

## Assignment 2

### Part 4

Course: Deep learning for texts and sequences.

Students: Lev Levin, id: 342480456

Lior Shimon, id: 341348498

#### Parameters for ner model:

##### From Pre-train Embedding layer:

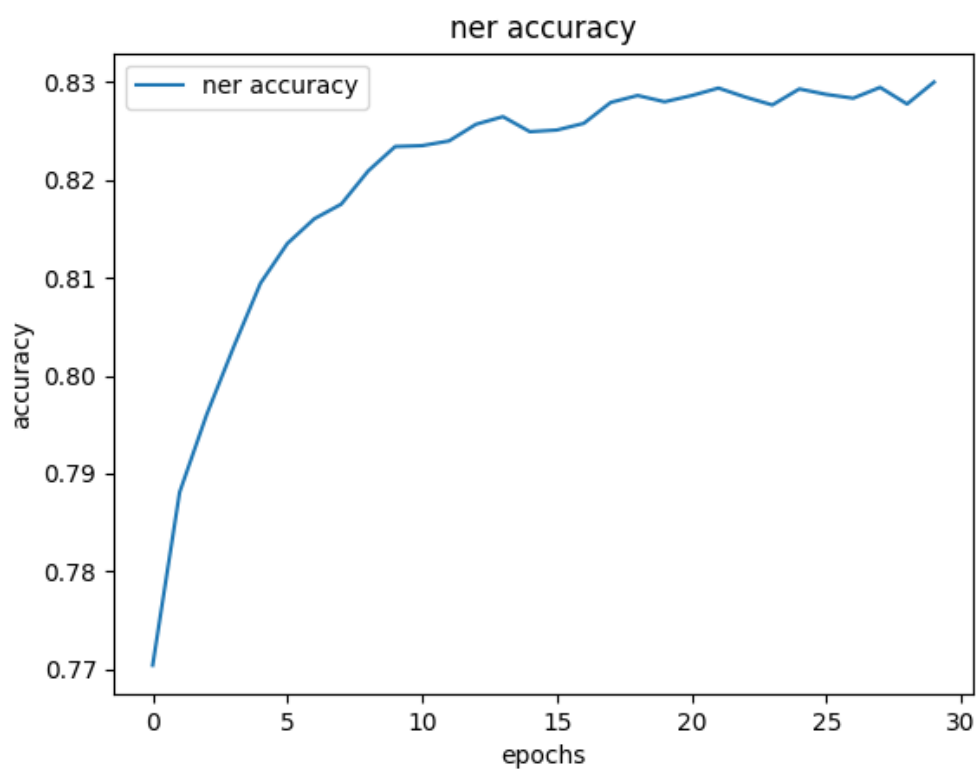
```
epochs = 30
hidden_size = 8
l_r = 0.01
torch.manual_seed(1)
batch_size = 1000
L2 regularization with coefficient: 1e-3
Optimizer: Adam
Random initialization of the Embedding layer: Uniform[-1:1]
```

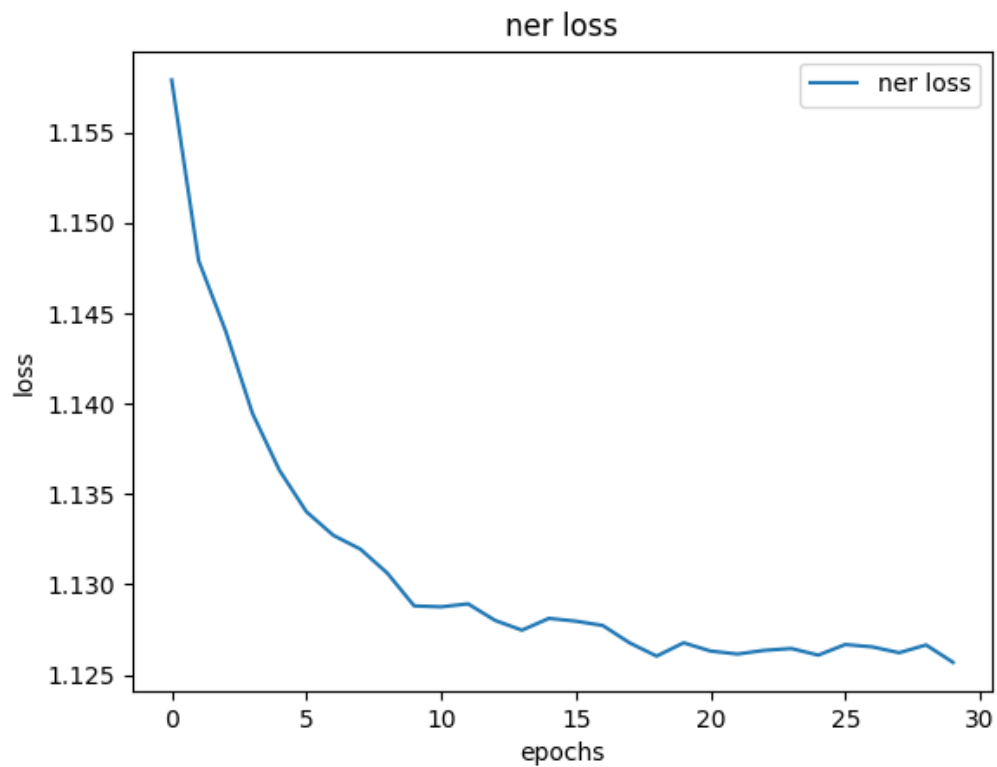
#### With those parameters we reached:

Epoch [30/30], Train\_Loss: 1.0551, Train\_Accuracy: 98.03%

Epoch [30/30] ,Dev\_Loss:1.1257, Dev\_Accuracy:83.00%

#### Graphs:





From randomly initialized Embedding layer:

```
epochs = 25
hidden_size = 6
l_r = 0.01
torch.manual_seed(1)
batch_size = 1000
L2 regularization with coefficient: 1e-3
```

Optimizer: Adam

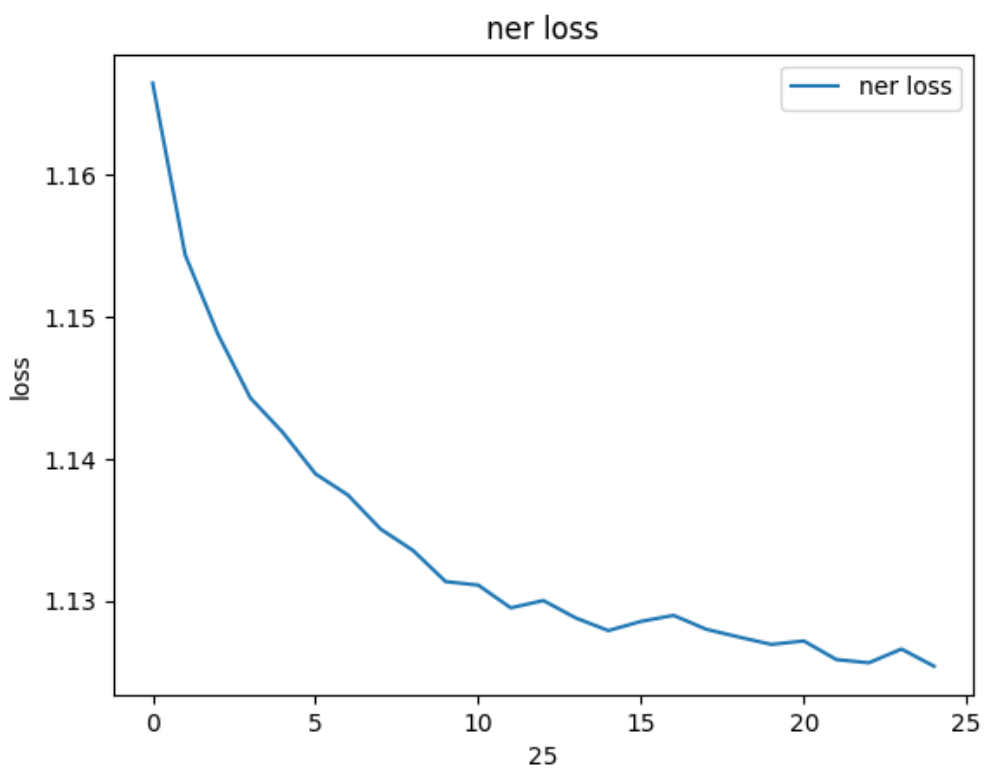
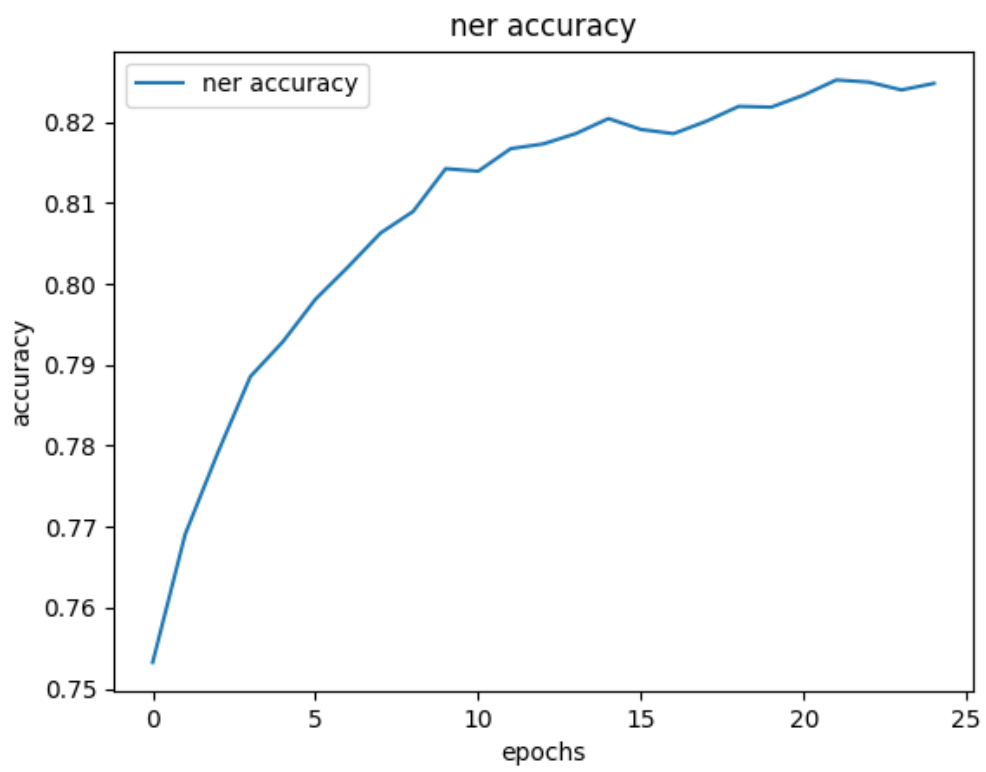
Random initialization of the Embedding layer: Uniform Distribution on [-1:1]

With those parameters we reached:

Epoch [25/25], Train\_Loss: 1.0562, Train\_Accuracy: 97.75%

Epoch [25/25] ,Dev\_Loss:1.1254, Dev\_Accuracy:82.48%

Graphs:



Parameters for pos model:

From Pre-train Embbeding layer:

Number of epochs : 24

Size of the hidden layer : 80

Learning rate : 0.01

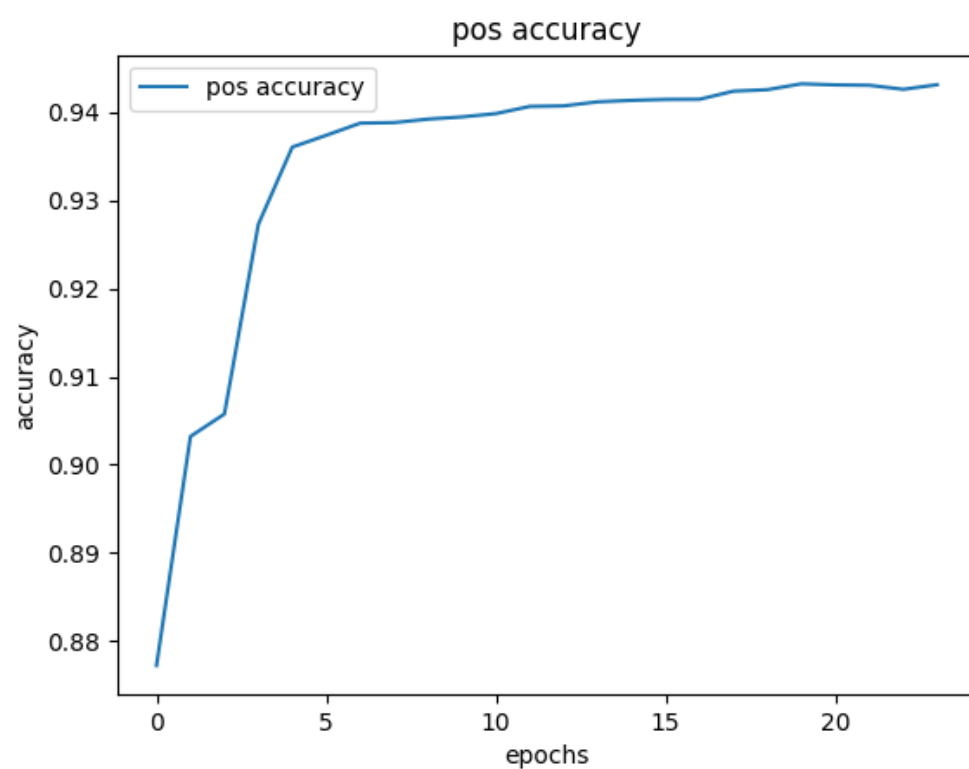
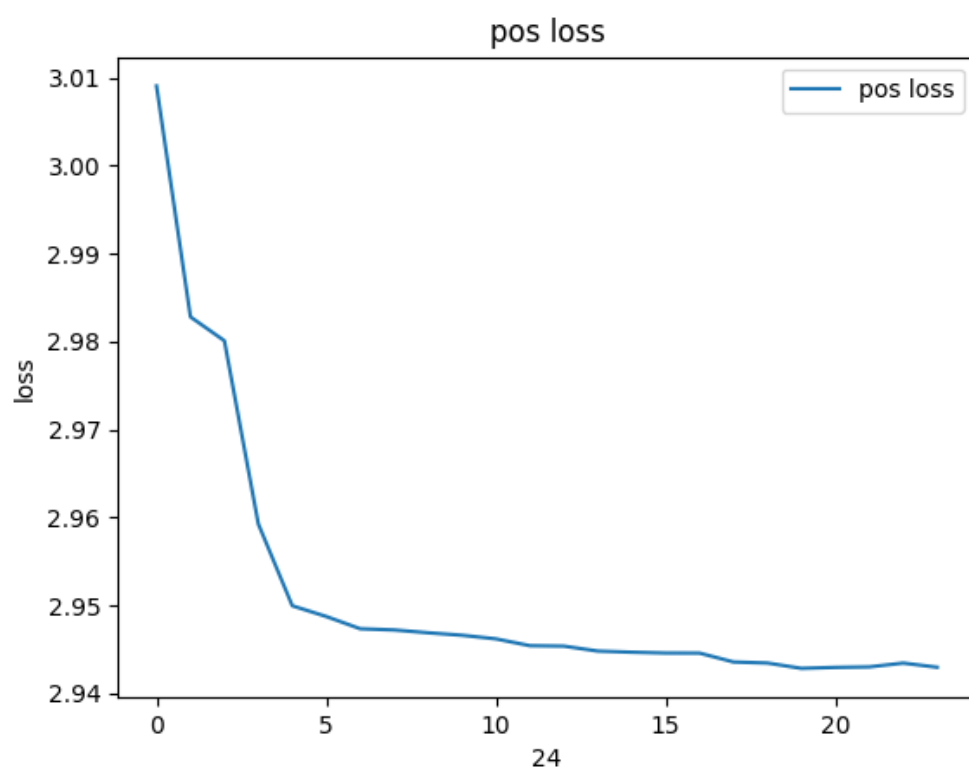
Batch size: 1000

Optimizer: Adam

**With those parameters we reached:**

Epoch [24/24], Train\_Loss: 2.9150, Train\_Accuracy: 97.11%

Epoch [24/24] ,Dev\_Loss:2.9430, Dev\_Accuracy:94.31%



From randomly initialized Embedding layer:

Number of epochs : 15

Size of the hidden layer : 60

Learning rate : 0.001

Batch size: 1000

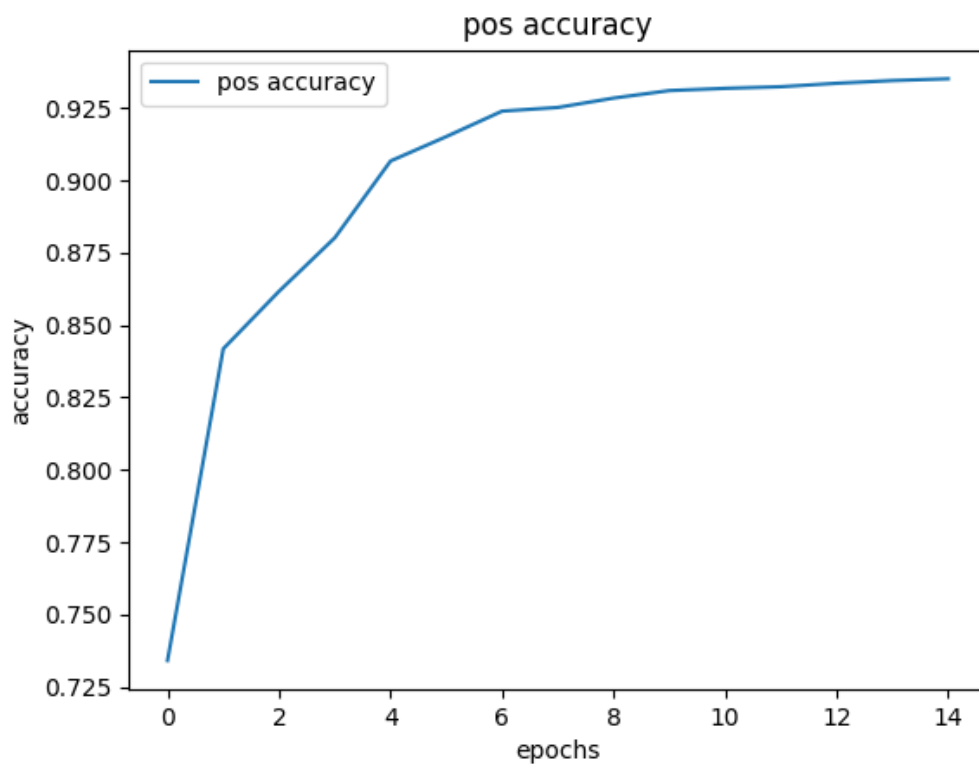
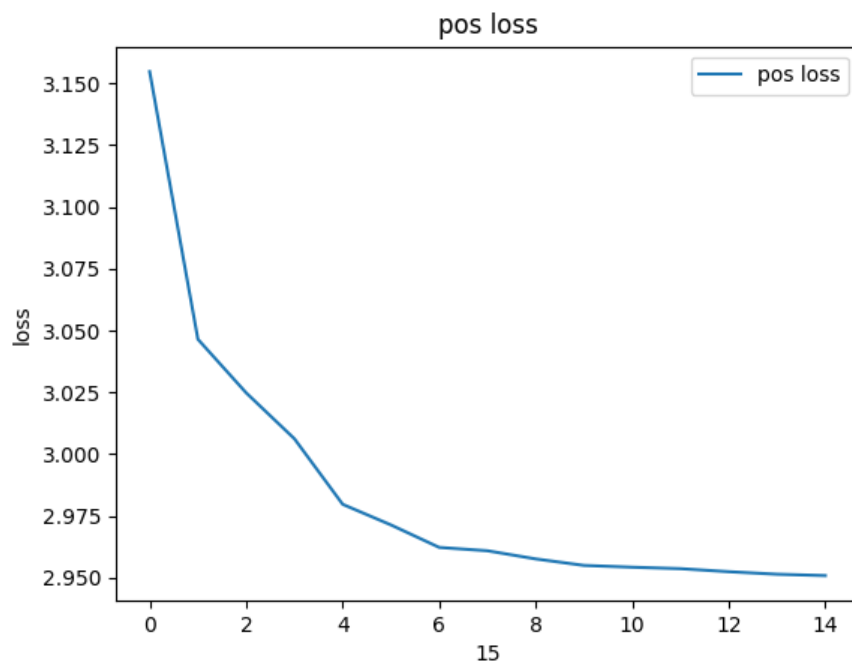
Optimizer: Adam

Random initialization of the Embedding layer: Uniform Distribution on [-1:1]

**With those parameters we reached:**

Epoch [15/15], Train\_Loss: 2.9218, Train\_Accuracy: 96.41%

Epoch [15/15] ,Dev\_Loss:2.9509, Dev\_Accuracy:93.51%



Our thought and our choice for the sub-words units representation:

First off all, we notice that prefixes and suffixes of all the window surrounding the word can also give information to the network about this word we want to predict.

(for example in POS tags, in the sentence "I am doing sport" after the suffix -ing- the next word is probably a Name NN).

That's the reason why we decided to choose to represent all the prefixes and the suffixes of the words in the Embedding matrix of 50 length vectors.

Thus, for each window of 5 words, we concatenated the vector representation of each words together, and concatenated the vector representation of the prefixes of each word and concatenated the vector representation of the suffixes of each words as well, then we summed those 3 freshly created vector, and give them as input to our network.

Note:

About the vector representation of the sub-words unit, we thought about using only the prefix and suffix of the word we want to tag.

The first option was to represent this prefix and suffix as 50 vector length representation as well, and concatenate them with 0s so the tensor size will all be of length 250. But in this case, summing a "regular tensor" with 2 tensors of the form 00..0xxxx00...0 was not natural and would loss a lot of information we want to give to the net.

Another option a 250-length vector each, but this would create 2 different sizes of embedding matrix. So we did not realized this representation neither.

For words that are less than 3, we added special suffixes <suf> and prefixes <pref> that indicate that there are no prefix/suffix in size of 3 and we allocate random vectors for them. It helps to handle cases with words such "a," , "on", which still contain some useful information.

### **Analysis about the results :**

There were slight increase in accuracy that resulted from adding sub units (in ner tagging by 1.5% average improvement) . We have concluded from it, that sub units(prefixes and suffixes) give useful information for named entity recognition tag and pos tagging.. Also, we have noticed that the accuracy increase significantly quicker (see the graphs) . According to our opinion: there are two reasons for faster convergence:

First one is that whenever we used pre trained vectors, the model converged quicker. It happens because model start learn from some good starting point. If we use randomly initialized vectors, it have to relearn the same things, therefore it takes more time to reach the same accuracy.



The second reason is a using of suffixes and prefixes. Model gets more useful information about a word per one training example. That's helps to it to correct weights more correctly per one example.

Between the pre-Embedding layer and the sub-word units, we didn't find a complementarity between them since there are no better result.