**Faculty of Computer Science, Dalhousie University**      *27-Nov-2024*

**CSCI 4152/6509 — Natural Language Processing**

**Lecture 21: CYK Algorithm and PCFGs**

Location: Carleton Tupper Building Theatre C      Instructor: Vlado Keselj
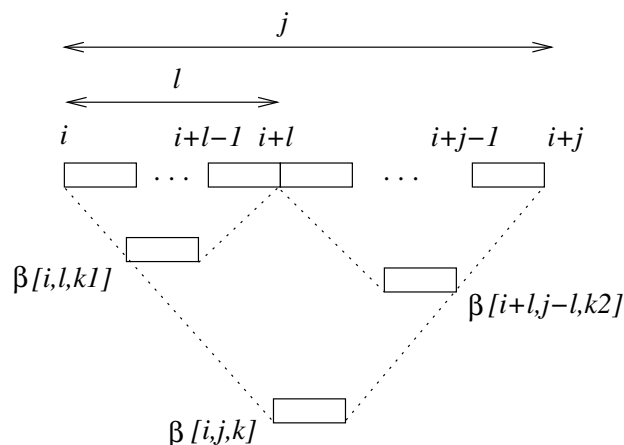Time:      16:05 – 17:25

**Previous Lecture**

- Phrase structure in English (continued):
  - NP, VP, PP, ADJP, ADVP
- Heads and dependency, dependency tree
- **CYK Chart Parsing Algorithm**
- Chomsky Normal Form (CNF)
- CYK algorithm example

**Implementation**

The example implies that we need to use a two-dimensional table to store chart entries. Using a two-dimensional table is a possible solution, but in that case the table entries would be quite complex since each of them needs to store a set of non-terminals. To make the solution simpler, we can use a three-dimensional table, such that the third dimension corresponds to all different non-terminals.

**Explanation of Index Use in CYK**



**CYK Algorithm**

Let all nonterminals be: $N^1, \ldots N^m$.

In the standard CYK algorithm, we have a two dimensional table $\beta$ in which only the entries $\beta_{ij}$, $1 \leq i \leq i+j-1 \leq n$, are used. Each entry $\beta_{ij}$ contains a set of nonterminals that can produce substring $w_i \ldots w_{i+j-1}$ using the grammar rules, i.e., $\beta_{ij} = \{N | N \Rightarrow^* w_i \ldots w_{i+j-1}\}$.

If we enumerate all nonterminals: $N^1$, $N^2$, ..., $N^m$, then each set of nonterminals $\beta_{ij}$ can be represented by extending $\beta$ to be a 3-dimensional table $\beta_{ijk}$, in which $\beta_{ijk} = 1$ means that $N^k$ can produce substring $w_i \ldots w_{i+j-1}$, and $\beta_{ijk} = 0$ that it cannot.

---

**Algorithm 1** CYK Parsing Algorithm

---

**Require:** sentence $= w_1 \ldots w_n$, and a CFG in CNF with nonterminals $N^1 \ldots N^m$,
    $N^1$ is the start symbol
**Ensure:** parsed sentence
 1: allocate matrix $\beta \in \{0, 1\}^{n \times n \times m}$ and initialize all entries to 0
 2: **for** $i \leftarrow 1$ to $n$ **do**
 3:     **for all** rules $N^k \to w_i$ **do**
 4:         $\beta[i, 1, k] \leftarrow 1$
 5: **for** $j \leftarrow 2$ to $n$ **do**
 6:     **for** $i \leftarrow 1$ to $n - j + 1$ **do**
 7:         **for** $l \leftarrow 1$ to $j - 1$ **do**
 8:             **for all** rules $N^k \to N^{k_1} N^{k_2}$ **do**
 9:                 $\beta[i, j, k] \leftarrow \beta[i, j, k]$ OR $(\beta[i, l, k_1]$ AND $\beta[i + l, j - l, k_2])$
10: **return** $\beta[1, n, 1]$

---

The line $\beta[i, j, k] \leftarrow \beta[i, j, k]$ OR $(\beta[i, l, k_1]$ AND $\beta[i + l, j - l, k_2])$ in the algorithm is essentially is a shorthand expression for:

  **if** $\beta[i, l, k_1]$ AND $\beta[i + l, j - l, k_2]$ **then**
    $\beta[i, j, k] \leftarrow 1$

# 26   Efficient Inference in PCFG Model

**Efficient Inference in PCFG Model**

Let us consider the marginalization task:

P(sentence) $=?$

If 'sentence' is the following sequence of words: $w_1 w_2 \ldots w_n$, then P(sentence) is the following conditional probability:

$$\text{P(sentence)} = \text{P}(w_1 w_2 \ldots w_n | S)$$

i.e., it is the probability of generating the sentence given that we start from $S$, i.e, it is $\text{P}(S \Rightarrow^* w_1 \ldots w_n)$.

An obvious way to calculate this marginal probability is to find all parse trees of a sentence and sum their probabilities, i.e:

$$\text{P(sentence)} = \sum_{t \in T} \text{P}(t),$$

where $T$ is the set of all parse trees of the sentence 'sentence'. However, this may be very inefficient. We also need a way to find all parse trees.

As an example illustrating that the above direct approach may lead to an exponential algorithm, consider a CFG with only two rules $S \Rightarrow S\ S$ and $S \Rightarrow a$. The sentences $a^n$ have as many parse trees as there are binary trees with $n$ leaves, which is a well-known Catalan number, $\approx \frac{4^n}{n^{3/2}\sqrt{\pi}}$ as $n \to \infty$.

An algorithm for efficient marginalization can be derived from the CYK algorithm.

**PCFG Marginalization**

The CKY algorithms is adapted to solve the problem of efficient PCFG marginalization, but replacing entries of the table $\beta$ with numbers between 0 and 1. These numbers are called <u>inside probabilities,</u> and they represent the following probabilities:

$$\beta[i, j, k] = \mathrm{P}(w_i \ldots w_{i+j-1}|N^k)$$

So, $\beta[i, j, k]$ is the probability that the string $w_i \ldots w_{i+j-1}$ is generated in a derivation where the starting non-terminal is $N^k$. Algorithm 2 is the probabilistic CYK algorithm for calculating P(sentence).

---

**Algorithm 2** Probabilistic CYK for P(sentence)

---

**Require:** sentence $= w_1 \ldots w_n$, and a PCFG in CNF with nonterminals $N^1 \ldots N^m$, $N^1$ is the start symbol
**Ensure:** P(sentence) is returned
 1: allocate $\beta \in \mathbb{R}^{n \times n \times m}$ and initialize all entries to 0
 2: **for** $i \leftarrow 1$ to $n$ **do**
 3:    **for all** rules $N^k \rightarrow w_i$ **do**
 4:       $\beta[i, 1, k] \leftarrow \mathrm{P}(N^k \rightarrow w_i)$
 5: **for** $j \leftarrow 2$ to $n$ **do**
 6:    **for** $i \leftarrow 1$ to $n - j + 1$ **do**
 7:       **for** $l \leftarrow 1$ to $j - 1$ **do**
 8:          **for all** rules $N^k \rightarrow N^{k_1} N^{k_2}$ **do**
 9:             $\beta[i, j, k] \leftarrow \beta[i, j, k]+ \mathrm{P}(N^k \rightarrow N^{k_1} N^{k_2}) \cdot \beta[i, l, k_1] \cdot \beta[i + l, j - l, k_2]$
10: **return** $\beta[1, n, 1]$

---

**PCFG Marginalization Example (grammar)**

| S | $\rightarrow$ | NP VP | /1 | VP | $\rightarrow$ | V NP | /.5 | N | $\rightarrow$ | time | /.5 |
|---|---|---|---|---|---|---|---|---|---|---|---|
| NP | $\rightarrow$ | time | /.4 | VP | $\rightarrow$ | V PP | /.5 | N | $\rightarrow$ | arrow | /.3 |
| NP | $\rightarrow$ | N N | /.2 | PP | $\rightarrow$ | P NP | /1 | N | $\rightarrow$ | flies | /.2 |
| NP | $\rightarrow$ | D N | /.4 | | | | | D | $\rightarrow$ | an | /1 |
| V | $\rightarrow$ | like | /.3 | | | | | | | | |
| V | $\rightarrow$ | flies | /.7 | | | | | | | | |
| P | $\rightarrow$ | like | /1 | | | | | | | | |

**PCFG Marginalization Example (chart)**



$\beta$ [1,1,.]      $\beta$ [1,2,.]      $\beta$ [1,3,.]      $\beta$ [1,4,.]      $\beta$ [1,5,.]

|   | time | | flies | | like | | an | | arrow | |
|---|---|---|---|---|---|---|---|---|---|---|

NP: 0.4
N: 0.5

V: 0.7
N: 0.2

V: 0.3
P: 1

D: 1

N: 0.3

$\beta$ [2,1,.]  →  NP: 0.02                           NP: 0.12

D  N  P(NP–>D N)
1 x 0.3 x 0.4 = 0.12

N  N    P(NP–> N N)
0.5 x 0.2 x 0.2 = 0.02

PP: 0.12
VP: 0.018

P  NP  P(PP–>P NP)
1 x 0.12 x 1 = 0.12
V  NP  P(VP–>V NP)
0.3 x 0.12 x 0.5 = 0.018

VP: 0.042

V  PP  P(VP–>V PP)
0.7 x 0.12 x 0.5 = 0.042

S: 0.0168
0.01716

NP  VP  P(S–>NP VP)
0.4 x 0.042 x 1 = 0.0168
add
NP  VP P(S–>NP VP)
0.02 x 0.018 x 1 = 0.00036
0.0168+0.00036=0.01716

P(time flies like an arrow) =

= 0.01716

**Conditioning**

The conditioning computational problem in the PCFG model becomes the task of finding the conditional probability P(tree|sentence), for a particular sentence and a particular parse three of the given sentence. Using the definition of the conditional probability, we have:

$$P(\text{tree}|\text{sentence}) = \frac{P(\text{tree}, \text{sentence})}{P(\text{sentence})}$$

and since the sentence is a part of the parse tree, we can further write:

$$P(\text{tree}|\text{sentence}) = \frac{P(\text{tree}, \text{sentence})}{P(\text{sentence})} = \frac{P(\text{tree})}{P(\text{sentence})}$$

*Slide notes:*

> **Conditioning**
> - Conditioning in the PCFG model: P(tree|sentence)
> - Use the formula:
>
> $$P(\text{tree}|\text{sentence}) = \frac{P(\text{tree}, \text{sentence})}{P(\text{sentence})} = \frac{P(\text{tree})}{P(\text{sentence})}$$
>
> - P(tree) — directly evaluated
> - P(sentence) — marginalization

P(tree) is calculated by multiplying probabilities of all rules in the tree, and P(sentence) is calculated by the Algorithm 2 used for marginalization.

## Completion

The <u>completion</u> task becomes the parsing problem; i.e., the problem of finding the most probably parse tree give the sentence, which can be expressed as:

$$\arg\max_{\text{tree}} P(\text{tree}|\text{sentence})$$

*Slide notes:*

---

**Completion**

– Finding the most likely parse tree of a sentence:

$$\arg\max_{\text{tree}} P(\text{tree}|\text{sentence})$$

– Use the CYK algorithm in which line 9 is replaced with:

   9: $\beta[i,j,k] \leftarrow \max(\beta[i,j,k], P(N^k \rightarrow N^{k_1}N^{k_2}) \cdot \beta[i,l,k_1] \cdot \beta[i+l,j-l,k_2])$

– Return the most likely tree

---

The most probable completion is computed in a similar way to the marginalization algorithm (Algorithm 2). The difference is that the line 9 is replaced by the line

   9: $\beta[i,j,k] \leftarrow \max(\beta[i,j,k], P(N^k \rightarrow N^{k_1}N^{k_2}) \cdot \beta[i,l,k_1] \cdot \beta[i+l,j-l,k_2])$

Additionally in step 10, we are not just interested in $\beta[1,n,1]$, which is the probability of the most probable tree, but we also want to obtain the actual tree.

---

**Algorithm 3** CYK-based Completion Algorithm for $\arg\max_t P(t|\text{sentence})$

---

**Require:** sentence $= w_1 \ldots w_n$, and a PCFG in CNF with nonterminals $N^1 \ldots N^m$, $N^1$ is the start symbol
**Ensure:** The most likely parse tree is returned
1: allocate $\beta \in \mathbb{R}^{n \times n \times m}$ and initialize all entries to 0
2: **for** $i \leftarrow 1$ to $n$ **do**
3:    **for all** rules $N^k \rightarrow w_i$ **do**
4:       $\beta[i,1,k] \leftarrow P(N^k \rightarrow w_i)$
5: **for** $j \leftarrow 2$ to $n$ **do**
6:    **for** $i \leftarrow 1$ to $n - j + 1$ **do**
7:       **for** $l \leftarrow 1$ to $j - 1$ **do**
8:          **for all** rules $N^k \rightarrow N^{k_1}N^{k_2}$ **do**
9:             $\beta[i,j,k] \leftarrow \max(\beta[i,j,k], P(N^k \rightarrow N^{k_1}N^{k_2}) \cdot \beta[i,l,k_1] \cdot \beta[i+l,j-l,k_2])$
10: **return** Reconstruct$(1,n,1,\beta)$

---

The tree can be reconstructed from the table using algorithm 4:

10: **return** Reconstruct$(1,n,1,\beta)$

## PCFG Completion Example (grammar)

| S | $\rightarrow$ | NP VP | /1 | VP | $\rightarrow$ | V NP | /.5 | N | $\rightarrow$ | time | /.5 |
| NP | $\rightarrow$ | time | /.4 | VP | $\rightarrow$ | V PP | /.5 | N | $\rightarrow$ | arrow | /.3 |
| NP | $\rightarrow$ | N N | /.2 | PP | $\rightarrow$ | P NP | /1 | N | $\rightarrow$ | flies | /.2 |
| NP | $\rightarrow$ | D N | /.4 | | | | | D | $\rightarrow$ | an | /1 |
| V | $\rightarrow$ | like | /.3 | | | | | | | | |
| V | $\rightarrow$ | flies | /.7 | | | | | | | | |
| P | $\rightarrow$ | like | /1 | | | | | | | | |

---

**Algorithm 4** Reconstruct$(i, j, k, \beta)$

---

**Require:** $\beta$ — table from CYK, $i$ — index of the first word, $j$ — length of sub-string sentence, $k$ — index of non-terminal
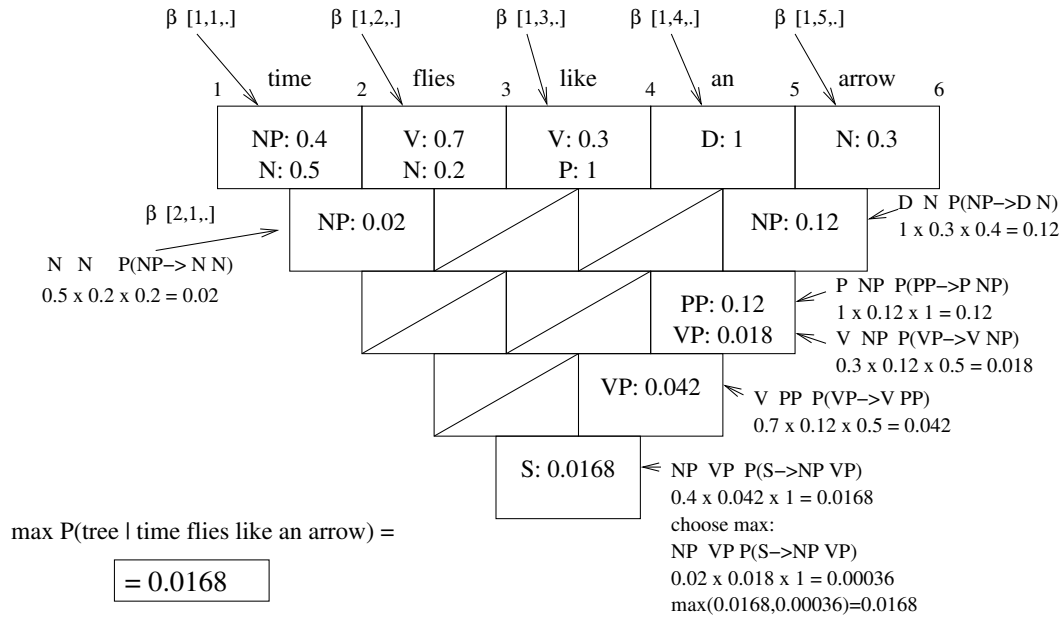
**Ensure:** a most probable tree with root $N^k$ and leaves $w_i \ldots w_{i+j-1}$ is returned

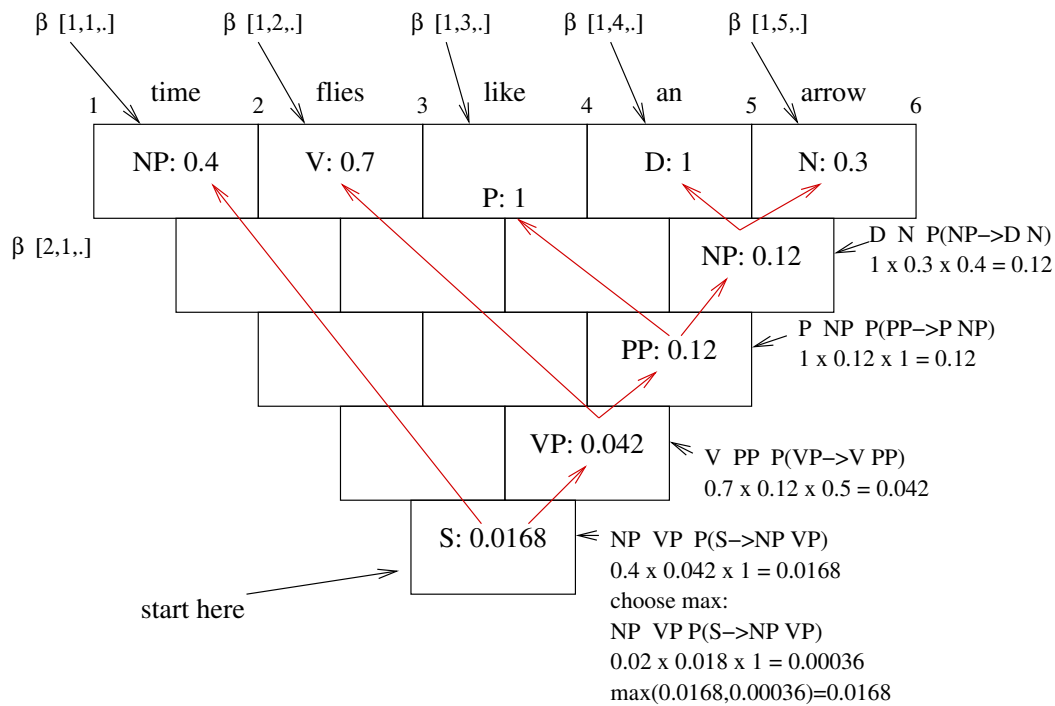1: **if** $j = 1$ **then**
2:     **return** tree with root $N^k$ and child $w_i$
3: **for** $l \leftarrow 1$ to $j - 1$ **do**
4:     **for all** rules $N^k \rightarrow N^{k_1} N^{k_2}$ **do**
5:         **if** $\beta[i, j, k] = \mathrm{P}(N^k \rightarrow N^{k_1} N^{k_2}) \cdot \beta[i, l, k_1] \cdot \beta[i + l, j - l, k_2]$ **then**
6:             create a tree $t$ with root $N^k$
7:             $t.left\_child \leftarrow$ Reconstruct$(i, l, k_1, \beta)$
8:             $t.right\_child \leftarrow$ Reconstruct$(i + l, j - l, k_2, \beta)$
9:             **return** $t$

---

## PCFG Completion Example (chart)
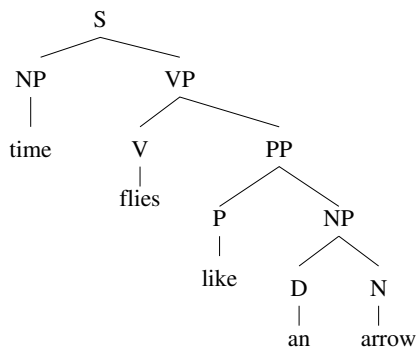


max P(tree | time flies like an arrow) =

= 0.0168

**PCFG Completion Example (tree reconstruction)**

β [1,1,.]          β [1,2,.]          β [1,3,.]          β [1,4,.]          β [1,5,.]

1      time      2      flies      3      like      4      an      5      arrow      6

| NP: 0.4 | V: 0.7 | | D: 1 | N: 0.3 |

P: 1

β [2,1,.]

NP: 0.12          D N P(NP–>D N)
                  1 x 0.3 x 0.4 = 0.12

PP: 0.12          P NP P(PP–>P NP)
                  1 x 0.12 x 1 = 0.12

VP: 0.042         V PP P(VP–>V PP)
                  0.7 x 0.12 x 0.5 = 0.042

S: 0.0168         NP VP P(S–>NP VP)
                  0.4 x 0.042 x 1 = 0.0168
                  choose max:
start here        NP VP P(S–>NP VP)
                  0.02 x 0.018 x 1 = 0.00036
                  max(0.0168,0.00036)=0.0168

**PCFG Completion Example (final tree)**

The most probable three:

```
                    S
         _____|_____
        NP                      VP
        |              _____|_____
       time           V                   PP
        |             |            _____|_____
      time          flies         P              NP
                                  |          _____|_____
                                 like        D          N
                                  |          |          |
                                 like        an       arrow
```

**Topics related to PCFGs**

– An interesting open problem is whether the inference in PCFGs can be reduced to a message-passing-style algorithm as used in Bayesian Networks?

**Issues with PCFGs**

The Probabilistic Context-Free Grammars were shown to perform quite well in parsing English, but usually with some additional mechanisms to address certain issues. Two most prominent issues in using PCFGs to parse nature languages are the inability of PCFGs to capture structural and lexical dependencies.

**Structural dependencies**    are rule dependencies on the position in a parse tree. For example, pronouns occur more frequently as subjects than objects in sentences, so the rule choice between NP → PRP and NP → DT NN should depend on the position of a noun phrase in a tree. Generally, NL parse trees are usually deeper at their right side than the left side, and this propertly is typically not modeled well with PCFGs.

**Lexical dependencies**    are rule dependencies on the words that are eventually derived from those rules, particularly phrase head words. As an example, the PP-attachment problem is resolved based on the rule probabilities of the rules applied higher in the parse tree, such as NP → NP PP and VP → VB NP PP, while they truly frequently depend on the verb being used and other word, particularly head words.

*Slide notes:*

---
**PP-Attachment Example**

- – Consider sentences:
    - – "Workers dumped sacks into a bin." and
    - – "Workers dumped sacks of fish."
- – and rules:
    - – NP → NP PP
    - – VP → VBD NP
    - – VP → VBD NP PP
---

As an example, let us consider simple sentences:

- – "Workers dumped sacks into a bin." (from [JM]), and
- – "Workers dumped sacks of fish."

At some level of parsing, we can see both of these sentences as:

- – NP VBD NP PP

and now the question is whether the "NP PP" should be combined to make an NP, or the sequence "VBD NP PP" should be combined to make a VP. In a PCFG this will depend only on the probability of the rules: NP → NP PP, VP → VBD NP, and VP → VBD NP PP. However, we can see that the probailities should actually depend on affinity of the verb 'dump' and preposition 'into' on one side, and the noun 'sacks' and the preposition 'of' on other side.

**A Solution: Probabilistic Lexicalized CFGs**

- – use heads of phrases
- – expanded set of rules, e.g.:

$$VP(dumped) \rightarrow VBD(dumped) \ NP(sacks) \ PP(into)$$

- – large number of new rules
- – sparse data problem
- – solution: new independence assumptions
- – proposed solutions by Charniak, Collins, etc. around 1999