

(概率)上下文无关文法

上下文无关文法(CFG)

上下文无关文法，Context Free Grammar，CFG是一个四元组：

$G = (N, T, P, S)$ ，其中

- N 是**非终结符**（Nonterminals）的有限集合；
- T 是**终结符**（Terminals）的有限集合，且 $N \cap T = \Phi$ ；
- P 是**产生式**（Productions）的有限集合，形如：
 $A \rightarrow \alpha$ ，其中 $A \in N$ （左部）， $\alpha \in (N \cup T)^*$ （右部），
若 $\alpha = \varepsilon$ ，则称 $A \rightarrow \varepsilon$ 为空产生式（也可以记为 $A \rightarrow$ ）；
- S 是非终结符，称为文法的**开始符号**（Start symbol）。

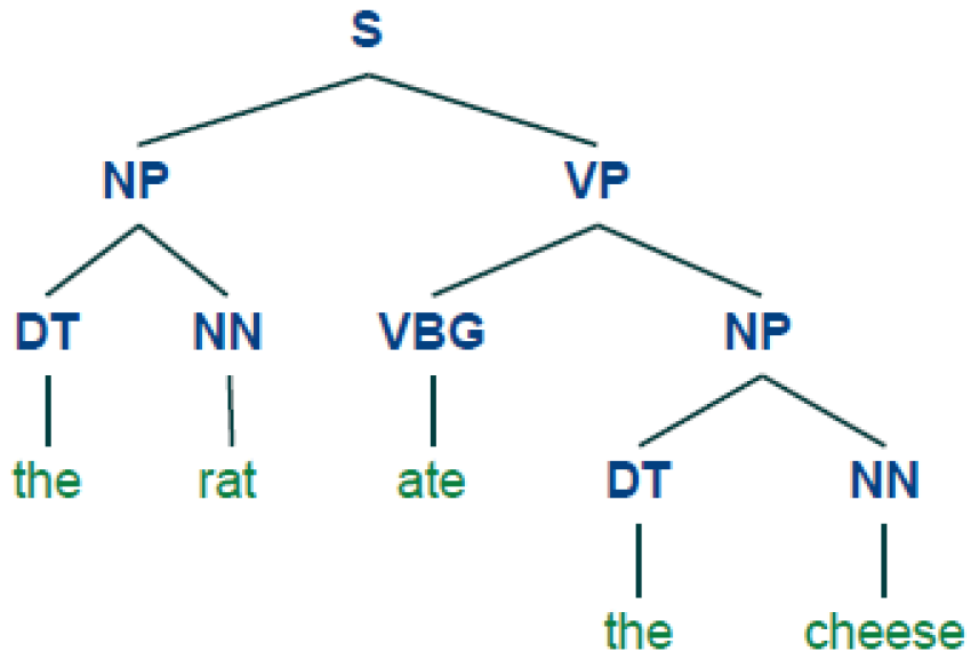
- 终结符集: rat, the, ate, cheese
- 非终结行集: S, NP, VP, DT, VBD, NN
- 产生式集合:
 - $S \rightarrow NP VP$
 - $NP \rightarrow DT NN$
 - $VP \rightarrow VBD NP$
 - $DT \rightarrow the$
 - $NN \rightarrow rat$
 - $NN \rightarrow cheese$
 - $VBD \rightarrow ate$
- 开始符号: S

$S \rightarrow NP VP$
 $NP \rightarrow DT NN$
 $VP \rightarrow VBD NP$
 $DT \rightarrow the$
 $NN \rightarrow rat$
 $NN \rightarrow cheese$
 $VBD \rightarrow ate$

- 求句子the rat ate the cheese的分析树

$S \rightarrow NP VP$
 $NP \rightarrow DT NN$
 $VP \rightarrow VBD NP$
 $DT \rightarrow the$
 $NN \rightarrow rat$
 $NN \rightarrow cheese$
 $VBD \rightarrow ate$

- 求句子the rat ate the cheese的分析树



句法分析 (Syntactic Parsing)

- 句法解析就是指在给定词串的情况下，例如the rat ate the cheese，来识别其可能的句法结构。

S → NP

S → VP

VP → VC NP

NP → DJ NP

NP → noun

DJ → VP de

DJ → NP de

VC → vt adj

VC → VC utl

咬 死 了 猎 人 的 狗

vt → 咬

adj → 死

utl → 了

noun → 猎人

noun → 狗

de → 的

S → NP

S → VP

VP → VC NP

NP → DJ NP

NP → noun

DJ → VP de

DJ → NP de

VC → vt adj

VC → VC utl

vt → 咬

adj → 死

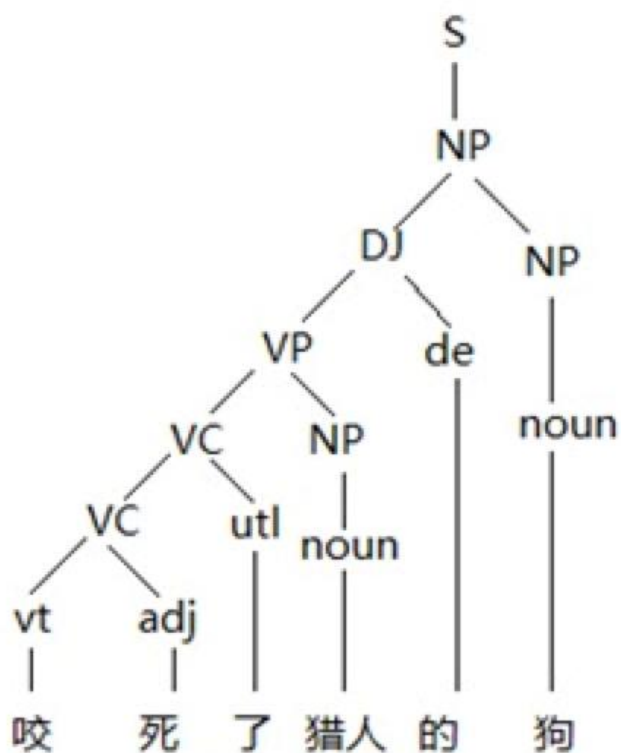
utl → 了

noun → 猎人

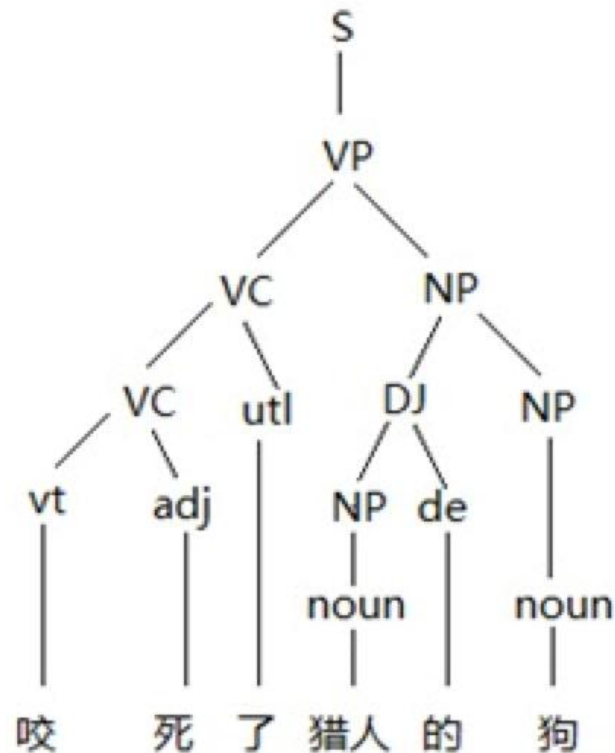
noun → 狗

de → 的

咬 死 了 猎 人 的 狗



(1)



(2)

怎么判断哪个是合理的??

概率上下文无关文法(CFG)

概率上下文无关文法, Probabilistic Context Free Grammar, CFG是一个四元组:

$G = (N, T, P, S)$, 其中

- N 是非终结符 (Nonterminals) 的有限集合;
- T 是终结符 (Terminals) 的有限集合, 且 $N \cap T = \Phi$;
- P 是产生式 (Productions) 的有限集合, 形如:
 $A \rightarrow \alpha [p]$, 其中 $A \in N$ (左部), $\alpha \in (N \cup T)^*$ (右部),
若 $\alpha = \varepsilon$, 则称 $A \rightarrow \varepsilon$ 为空产生式 (也可以记为 $A \rightarrow$),
 p 是0~1之间的值, 表示该产生式的概率 $P(\alpha|A)$;
- S 是非终结符, 称为文法的开始符号 (Start symbol)。

- $A \rightarrow \alpha$ [p]
 - 产生式左边 LHS (left hand side)
 - 产生式右边 RHS (right hand side)
-
- ▶ $NP \rightarrow DT\ NN$ [p = 0.45]
 - ▶ $NN \rightarrow leprechaun$ [p = 0.0001]

对于一个给定的LHS要满足加和等于1。

$S \rightarrow NP VP$	[.80]
$S \rightarrow Aux NP VP$	[.15]
$S \rightarrow VP$	[.05]
$NP \rightarrow Pronoun$	[.35]
$NP \rightarrow Proper-Noun$	[.30]
$NP \rightarrow Det Nominal$	[.20]
$NP \rightarrow Nominal$	[.15]
$Nominal \rightarrow Noun$	[.75]
$Nominal \rightarrow Nominal Noun$	[.20]
$Nominal \rightarrow Nominal PP$	[.05]
$VP \rightarrow Verb$	[.35]
$VP \rightarrow Verb NP$	[.20]
$VP \rightarrow Verb NP PP$	[.10]
$VP \rightarrow Verb PP$	[.15]
$VP \rightarrow Verb NP NP$	[.05]
$VP \rightarrow VP PP$	[.15]
$PP \rightarrow Preposition NP$	[1.0]

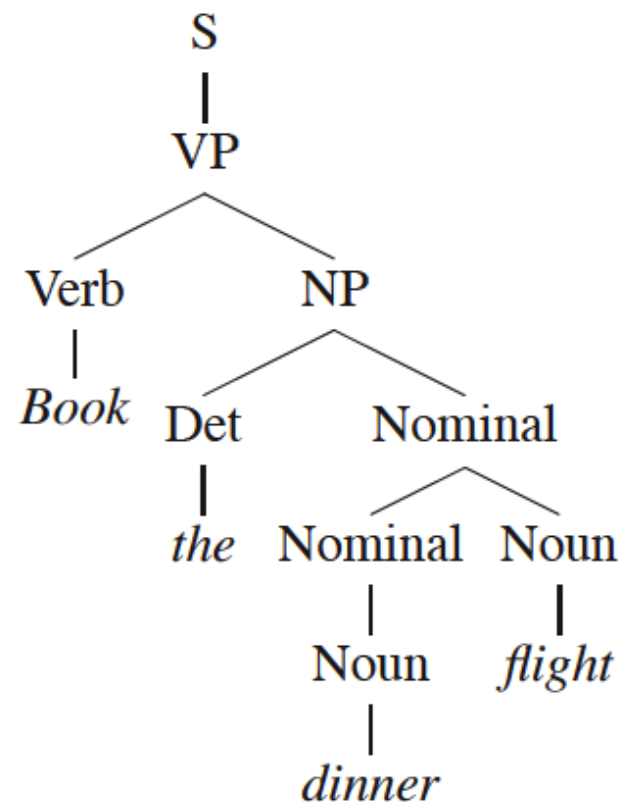
$Det \rightarrow that$	[.10]		a	[.30]		the	[.60]
$Noun \rightarrow book$	[.10]		$flight$	[.30]			
			$meal$	[.15]		$money$	[.05]
			$flights$	[.40]		$dinner$	[.10]
$Verb \rightarrow book$	[.30]		$include$	[.30]			
			$prefer$	[.40]			
$Pronoun \rightarrow I$	[.40]		she	[.05]			
			me	[.15]		you	[.40]
$Proper-Noun \rightarrow Houston$	[.60]						
			TWA	[.40]			
$Aux \rightarrow does$	[.60]		can	[.40]			
$Preposition \rightarrow from$	[.30]		to	[.30]			
			on	[.20]		$near$	[.15]
			$through$	[.05]			

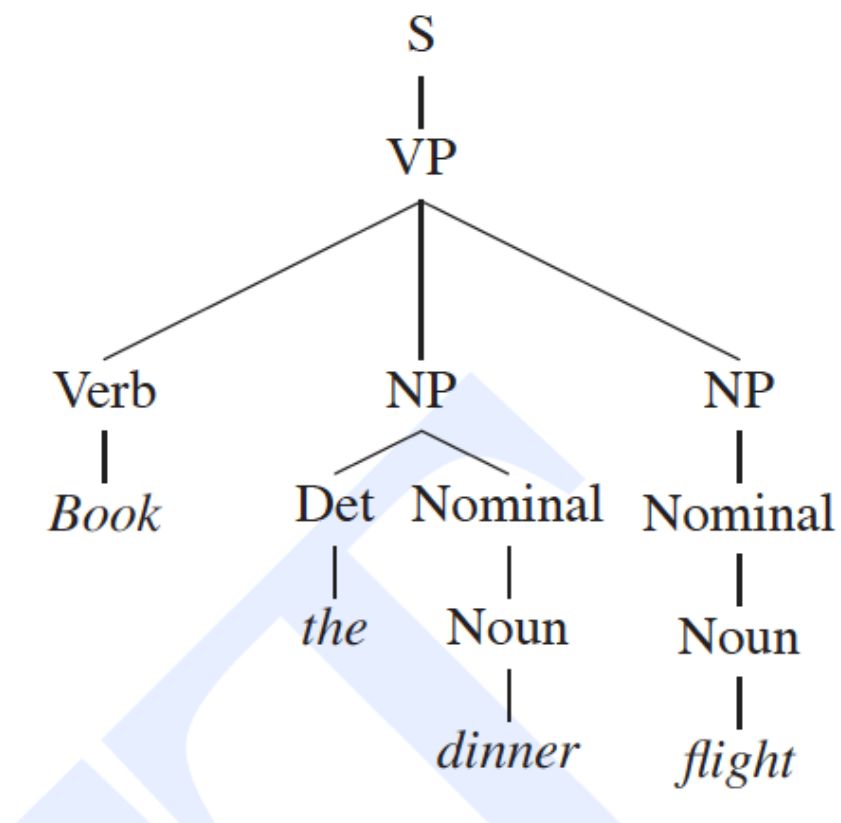
$$P(T) = \prod_{i=1}^n P(\text{RHS}_i | \text{LHS}_i)$$

句法树T的概率

一般来说，概率值大的更可能是正确的句法树。

$$\begin{aligned}
P(T) = & P(S \rightarrow VP) \times \\
& P(VP \rightarrow \text{Verb NP}) \times \\
& P(\text{Verb} \rightarrow \textit{Book}) \times \\
& P(NP \rightarrow \text{Det Nominal}) \times \\
& P(\text{Det} \rightarrow \textit{the}) \times \\
& P(\text{Nominal} \rightarrow \text{Nominal Noun}) \times \\
& P(\text{Noun} \rightarrow \textit{dinner}) \times \\
& P(\text{Noun} \rightarrow \textit{flight}) = 2.2 \times 10^{-6}
\end{aligned}$$





$$P = 6.1 \times 10^{-7}$$

- 给定一个PCFG (概率上下文无关文法)和一个句子, 如何构造该句子概率最大的句法树?
- CNF (Chomsky Normal Form)
- CYK算法

CNF (Chomsky Normal Form)

- 如果一个上下文无关文法的每个产生式的形式为：

$$A \rightarrow BC \text{ 或 } A \rightarrow a$$

即规则的右部或者是两个非终结符或者是一个终结符.

- 任何CFG都可以转变成一个弱等价的Chomsky范式语法。

CYK算法 (CKY)

- 给定一个句子 $s = w_1 w_2, \dots, w_n$, 和一个上下文无关文法PCFG, $G=(T, N, S, R)$;
- 定义一个跨越单词 i 到 j 的概率最大的语法成分 π :

$$\pi(i,j,X) \quad (i, j \in 1 \dots n, X \in N),$$

- 目标是找到一个属于 $\pi [1, n, S]$ 的所有树中概率最大的那棵。

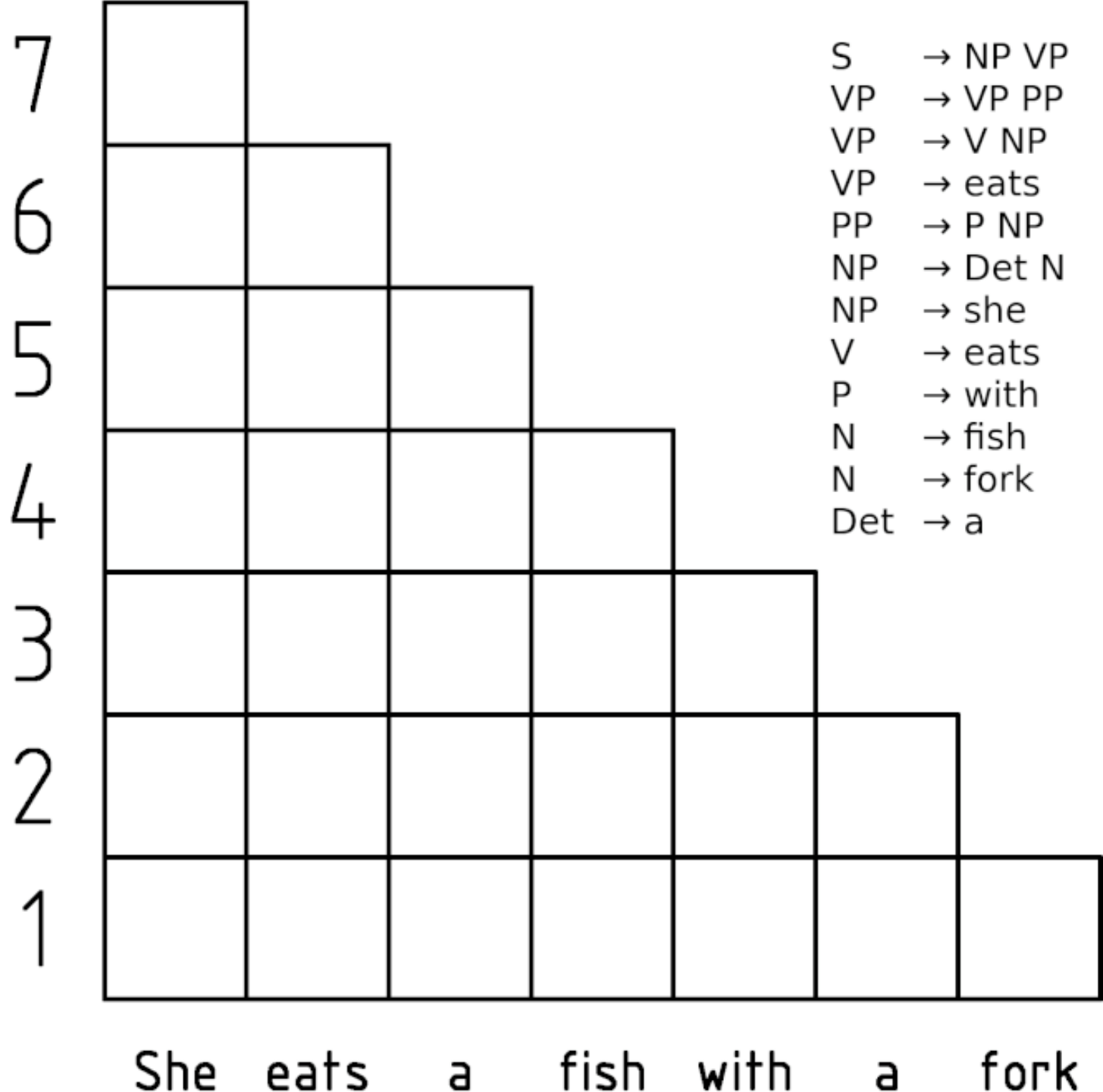
CYK算法用于PCFG下的句法分析：

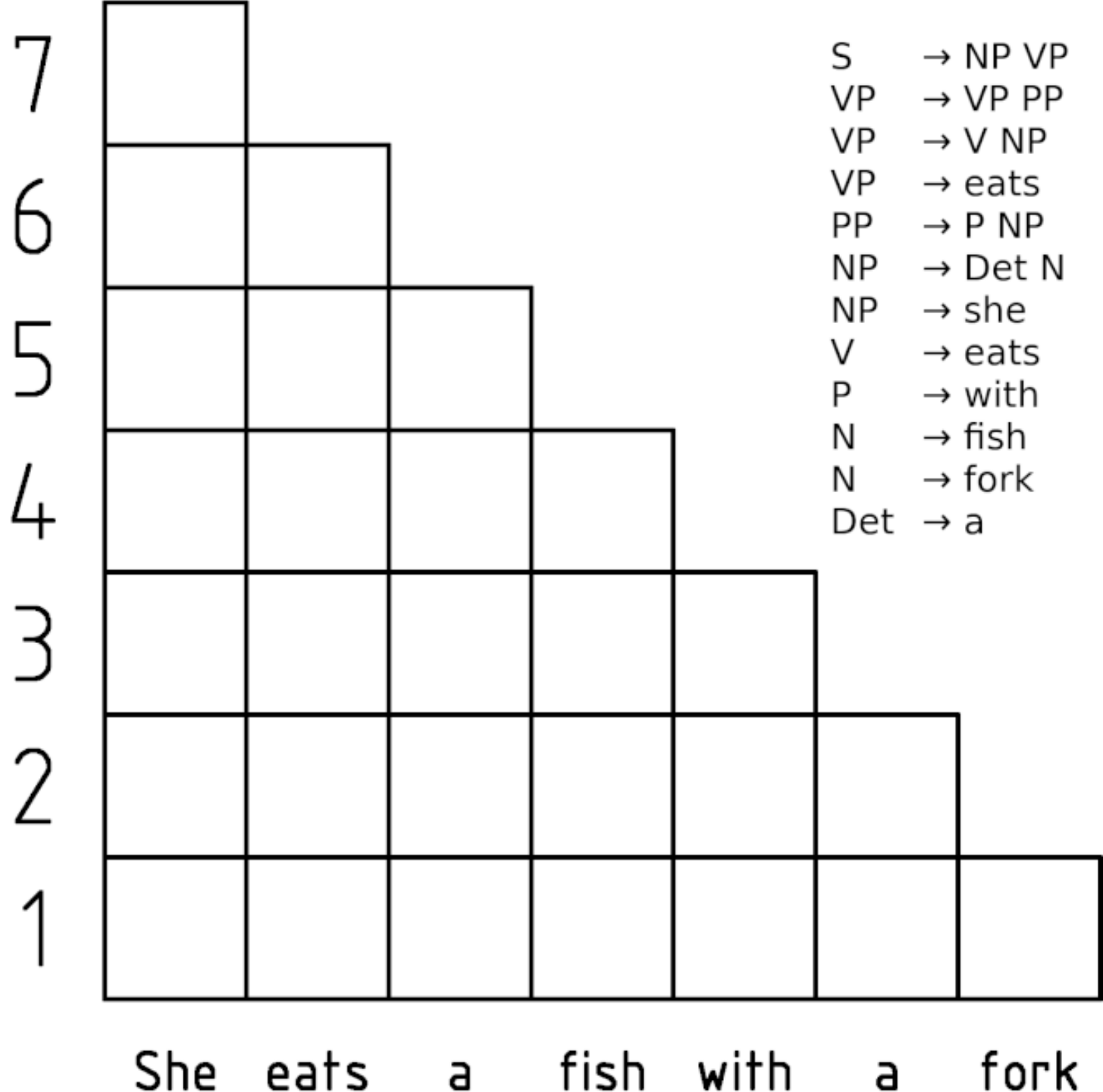
- 基本定义： for all $i=1, \dots, n$, $X \in N$

$$\pi(i, i, X) = q(X \rightarrow w_i) \text{ (if } X \rightarrow w_i \text{ 没有出现在语法中, 则定义 } q(X \rightarrow w_i) = 0)$$

- 递归定义： for all $i=1, \dots, n$, $j=(i+1), \dots, n$, $X \in N$

$$\pi(i, j, X) = \max(q(X \rightarrow YZ) \times \pi(i, k, Y) \times \pi(k+1, j, Z)) \text{ (} i \leq k \leq j-1)$$





S → NP VP	0.9		
S → VP	0.1	N → <i>people</i>	0.5
VP → V NP	0.5	N → <i>fish</i>	0.2
VP → V	0.1	N → <i>tanks</i>	0.2
VP → V @VP_V	0.3	N → <i>rods</i>	0.1
VP → V PP	0.1	V → <i>people</i>	0.1
@VP_V → NP PP	1.0	V → <i>fish</i>	0.6
NP → NP NP	0.1	V → <i>tanks</i>	0.3
NP → NP PP	0.2	P → <i>with</i>	1.0
NP → N	0.7		
PP → P NP	1.0		

求 fish people fish tanks的最优分析树.