# Constituency Parsing with Probabilistic Context-Free Grammars

**Raghu Chittersu**
Department of Computation and Data Sciences
Indian Institute of Sciences, Bengaluru
chittersuv@iisc.ac.in

## 1 Introduction

Constituency parsing extracts constituency-based parse tree from a given sentence. It represents syntactic structure of the sentence.

Constituency parsing has many applications like Parts of speech tagging and other text mining tasks. Context Free grammar is essential for generation of parse trees. Since there are ambiguities one uses the tree with highest generation probability. CYK algorithm is a Dynamic Programming type algorithm which can be used to generate parse tree given Probabilistic Context Free Grammar

In this assignment constituency parsing is implemented using PCFG and CYKParser. Smoothing and prior probabilities are used to handle the sparsity of training data. Then the model's performance is various methods along with comparing with some of the best parsers.

The code is available at https://github.com/rv-chittersu/CYK-Parser.

## 2 Data

In this assignment, nltk's treebank corpus is used as data-set. It contains 3914 parsed and raw sentences in 200 files.

The production rules from parsed sentences of train set are used to build the PCFG. The raw sentences from test set are used to predict the parse tree and evaluation.

## 3 Methodology

The complete implementation of Constituency Parsing involves following components
(1) Probabilistic Context Free Grammar
(2) CYK Parser

### 3.1 Probabilistic Context Free Grammar

The PCFG is made of
(1) set of non-terminal symbols $M$,
(2) set of terminal symbols $T$,
(3) set of production rules $R$,
(4) start symbol $S$ and
(5) probabilities for with production rules $P$

The CYK Parser requires the PCFG to be in Chomsky Normal Form i.e each production should satisfy following rules

$$X -> Y\ Z$$

$$X -> t$$

where $X, Y, Z$ are Non Terminals and t is Terminal.

### 3.2 CYK Parser

The CYK Algorithm uses PCFG which is in Chomsky Normal Form to build a parse tree from given sentence The CYK Algorithm follows Dynamic Program Paradigm

The recurrence relation for probability of assigning a segment from $i$ to $j$ a non terminal $X$ is given by

$$P_X(i,j) = \max_{i<k<j}(P_Y(i,k)*P_Z(k,j)*P(X->YZ))$$

In bottom up approach or memoization method of dynamic program we start from sub-string of length 1 and keep adding rules satisfying above recurrence relation to the parse table till we reach string length. One should also keep track of the rules that generated maximum probability at each position through back pointers.

Once the parsing is done. The tree is generated recursively from Start Symbol $S$ in final cell in the table through back pointers.

## 4 Handling Sparsity

Since the number words in vocabulary are limited. If new token show up while parsing a sentence the above algorithm will fail to parse as there are no production rule that lead to the token.

**Method1** If token $t$ is OOV

$$P(X->t) = \frac{\#(X->Terminal)}{\#(Terminals)}$$

The another way of looking at above probability is probability of X generating being the lhs given the rhs of the production a non terminal.

We augment PCFG with above rule if OOV token $t$ appears in the sentence.

**Method2**

$$P(X->t) = \frac{1}{\#X + 1}$$

Along with above method add one smoothing is used to get the probability for unknown at each Terminal generating non terminal.

## 5 Experiment Setup

### 5.1 Pre Processing

The important preprocessing step is to convert the parse tree from treebank dataset to Chomsky Normal Form. It is handles using nltk module.

The other alterations done are converting strings to lower case and numbers to special NUM_TOKEN

### 5.2 Training

During training phase the parse tress from treebank training data set are processed to produce occurances of Production Rules. The produces output is saved to checkpoint file and will be used in further steps.

### 5.3 Testing

The testing is done by parsing the sentences in test set dervied from treebank datatset. The task is accomplished by multiple processes where each of them writes gold and generated parse tree to individual files. Then all the files by multiple processes are stitched to generate final gold and generated results

### 5.4 Evaluation

The performance of the model is computed based on comparing predicted parses and gold parses from nltk dataset. Metrics such as Bracketing Recall, Bracketing Precision and Bracketing FMeasure, Tagging accuracy is reported along with other metrics. The results are obtained using NYU's program EVALB (`https://nlp.cs.nyu.edu/evalb/`)

## 6 Results

The results of various metrics for above experiment with two different sparsity methods can be found below.

Table 1: Results on Test Data for method 1

| Metric | Result |
|---|---|
| **Bracketing Recall** | 44.32 |
| **Bracketing Precision** | 42.04 |
| **Bracketing FMeasure** | 43.15 |
| Complete match | 1.91 |
| Average crossing | 11.66 |
| No crossing | 4.31 |
| 2 or less crossing | 11.72 |
| **Tagging accuracy** | 84.30 |

The results are on 418 sentences from text set derived from treebank dataset.

Table 2: Results on Test Data for method 2

| Metric | Result |
|---|---|
| **Bracketing Recall** | 44.69 |
| **Bracketing Precision** | 41.87 |
| **Bracketing FMeasure** | 43.32 |
| Complete match | 1.91 |
| Average crossing | 10.93 |
| No crossing | 3.18 |
| 2 or less crossing | 11.46 |
| **Tagging accuracy** | 82.70 |

The results are on 418 sentences from text set derived from treebank dataset.

It is difficult to compare these two models as scores are very similar. The first model has the better Tagging Accuracy where as second model got better Bracketing FMeasure

## 7 Comparison with Stanford Parser

**CASE-1**

| Sentence | Rose is red . |
|---|---|
| Gold Parse | (ROOT (S (S (VP (VBD rose))) (VP (VBZ is) (ADJP (JJ red))) (. .))) |
| Predicted | (S<VP-.>(VP (VBD Rose) (VP (VBZ is) (ADJP (JJ red)))) (. .)) |

The mismatch is due to the absence of following production rules

$$S + VP - > VBD$$

$$S - > S + VP \ S| < VP - . >$$

adding them to rules with high probability fixed the ambiguity.

**CASE-2**

| Sentence | He knows. |
|---|---|
| Gold Parse | (ROOT (ROOT (S (NP (PRP He)) (VP (VBZ knows)) (. .))) |
| Predicted | (S (NP (PRP He)) (S<VP-.>(VP (VBZ knows)) (. .))) |

The mismatch is due to the relative probabilities of following production rules

$$S - > NP \ VP$$

$$S - > NP \ S| < VP - . >$$

pdating the rule with high probability fixed the ambiguity.

## 8 Conclusion

In this assignment a Probabilistic Context Free Grammar is built from training set. Then CYK algorithm is used along with smoothing to generate Constituency Parse of given sentence. Then the results produces is compared with the Stanford Parser.