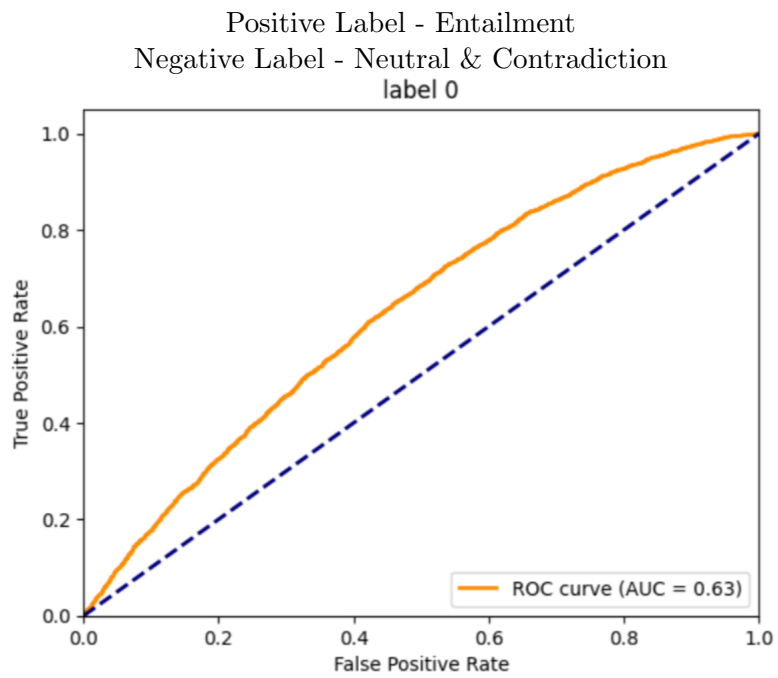


Figure 7: ROC Curve for MultiNLI Dataset

Positive Label - Neutral
Negative Label - Entailment & Contradiction

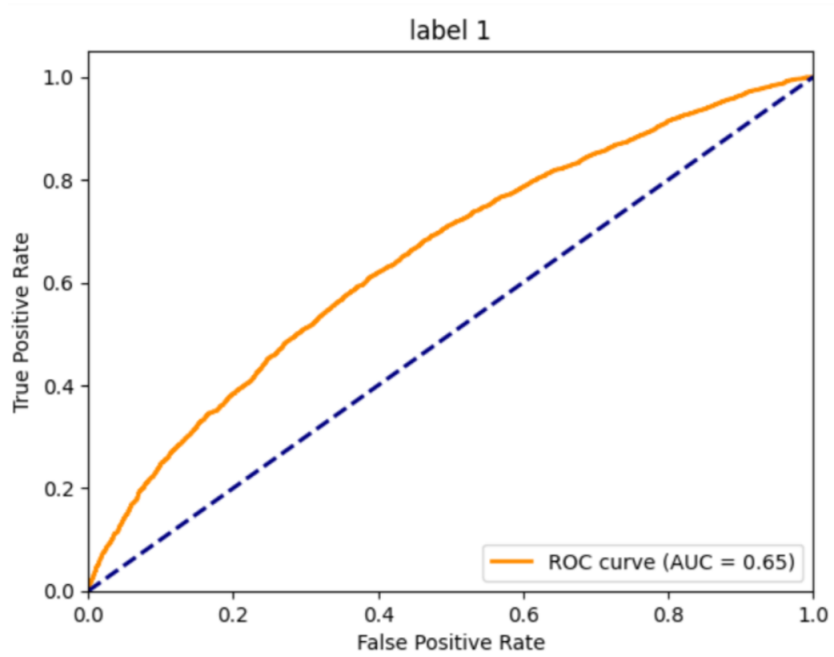


Figure 8: ROC Curve for MultiNLI Dataset

Positive Label - Contradiction
Negative Label - Entailment & Neutral

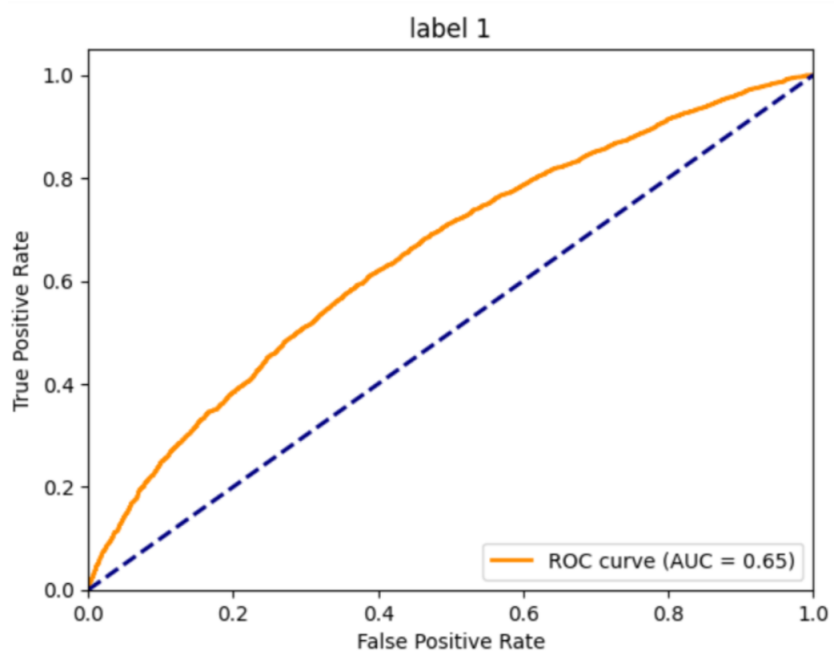


Figure 9: ROC Curve for MultiNLI Dataset

4 Conclusion

In this assignment, we analyze the performance of a deep learning model on two distinct NLP tasks: sentiment analysis and natural language inference (NLI). The aim is to assess the model's ability to accurately classify sentiment in movie reviews and determine the relationships between pairs of sentences. We utilize ELMo embeddings, which capture rich contextual information, to enhance the model's understanding of the text.

The sentiment analysis task involves classifying movie reviews as positive or negative. We evaluate the model's performance on the Stanford Sentiment Treebank (SST) dataset, which provides labeled movie reviews. For the NLI task, we use the Multi-NLI dataset, which requires determining the relationship (entailment, contradiction, or neutral) between two sentences.

Our experiments reveal that the sentiment analysis task achieves a remarkable test accuracy of 71% on the SST dataset. Additionally, the model exhibits an impressive AUC score of 0.80 for the receiver operating characteristic (ROC) curve. These findings demonstrate the model's ability to effectively distinguish between positive and negative movie reviews. The high accuracy and AUC score can be attributed to the clear distinction between the two sentiment classes in the dataset.

In contrast, the NLI task presents a more challenging scenario, where the model struggled to accurately classify the relationships between sentence pairs. The performance on the Multi-NLI dataset yields a lower test accuracy of 46% and an AUC score of 0.68 for the ROC curve. This suggests that the model faced difficulties in capturing the subtle differences that exist between sentences for entailment, contradiction, and neutral relationships. The complexity of the NLI task, coupled with the large number of samples in the dataset, presented a challenging learning environment for the model.

Our analysis highlights the effectiveness of ELMo embeddings in sentiment analysis, where the model excelled in accurately classifying positive and negative movie reviews. The contextual information captured by ELMo embeddings proved instrumental in distinguishing sentiment patterns. However, in the NLI task, the model struggled to achieve comparable performance, indicating the difficulty in accurately determining the relationships between sentences.

Further improvements could be explored, such as utilizing more advanced architectures or incorporating additional contextual information beyond ELMo embeddings. Additionally, exploring larger datasets or fine-tuning the model on domain-specific data might enhance its performance on the NLI task.