# Parsing (part 3)

Chapter 13 J&M'09

## Parsing

**CKY** (or CYK)
Named after John Cocke, Daniel Younger and Tadao Kasami.

- Passive chart parser
- Very efficient (runs in polynomial time; $n^3 \times |G|$)
- Requires grammar transformation.

**CKY requires grammars in Chomsky Normal Form (CNF)**

In CNF, all rules must be of one of the following forms:

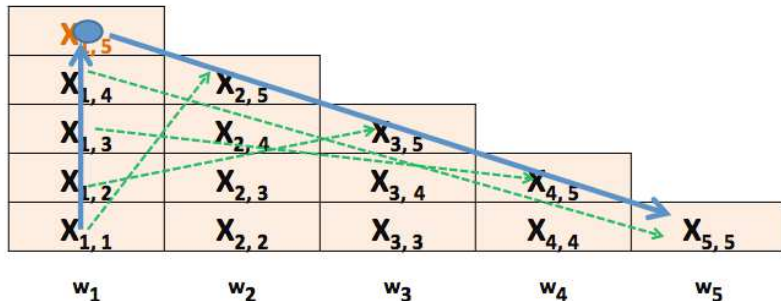$X \rightarrow Y Z$

$K \rightarrow w$ (where '$w$' is a word token)

This **binarization** step is crucial for efficient parsing.

# *Parsing*

**Conversion to CNF**

- Copy all conforming rules to new grammar

- Eliminate unit productions:
  Rules like **NP → PN** and **PN → Tom**
  become a single rule **NP → Tom**

- Convert branching rules:
  Rules like **VP → DTV NP NP** become

  **VP → DVP NP**
  **DVP → DTV NP**

# *Parsing*

For input $w_1 \ldots w_t$, build parse triangle:



Each row corresponds to a string of ascending length.
Each cell corresponds to all of the possible categories for the corresponding span. For example, in

(1) Doves dove

Each cell of row 1 would be {N,V}.

# Parsing

Input: *Tom saw a friend from Australia*

| Span length | | | | | | |
|---|---|---|---|---|---|---|
| 6 | S | | | | | |
| 5 | – | VP | | | | |
| 4 | S | – | NP | | | |
| 3 | – | VP | – | N | | |
| 2 | – | – | NP | – | PP | |
| 1 | NP | TV, N | DT | N,TV | P | NP |
| | Tom | saw | a | friend | from | Australia |
| | 0 | 1 | 2 | 3 | 4 | 5 | 6 |

$S \rightarrow NP\ VP$    $VP \rightarrow TV\ NP$    $VP \rightarrow VP\ PP$
$NP \rightarrow DT\ N$    $N \rightarrow N\ PP$    $PP \rightarrow P\ NP$
$NP \rightarrow Tom$    $TV \rightarrow saw$    $DT \rightarrow a$
$NP \rightarrow Australia$    $N \rightarrow saw$    $TV \rightarrow friend$
$N \rightarrow friend$    $P \rightarrow from$

# Parsing

Sometimes, the chart is shown as a parse triangle:

# Parsing

Sometimes, the chart is rotated:



$$w_1 \qquad w_2 \qquad w_2 \qquad w_4$$

## Parsing

**Exercise:**

Show that 'b a a b a' is parsed by CKY and the following grammar.

S → A B
S → B C
A → B A
A → a
B → C C
B → b
C → A B
C → a

# *Statistical Parsing*

**Probabilistic Context-free Grammar** (PCFG)

Rules are augmented with a probability:

| | | |
|---|---|---|
| NP → PN | [0.35] |
| NP → PRN | [0.30] |
| NP → DT N | [0.20] |
| NP → N | [0.15] |
| DT → *that* | [0.10] |
| DT → *a* | [0.30] |
| DT → *the* | [0.60] |

The total probability for rules with the same left-hand side is 1:

$$\sum_{\beta} P(X \rightarrow \beta) = 1$$

Where:

- $\beta$ is any sequence of (terminal/non-terminal) symbols.
- $P(X \rightarrow \beta)$ = probability of rule $X \rightarrow \beta$
  For example: $P(\text{NP} \rightarrow \text{PN}) = 0.35$

## Statistical Parsing

$T$ is a parse tree and $S = w_1...w_n$ a token sequence. Bayes' rule:

$$P(T|w_1...w_n) = \frac{P(w_1...w_n|T) \times P(T)}{P(w_1...w_n)}$$

But since $P(w_1...w_n|T)$ is always 1, then:

$$P(T|w_1...w_n) = \frac{P(T)}{P(w_1...w_n)}$$

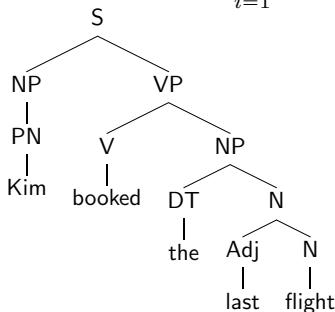In particular, we are interested in the most likely $T$ for $S$:

$$\hat{T} = argmax_{T:S=yield(T)} \frac{P(T)}{P(w_1...w_n)}$$

But since $P(w_1...w_n)$ is constant over all $T$'s for a given $S$ then:

$$\hat{T} \approx argmax_{T:S=yield(T)} P(T)$$

## Statistical Parsing

$$P(T) = \prod_{i=1}^{n} P(X_i \rightarrow Y_1...Y_n)$$



| Rule | P | Rule | P | Rule | P | Rule | P |
|------|---|------|---|------|---|------|---|
| S → NP VP | 0.8 | S → VP | 0.2 | VP → V | 0.3 | VP → V NP | 0.5 |
| VP → V NP NP | 0.2 | NP → DT N | 0.2 | NP → PN | 0.35 | NP → PRN | 0.3 |
| NP → N | 0.15 | N → Adj N | 0.4 | N → evening | 0.2 | N → flight | 0.1 |
| PN → Kim | 0.15 | V → booked | 0.4 | Adj → last | 0.1 | DT → the | 0.4 |

$$P(T) = 0.8 \times 0.35 \times 0.15 \times 0.5 \times 0.4 \times 0.2 \times 0.4 \times 0.4 \times 0.1 \times 0.1 = 2.7 \times 10^{-6}$$

# Statistical Parsing

We can use Treebanks to estimate the probabilities for CFG rules:

$$P(X \rightarrow Y_1...Y_n) = \frac{Count(X \rightarrow Y_1...Y_n)}{Count(X)}$$

**Example:**

$[_S$ $[_{NP}$ $[_{DT}$ This] $[_N$ text]] $[_{VP}$ $[_V$ is] $[[_{Adv}$ just ] $[_{NP}$ $[_{DT}$ an] $[_N$ example]]]]]. $[_S$ $[_{NP}$ $[_{PRN}$ I]] $[_{VP}$ $[_V$ made] $[_{NP}$ $[_{PRN}$ it]] $[_{RP}$ up]]].

$P(\text{NP} \rightarrow \text{DT N}) = \frac{2}{4} = 0.5$

# Statistical Parsing

Penn Treebank:

```
1  (S (NP-SBJ-1 Jones)
2      (VP followed
3          (NP him)
4          (PP-DIR into
5                (NP the front room))))
```

**Treebanks**

- Linguistic Data Consortium (LDC)
- European Language Resources Association (ELRA)
- Stanford list
- NLTK data
- Others

# Statistical Parsing

Or you can use a statistical parser to automatically parse a corpus you wish to use.

- Stanford Parser:

```
1  java -cp stanford-parser.jar:stanford-parser-3.4.1-models.jar
2  edu.stanford.nlp.parser.lexparser.LexicalizedParser -outputFormat penn
3  edu/stanford/nlp/models/lexparser/englishPCFG.ser.gz  input.txt >
4  output.txt
```

- Berkeley Parser: (demo)

```
1  java -jar berkeleyParser/BerkeleyParser-1.7.jar -gr
2  berkeleyParser/eng_sm6.gr < input.txt > output.txt
```

- More here

**Parsing PCFG with the CKY parser**

We need:

- A CNF grammar augmented with probabilities
- A way to resolve ambiguity:
  if there are two categories of the same type in the same span,
  then discard the less likely one.

| | | | | | | |
|---|---|---|---|---|---|---|
| S[.7×.1×.00001=.0000007] | | | | | | |
| S[.7×.1×.0001=.000007] | VP[.5×.2×.0001=.00001] VP[.1×.006×.01=.000006] | | | | | |
| – | VP[.5×.2×.001=.0001] VP[.1×.006×.072=.00001] | NP[.5×.4×.0009=.0001] | | | | |
| S[.7×.1×.006=.0004] | – | NP[.5×.4×.006=.001] | N[.3×.3×.01=.0009] | | | |
| – | VP[.5×.2×.06=.006] | – | N[.3×.3×.072=.006] | PP[.9×.4×.03=.01] | | |
| – | – | NP[.5×.4×.3=.06] | – | PP[.9×.4×.2=.072] | NP[.5×.2×.3=.03] | |
| NP[.1] | Adj[.1] TV[.2] | DT[.4] | N[.3] | P[.4] | DT[.2] NP[.2] | TV[.05] N[.3] |
| Mary | attacked | a | farmer | with | her | axe |

$S \rightarrow NP\ VP$ [.7]   $VP \rightarrow DVP\ NP$ [.3]   $VP \rightarrow VP\ PP$ [.1]

$VP \rightarrow TV\ NP$ [.5]   $NP \rightarrow DT\ N$ [.5]   $NP \rightarrow Mary$ [.1]

$NP \rightarrow her$ [.2]   $DT \rightarrow a$ [.4]   $DT \rightarrow her$ [.2]

$N \rightarrow Adj\ N$ [.1]   $N \rightarrow N\ PP$ [.3]   $N \rightarrow farmer$ [.3]

$N \rightarrow axe$ [.3]   $Adj \rightarrow attacked$ [.1]   $PP \rightarrow P\ NP$ [.9]

$P \rightarrow with$ [.4]   $TV \rightarrow attacked$ [.2]   $TV \rightarrow axe$ [.05]

# Statistical Parsing

**Evaluating PCFGS** (PARSEVAL)

How do the constituents in the hypothesis parse tree match the constituents in a hand-labeled 'gold standard' (reference) parse tree?

$T_C$ = set of constituents for S according to reference
$T_G$ = set of constituents hypothesized for S

1. **labeled precision**

$$\frac{|T_C \cap T_G|}{|T_G|} = \frac{\#\text{of correctly identified constituents}}{\#\text{of constituents hypothesized}} = \frac{t_p}{t_p + f_p}$$

2. **labeled recall**

$$\frac{|T_C \cap T_G|}{|T_C|} = \frac{\#\text{of correctly identified constituents}}{\#\text{of constituents in reference}} = \frac{t_p}{t_p + f_n}$$

**Rule of thumb**: as precision increases, recall drops and vice versa.

Often precision and recall are reported as a single number:

$$\textbf{F-measure}: F_\beta = \frac{(\beta^2 + 1) \times P \times R}{\beta^2 \times P + R}$$

$\beta > 1$ favors Recall
$\beta < 1$ favors Precision
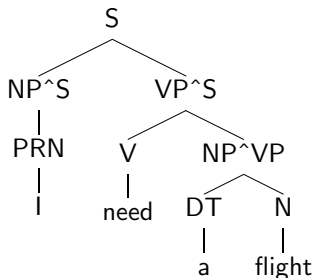
# Statistical Parsing

**Poor independence assumption**

Linguistic structures are not independent from each other. For example, the distribution of NP → PN is unbalanced:

91% of subject phrases are pronouns

34% of object phrases are pronouns

Solution: add information about the mother node in the daughters

*Parent annotation*: NP^S → PRN [.91] vs. NP^VP → PRN [.34]

```
                    S
           ┌────────┴────────┐
        NP^S              VP^S
          │            ┌────┴────┐
        PRN            V       NP^VP
          │            │      ┌──┴──┐
          I          need    DT     N
                             │      │
                             a    flight
```

# Statistical Parsing

PROBLEMS WITH PCFGS (continued)

**Lack of lexical conditioning**
Lexical items are important to resolve attachment ambiguities.

(2) a. *[Sam [dumped [the box into the bin]]].

   b. [Sam [dumped [the box] [into the bin]]].

(3) a. [Sam [dumped [the box in the bin]]].

   b. [Sam [dumped [the box] [in the bin]]].

(4) a. Sam likes [[green vegetables] and [music]].

   b. *Sam likes [green [vegetables and music]].

(5) a. *I need some [fresh [air and sunshine]].

   b. I need some [[fresh air] and [sunshine]].

Solution: add information about the token in the mother node

# Statistical Parsing

**Dealing with the lack of lexical conditioning**

*grammar lexicalizaton*