



INTERNATIONAL INSTITUTE OF
INFORMATION TECHNOLOGY

H Y D E R A B A D

INTRODUCTION TO NLP (CS7.401)

Assignment 2

Submitted by:

Bhanuj Gandhi

2022201068

March 8, 2023

Contents

1	Neural POS Tagging	1
1.1	Introduction	1
1.2	Dataset Overview	1
2	Methodology	2
3	Results	4

1 Neural POS Tagging

Q) Design, implement, and train a neural sequence model (RNN, LSTM, GRU, etc.) of your choice to (tokenize and) tag a given sentence with the correct part-of-speech tags. For example, given the input

Mary had a little lamb

your model should output

Mary NOUN
had VERB
a DET
little ADJ
lamb NOUN

Note that the part-of-speech tag is separated from each word by a tab `\t` character.

1.1 Introduction

Neural POS tagging is a method for automatically assigning parts-of-speech (POS) tags to words in a sentence using neural networks. POS tagging is an important task in natural language processing (NLP) because it helps in understanding the syntactic structure of sentences and is used in many downstream NLP tasks, such as text classification and named entity recognition.

1.2 Dataset Overview

The dataset used in this task is the *UD English-Atis* treebank, version 2.11, which includes data specific to the ATIS (Airline Travel Information System) domain. The dataset consists of the following files:

- `en_atis-ud-train.conllu`: This file contains the training set with annotated part-of-speech tags and syntactic dependency relations.
- `en_atis-ud-dev.conllu`: This file contains the development set, which is used for tuning the model hyperparameters and evaluating its performance during training.
- `en_atis-ud-test.conllu`: This file contains the test set, which is used to evaluate the final performance of the trained model.

The files are in the CoNLL-U format, which is a plain text format for representing syntactic dependency trees with annotations.

Each line in these files corresponds to a token in a sentence, and the fields in each line are separated by tabs. The first ten fields contain information about the token, including its index, word form, lemma, part-of-speech tag, and features. The eleventh field contains the index of the token's head in the sentence, and the twelfth field contains the syntactic dependency relation between the token and its head.

2 Methodology

Steps followed to implement the Neural POS Tagger

1. Examine the dataset to understand the available data and its structure in detail.
2. Created sequence of the dataset in order to feed the model

```
1     def prepare_datasets(dataset):
2         mod_data = []
3         for idx in range(len(dataset)):
4             tempword = []
5             temptag = []
6             for jdx in range(len(dataset[idx])):
7                 tempword.append(dataset[idx][jdx]["form"])
8                 temptag.append(dataset[idx][jdx]["upos"])
9
10            mod_data.append([tempword, temptag])
11    return mod_data
12
```

3. Create vocabulary for Train Dataset (both for word tokens as well as Tag Tokens), for this I have used `torchtext` vocabulary builder, which lets us created a dictionary which handles the unknowns.

```
1     word_vocab = torchtext.vocab.build_vocab_from_iterator(new_list)
2     word_vocab.insert_token('<unk>', 0)
3     word_vocab.set_default_index(word_vocab['<unk>'])
4
```

4. Create model, `pytorch` let's you create model using class inheritance. I have used 3 layers to define my model.

- (a) Embedding Layer
- (b) Bidirectional LSTM Layer
- (c) Linear Layer

- I have Bidirectional LSTM in POS tagging of English language because in English language, the POS tag of a word can depend on the words that come before and after it. By using a bidirectional LSTM, we can take into account both the preceding and following words when predicting the POS tag for a particular word. This helps to capture the context of the sentence better and leads to more accurate POS tagging.
- I have incorporated Dropout in the model since the corpus being used is small, which increases the risk of overfitting. By dropping some of the weights and re-learning them in subsequent iterations, the model can prevent overfitting.

```
1     class LSTMTagger(nn.Module):
2         def __init__(
3             self,
4             word_embedding_dim,
5             word_hidden_dim,
6             vocab_size,
7             tagset_size,
8         ):
9             super(LSTMTagger, self).__init__()
10            self.word_hidden_dim = word_hidden_dim
```

```

11         self.word_embeddings = nn.Embedding(vocab_size,
word_embedding_dim)
12         self.lstm = nn.LSTM(word_embedding_dim, word_hidden_dim,
num_layers = 1, bidirectional = True)
13
14         self.hidden2tag = nn.Linear(word_hidden_dim*2, tagset_size)
15
16         self.dropout = nn.Dropout(0.1)
17
18     def forward(self, sentence):
19         embeds = self.dropout(self.word_embeddings(sentence))
20         lstm_out, _ = self.lstm(embeds.view(len(sentence), 1, -1))
21         tag_space = self.hidden2tag(lstm_out.view(len(sentence), -1))
22         tag_scores = F.log_softmax(tag_space, dim=1)
23         return tag_scores
24
25

```

5. To optimize the model during training, I have used the *cross entropy loss function* and the *Adam optimizer* to update the model's parameters based on the gradients calculated during backpropagation. The *cross entropy* loss function is used because it is commonly used for multi-class classification tasks. The Adam optimizer is a popular optimization algorithm that adapts the learning rate during training to improve convergence.

```

1     # Define the loss function as the Cross Entropy Loss
2     loss_function = nn.CrossEntropyLoss()
3
4     # We will be using an Adam optimizer
5     optimizer = optim.Adam(model.parameters(), lr=0.005)
6

```

6. I have experimented with different values of *hyperparameters* and have decided on some based on my observations. I noticed that the model performed better with smaller values for the embedding size and hidden layer, which may be due to the small size of the corpus. The following are the finalized *hyperparameters*

```

1     WORD_EMBEDDING_DIM = 64
2     WORD_HIDDEN_DIM = 64
3     EPOCHS = 40
4     BIDIRECTIONAL = True
5     DROPOUT = 0.1
6     LEARNING_RATE = 0.05
7

```

3 Results

After following the approach discussed above, I was able to achieve **98% accuracy on the test dataset**.

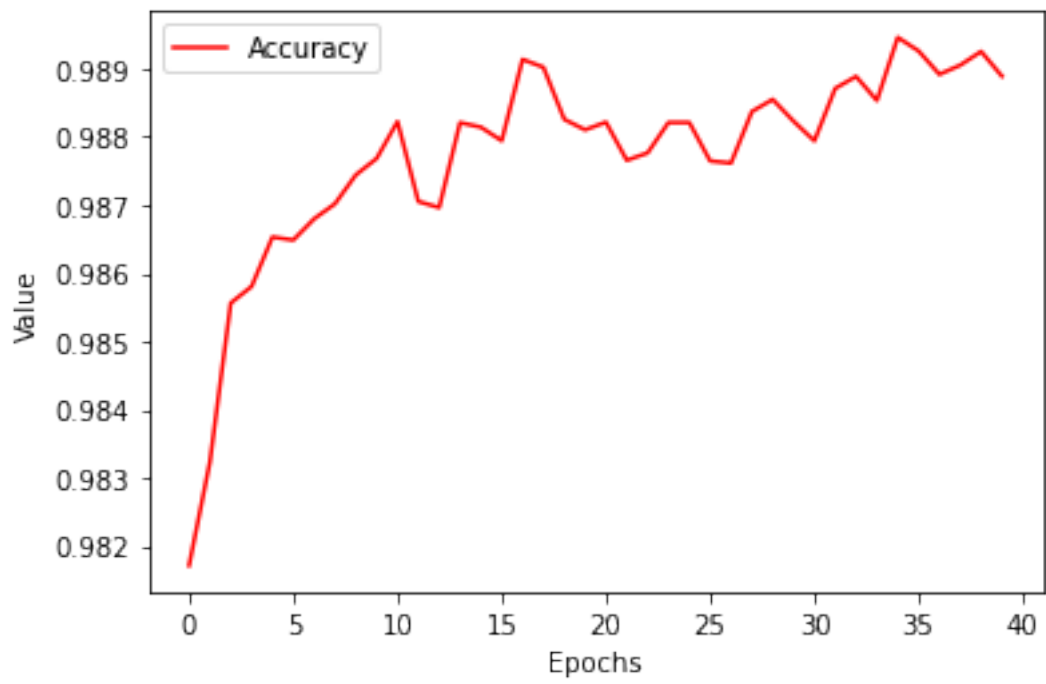
Below is the *classification report* of the trained model

	precision	recall	f1-score	support
ADJ	0.98	0.90	0.94	1632
ADP	0.97	0.99	0.98	10791
ADV	0.88	0.89	0.88	431
AUX	0.98	0.98	0.98	1732
CCONJ	1.00	0.99	0.99	751
DET	0.96	0.99	0.98	3805
INTJ	0.94	0.99	0.96	319
NOUN	0.99	0.99	0.99	8621
NUM	0.99	0.98	0.99	933
PART	0.90	0.99	0.94	366
PRON	1.00	0.96	0.98	3022
PROPN	1.00	1.00	1.00	11657
VERB	0.99	0.93	0.96	4595
accuracy			0.98	48655
macro avg	0.97	0.97	0.97	48655
weighted avg	0.98	0.98	0.98	48655

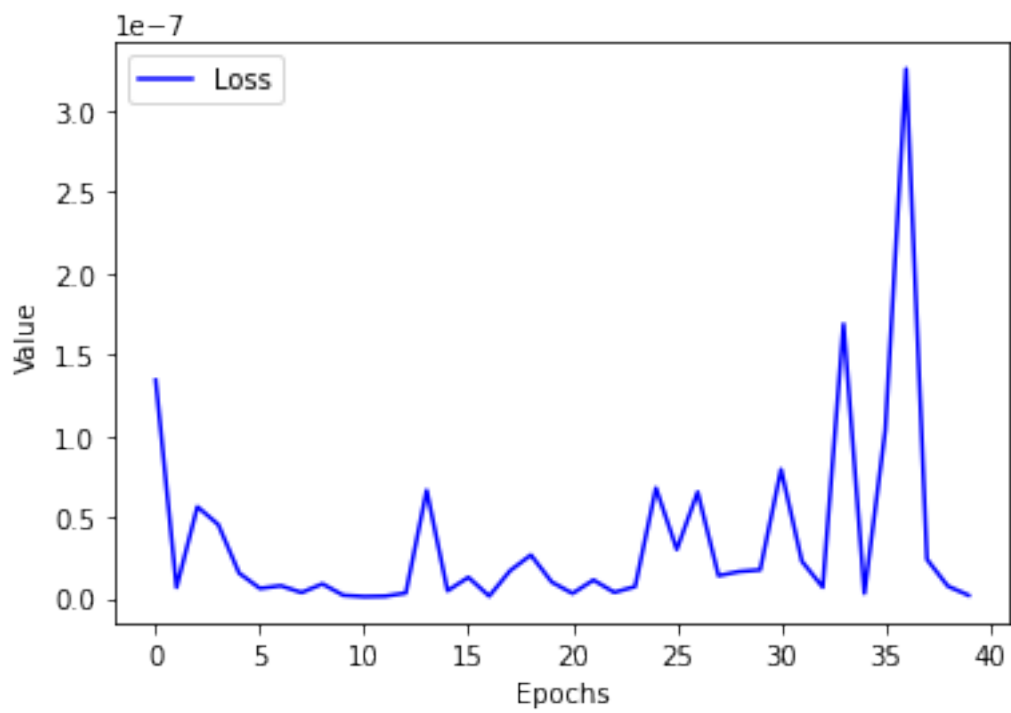
Figure 1: Classification Report

The figures above depict various metrics such as *precision*, *recall*, and *F1-score*, which provide a balance between precision and recall by computing their harmonic mean.

The below plots display the training and validation accuracy, as well as the loss during the training process. These plots provide an analysis of the model's performance with each epoch during training.

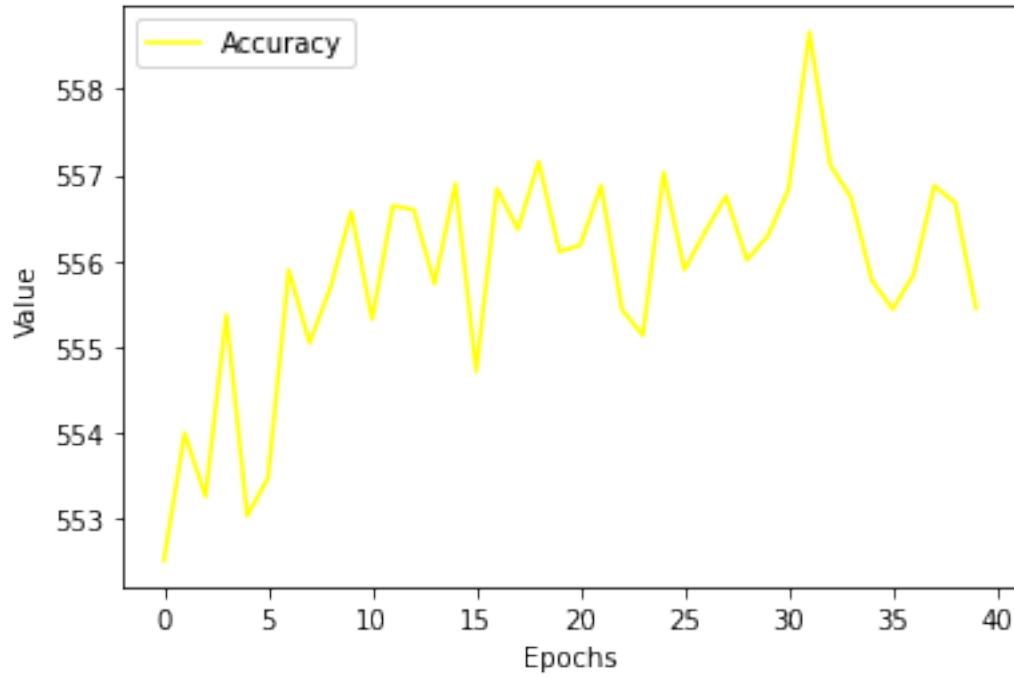


(a) Train Set Accuracy

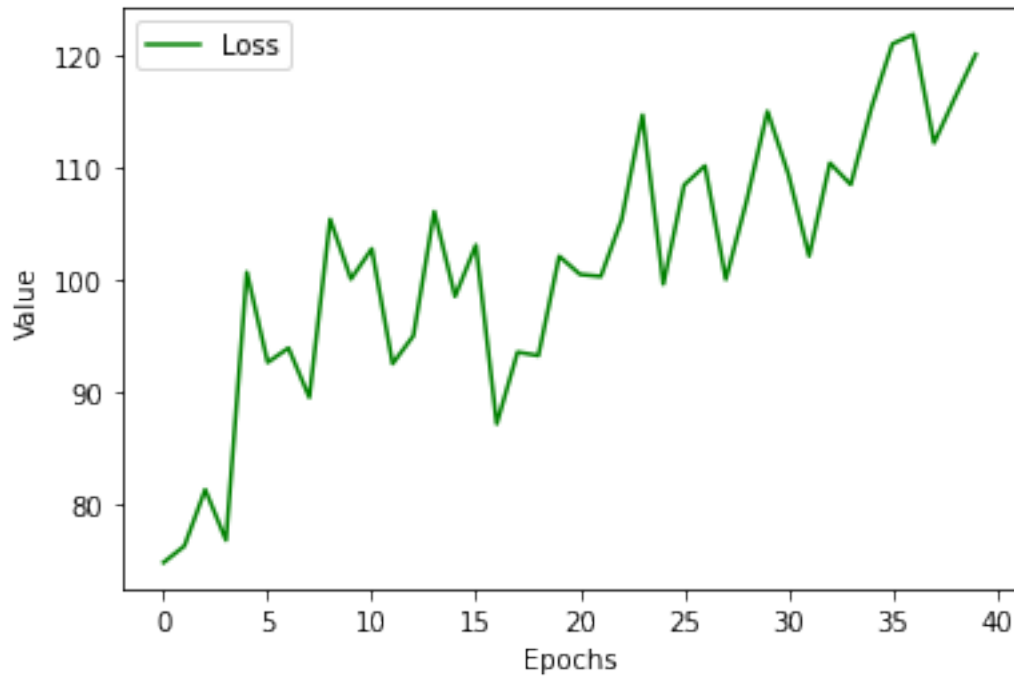


(b) Train Set Loss

Figure 2: Train Set Metrics



(a) Validation Set Accuracy



(b) Validation Set Loss

Figure 3: Validation Set Metrics

The results demonstrate that as the number of epochs increase, there is a reduction in overall loss and an increase in accuracy. This indicates that the model has been effectively trained.