

3 Task 1

First, preprocessing of the tree bank is done to replace single occurrence terminals with "unk" symbol. Then the new tree bank is binarized and then PCFG is created. The result on the test data are following. Instructions are written in readme file at github.

1. The precision score got is 0.846
2. Recall is 0.816
3. F-1 score is 0.831

4 Task 2

Two short sentences: Sentence 1: "This is a car."

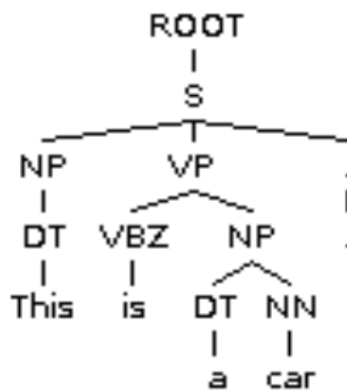


Figure 2: Short Sentence 1

Parsing by the algorithm:

[DT 'This'] [[VB 'is'] [VB 'a']] [[N 'car'] [.]]

Sentence 2: "Hello, there!"

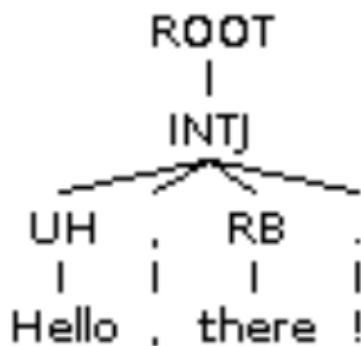


Figure 3: Short Sentence 2

Parsing by the algorithm:

[RB 'Hello'] [[[, ' , '] [RB 'there']] [. ' ! ']]

Two long sentences: Sentence 1: "I shot an elephant in my pajamas"

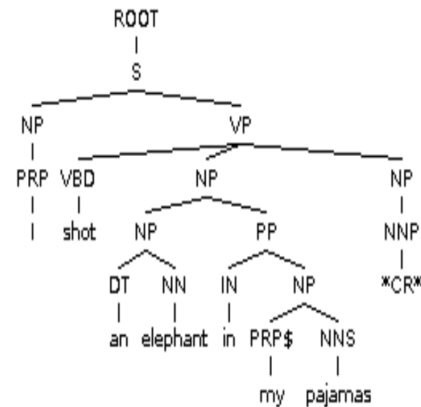


Figure 4: Long Sentence 1

Parsing by the algorithm:

[S [NP 'I'] [VP [V 'shot'] [NP [NP0 [Det 'an'] [N 'elephant']]] [PP [P 'in'] [NP [Det 'my'] [N 'pajamas']]]]]]

Sentence 2: "I am really upset that you always fight over small issues."

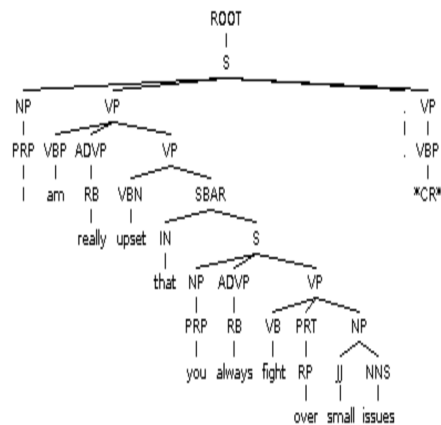


Figure 5: Long Sentence 2

Parsing by the algorithm:

[S ['I' PRP] [VP [V 'am'] [NP [NP0 'really'] [N 'upset']]] [[IN 'that'] [NP 'you'] [NP 'always'] [N 'fight'] [RP 'over'] [NP 'small'] [NNS 'issues'] [. ' ']]]

By changing the probabilities different tags can be obtained, but still one or two tags were wrong in many sentences.

References

Beatrice Santorini Ann Taylor, Mitchell Marcus. 1992. [The penn treebank project](#).

Roger Levy. 2015. [Treebanks and pcfgs](#).