



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 7, 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. Analysis of the dataset, carefully look at what all data is given and how it is formatted.
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 also used Dropout, as our corpus is very small and it is possible for model to *overfit*, thus dropping some of the weights and re-learn them in further iterations.

```
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,
12            word_embedding_dim)
13            self.lstm = nn.LSTM(word_embedding_dim, word_hidden_dim,
14            num_layers = 1, bidirectional = True)
15
16            self.hidden2tag = nn.Linear(word_hidden_dim*2, tagset_size)
17
18            self.dropout = nn.Dropout(0.1)
```

```

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 calculate loss and to calculate the *Gradient Descent* in each iteration, I have used Negative Loss Likelihood function and to optimize the learning rate I have used Adam as optimizer.

```

1     # Define the loss function as the Negative Log Likelihood loss (NLLLoss
2     )
3     loss_function = nn.NLLLoss()
4
5     # We will be using a simple SGD optimizer
6     optimizer = optim.Adam(model.parameters(), lr=0.01)

```

3 Results



