

NLP Assignment 3

Analysis (Experiments)

Before testing the experiments as mentioned in the assignment paper, I ran my model on the default configurations given and got the below scores:

| Experiments | UAS | UASno Punc | LAS | LASno Punc | UEM | UEMno Punc | Root | Loss |
|-----------------------|-----------------------|-----------------------|-----------------------|-----------------------|-----------------------|-----------------------|-----------------------|------------------------|
| Default Configuration | 70.536 181668 6 | 73.308 653139 7 | 66.370 865219 2 | 68.716 441530 5 | 10.823 529411 8 | 11.352 941176 5 | 64.470 588235 3 | 0.3614805 272221565 |

Then I tried increasing the learning rate and number of iterations to see how the model performs. I saw a significant increase in the performance of the model at iterations 2001 and learning rate 0.3, so I fixed these two hyperparameters and performed the below testing.

1. **Number of hidden layers** – As per the paper, only one hidden layer was implemented. Here we had to implement 2 more hidden layers.

Configuration: Default with learning rate 0.3, hidden size=150 and number of iterations 2001

| Experiments | UAS | UASno Punc | LAS | LASno Punc | UEM | UEMno Punc | Root | Loss |
|----------------|-----------------------|-----------------------|-----------------------|-----------------------|-----------------------|-----------------------|-----------------------|-------------------------|
| Hidden Layer | 80.088 241892 5 | 82.165 263098 4 | 76.927 487100 2 | 78.593 228960 6 | 20.823 529411 8 | 22.235 294117 6 | 82.058 823529 4 | 0.20614504 769444467 |
| Hidden Layer 1 | 81.267 293167 5 | 83.196 744475 2 | 78.353 316549 1 | 79.912 959927 7 | 22.588 235294 1 | 24.176 470588 2 | 82.470 588235 3 | 0.20979091 927409171 |

| | | | | | | | | |
|----------------|-------|-------|-------|------|------|-------|--------|-------|
| Hidden Layer 2 | 18.45 | 19.50 | 14.67 | 1.56 | 0.64 | 0.645 | 14.877 | 1.133 |
|----------------|-------|-------|-------|------|------|-------|--------|-------|

Observation: As we can see from the result, increase in hidden layer doesn't necessarily increase the accuracy a lot. By adding one hidden layer of 150 nodes, I saw an increase in the scores, but when I added a third layer of 100 nodes I saw a considerable decrease in the performance. The loss computed was very high as evident from the scores above. Also, in the paper it was mentioned that adding more and more hidden layers do not necessarily have an impact in the performance and that is quite evident from the results.

2. **Capturing interactions** – The paper used a cube non-linearity to account for interactions between different combinations of token features. We had to implement sigmoid, relu and tanh to see the performance.

Configuration: Default with learning rate 0.3 and number of iterations 2001

| Experiments | UAS | UASno Punc | LAS | LASno Punc | UEM | UEMno Punc | Root | Loss |
|-----------------|-----------------------|-----------------------|-----------------------|-----------------------|-----------------------|-----------------------|-----------------------|-------------------------|
| Cube Activation | 80.088 241892 5 | 82.165 263098 4 | 76.927 487100 2 | 78.593 228960 6 | 20.823 529411 8 | 22.235 294117 6 | 82.058 823529 4 | 0.20614504 769444467 |
| Relu | 78.904 205199 8 | 81.077 262194 1 | 75.656 205598 6 | 77.431 752670 5 | 18.764 705882 4 | 19.823 529411 8 | 75.176 470588 2 | 0.25332181 15568161 |
| Sigmoid | 68.561 956277 9 | 71.497 202283 4 | 62.818 755141 2 | 65.280 054258 7 | 9.0 | 9.5294 117647 1 | 60.647 058823 5 | 0.47858031 60071373 |
| Tanh | 78.208 739437 1 | 80.444 243486 1 | 74.831 118977 | 76.643 305262 | 17.235 294117 6 | 18.235 294117 6 | 76.058 823529 4 | 0.26513482 15341568 |

Observation: Clearly as the paper suggest the cube activation function outperforms the other functions.

3. **Parallel Layers each for word, Pos tag and deps–**

Configuration: Default with learning rate 0.3 and number of iterations 2001

| Experiments | UAS | UASnoPunc | LAS | LASnoPunc | UEM | UEMnoPunc | Root | Loss |
|------------------------|---------------|---------------|---------------|---------------|-----|---------------|---------------|-------------------|
| Parallel hidden layers | 80.6790138844 | 82.7982818064 | 77.6877632924 | 79.4240660148 | 21 | 22.2352941176 | 79.7647058824 | 0.238249968290329 |

Observation: By implementing separate hidden layers, the performance of the model did not improve. This might be because there is no interaction between the words, pos tags and labels and hence the decrease in scores.

4. Effect of fixing Word, POS and Dep Embedding's – Fix word embeddings and using the pre-trained model.

Configuration: Default with learning rate 0.3, hidden size=150 and number of iterations 2000

| Experiments | UAS | UASnoPunc | LAS | LASnoPunc | UEM | UEMnoPunc | Root | Loss |
|------------------|------|--------------|------|-----------|------|-----------|------|------|
| Fixed Embeddings | 40.8 | 48.731701803 | 34.5 | 37.8 | 2.87 | 2.90 | 33.5 | 0.65 |

Observation: I noticed a remarkable decrease in the in the scores by setting trainable as False. This is because the model is not learning new interactions and therefore the performance is not good.

5. Best Configuration –

I tried to increase the interactions to 4000 with learning rate as 0.3 to see how well the model works as with lesser iterations the loss value was not converging properly. The model accuracy increased significantly by doing so and the loss value also converged.

| Parameter | Value |
|--------------------|-------------|
| Maximum Iterations | 4001 |
| Batch Size | 10000 |
| Hidden Size | 200 |
| Embedding Size | 50 |
| Learning Rate | 0.3 |
| Display Step | 100 |

| | |
|---------------------|-----------------|
| Lambda Value | 10^{-8} |
| Validation Step | 200 |
| Activation Function | Cube Activation |
| Number of tokens | 48 |

| Experiments | UAS | UASnoPunc | LAS | LASnoPunc | UEM | UEMnoPunc | Root | Loss |
|-------------|--------------|---------------|---------------|---------------|---------------|---------------|---------------|---------------------|
| Best | 84.624971957 | 86.4663991409 | 82.0724381185 | 83.5867292149 | 27.4117647059 | 29.1176470588 | 84.4117647059 | 0.16356452628970147 |

- 6. Gradient Clipping** - Gradient Clipping: It is a method which clips the gradients or caps them to a threshold value to prevent it from getting too large or too small.

Usefulness: Gradient Clipping is common in recurrent neural networks where gradients are being propagated both in forward and backward mode.

There are two problems to look into:

1. **Vanishing gradient Issue:** It is seen that when the gradient of the loss function is multiplied with numbers less than 1 then the gradients tend to become vanishingly small and hence the name vanishing gradient problem.
2. **Exploding gradient Issue:** Sometimes the other way happens. The gradient of the loss function when multiplied with numbers greater than 1 then there is a possibility that the gradient explodes and hence the name exploding gradient problem.

In both the cases, we need to set a threshold so that the gradients do not become too large or too small. Hence, we use gradient clipping to clip the numbers to prevent them from the above two problems.

By removing gradient clipping:

Configuration: Default with learning rate 0.1 and number of iterations 1001

| Experiments | UAS | UASnoPunc | LAS | LASnoPunc | UEM | UEMnoPunc | Root | Loss |
|------------------------|-----|-----------|-----|-----------|-----|-----------|------|------|
| Parallel hidden layers | nan | nan | nan | nan | nan | nan | nan | nan |

Observation: The loss values after the first one (started with 4.51 at step 0) started to come as Nan after removing the gradient clipping.