

Uvod u obradu prirodnog jezika

14.1. Probabilističko parsiranje (Probabilistic parsing)

Branko Žitko

prevedeno od: Dan Jurafsky, Chris Manning

Gramatika strukture fraze

$S \rightarrow NP VP$

$VP \rightarrow V NP$

$VP \rightarrow V NP PP$

$NP \rightarrow NP NP$

$NP \rightarrow NP PP$

$NP \rightarrow N$

$NP \rightarrow \varepsilon$

$PP \rightarrow P NP$

$N \rightarrow \text{primati}$

$N \rightarrow \text{kape}$

$N \rightarrow \text{nose}$

$N \rightarrow \text{glavi}$

$V \rightarrow \text{primati}$

$V \rightarrow \text{kape}$

$V \rightarrow \text{nose}$

$P \rightarrow \text{na}$

primati kape nose

primati kape na glavi

Gramatike strukture fraze

- Kontekstno neovisne gramatike
Context Free Grammars (CFG)
- $G = (T, N, S, R)$
 - T – skup terminalnih simbola
 - N – skup neterminalnih simbola
 - S – početni simbol ($S \in N$)
 - R – skup pravila/produkcija oblika $X \rightarrow \gamma$
 $X \in N, \gamma \in (N \cup T)^*$
- Gramatika G generira jezik L

Gramatike strukture fraze u obradi prirodnog jezika

- $G = (T, C, N, S, L, R)$
 - T – skup terminalnih simbola
 - C – skup preterminalnih simbola
 - N – skup neterminalnih simbola
 - S – početni simbol ($S \in N$)
 - L – leksikon, skup elemenata obila $X \rightarrow x$
 $X \in N, x \in T$
 - R – skup pravila/produkcija oblika $X \rightarrow \gamma$
 $X \in N, \gamma \in (N \cup T)^*$
- Gramatika G generira jezik L

$S \rightarrow NP VP$ $VP \rightarrow V NP$ $VP \rightarrow V NP PP$ $NP \rightarrow NP NP$ $NP \rightarrow NP PP$ $NP \rightarrow N$ $NP \rightarrow \varepsilon$ $PP \rightarrow P NP$ $N \rightarrow \textit{primati}$ $N \rightarrow \textit{kape}$ $N \rightarrow \textit{nose}$ $N \rightarrow \textit{glavi}$ $V \rightarrow \textit{primati}$ $V \rightarrow \textit{kape}$ $V \rightarrow \textit{nose}$ $P \rightarrow \textit{na}$

primati kape nose

primati kape na glavi

Probabilistic CFG (PCFG)

- $G = (T, N, S, R, P)$
 - T – skup terminalnih simbola
 - N – skup neterminalnih simbola
 - S – početni simbol ($S \in N$)
 - R – skup pravila/produkcija oblika $X \rightarrow \gamma$
 $X \in N, \gamma \in (N \cup T)^*$
 - P – probabilistička funkcija
 $P : R \rightarrow [0, 1]$
$$\forall X \in N, \sum_{X \rightarrow \gamma \in R} P(X \rightarrow \gamma) = 1$$
- Gramatika G generira model jezika L

$$\sum_{\gamma \in T^*} P(\gamma) = 1$$

PCFG

SENT \rightarrow NP VP 1.0

VP \rightarrow V NP 0.6

VP \rightarrow V NP PP 0.4

NP \rightarrow NP NP 0.1

NP \rightarrow NP PP 0.2

NP \rightarrow N 0.7

~~NP \rightarrow ϵ 0.0~~

PP \rightarrow P NP 1.0

N \rightarrow primati 0.5

N \rightarrow kape 0.2

N \rightarrow nose 0.2

N \rightarrow glavi 0.1

V \rightarrow primati 0.1

V \rightarrow kape 0.6

V \rightarrow nose 0.3

S \rightarrow na 1.0

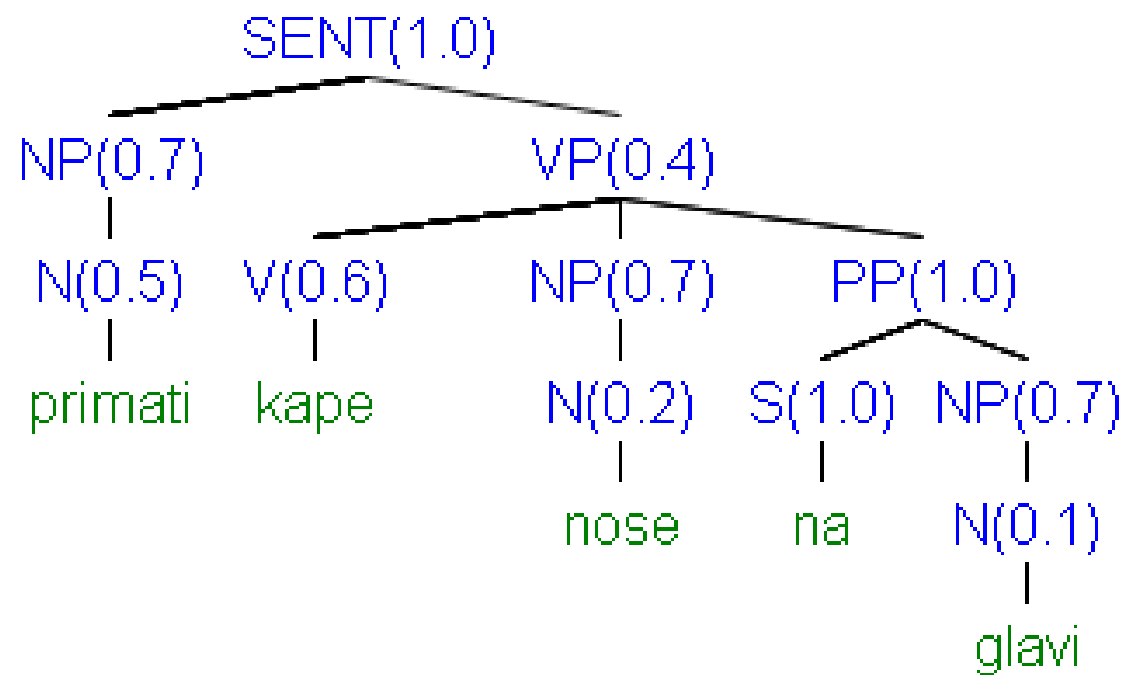
primati kape nose

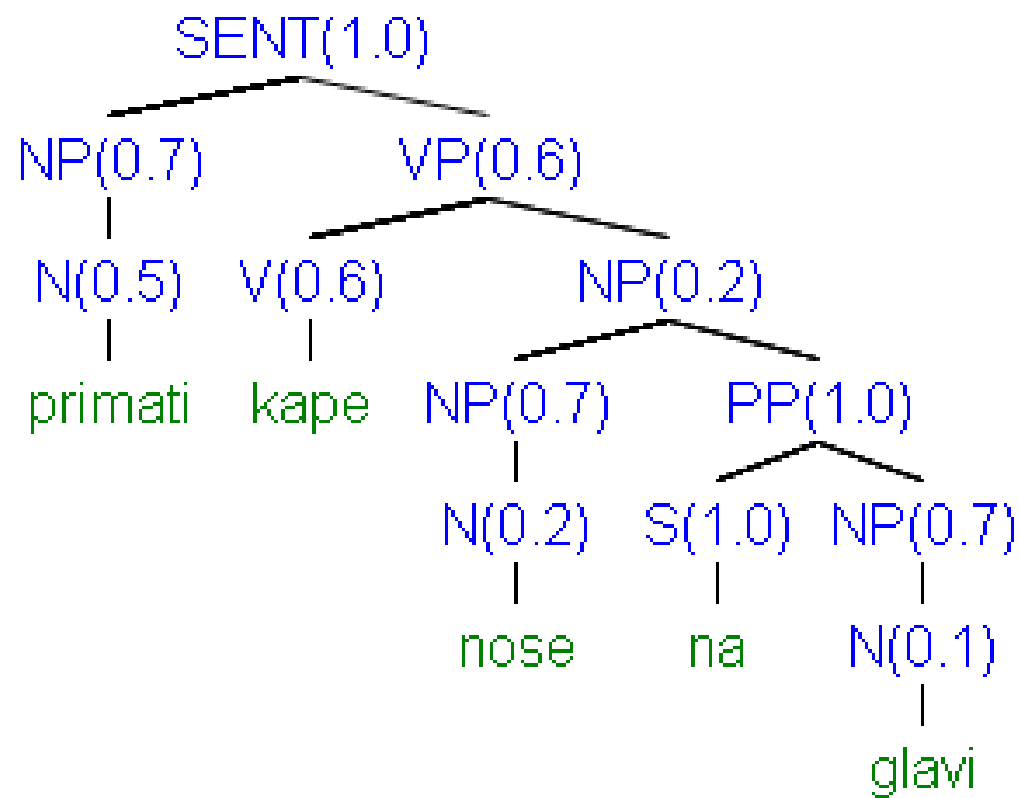
primati kape na glavi

Vjerojatnosti stabala i nizova riječi

- $P(t)$ – vjerojatnost stabla t je umnožak vjerojatnosti pravila korištenih za generiranje stabla
- $P(s)$ – vjerojatnost niza riječi s je suma vjerojatnosti stabala koji generiraju s

$$\begin{aligned} P(s) &= \sum_t P(s, t) \text{ gdje je } t \text{ stablo parsiranja od } s \\ &= \sum_t P(t) \end{aligned}$$





Vjerojatnosti stabla i niza riječi

- $s = \textit{primati kape nose na glavi}$
- $P(t_1) = 1.0 \times 0.7 \times 0.4 \times 0.5 \times 0.6 \times 0.7$ Spajanje na glagol
 $\times 1.0 \times 0.2 \times 1.0 \times 0.7 \times 0.1$
 $= 0.0008232$
- $P(t_2) = 1.0 \times 0.7 \times 0.6 \times 0.5 \times 0.6 \times 0.2$ Spajanje na imenicu
 $\times 0.7 \times 1.0 \times 0.2 \times 1.0 \times 0.7 \times 0.1$
 $= 0.00024696$
- $P(s) = P(t_1) + P(t_2)$
 $= 0.0008232 + 0.00024696$
 $= 0.00107016$

Uvod u obradu prirodnog jezika

14.2. Transformacija gramatike

Branko Žitko

prevedeno od: Dan Jurafsky, Chris Manning

Chomsky normalna forma

- Sva pravila su oblika $X \rightarrow YZ$ ili $X \rightarrow w$
 $X, Y, Z \in N, w \in T$
- Transformacija u ovu formu ne mijenja slabo generativno svojstvo CFG
 - Odnosno, prepoznaje isti jezika, ali možda s različitim stablima
- Prazna i unarna pravila se rekurzivno izbacuju
- n-arna pravila se dijele uvođenjem novih neterminala ($n > 2$)

$$S \rightarrow NP VP$$
$$VP \rightarrow V NP$$
$$VP \rightarrow V NP PP$$
$$NP \rightarrow NP NP$$
$$NP \rightarrow NP PP$$
$$NP \rightarrow N$$
$$NP \rightarrow \varepsilon$$
$$PP \rightarrow P NP$$
$$N \rightarrow \textit{primati}$$
$$N \rightarrow \textit{kape}$$
$$N \rightarrow \textit{nose}$$
$$N \rightarrow \textit{glavi}$$
$$V \rightarrow \textit{primati}$$
$$V \rightarrow \textit{kape}$$
$$V \rightarrow \textit{nose}$$
$$P \rightarrow \textit{na}$$

Chomsky – eliminacija ϵ

$S \rightarrow NP VP$

$S \rightarrow NP VP$
 $S \rightarrow VP$

$VP \rightarrow V NP$

$VP \rightarrow V NP$
 $VP \rightarrow NP$

$VP \rightarrow V NP PP$

$VP \rightarrow V NP PP$
 $VP \rightarrow V PP$

$NP \rightarrow NP NP$

$NP \rightarrow NP NP$
 $NP \rightarrow NP$

$NP \rightarrow NP PP$

$NP \rightarrow NP PP$
 $NP \rightarrow PP$

$NP \rightarrow N$

~~$NP \rightarrow \epsilon$~~

$PP \rightarrow P NP$

$PP \rightarrow NP PP$
 $PP \rightarrow P$

$N \rightarrow \textit{primati}$

$N \rightarrow \textit{kape}$

$N \rightarrow \textit{nose}$

$N \rightarrow \textit{glavi}$

$V \rightarrow \textit{primati}$

$V \rightarrow \textit{kape}$

$V \rightarrow \textit{nose}$

$P \rightarrow \textit{na}$

Chomsky – eliminacija unarnih pravila

$S \rightarrow NP VP$

~~$S \rightarrow VP$~~

$VP \rightarrow V NP$

$VP \rightarrow V$

$VP \rightarrow V NP PP$

$VP \rightarrow V PP$

$NP \rightarrow NP NP$

$NP \rightarrow NP$

$NP \rightarrow NP PP$

$NP \rightarrow PP$

$NP \rightarrow N$

$PP \rightarrow P NP$

$PP \rightarrow P$

$S \rightarrow V NP$
 $S \rightarrow V$
 $S \rightarrow V NP PP$
 $S \rightarrow V PP$

$N \rightarrow \textit{primati}$

$N \rightarrow \textit{kape}$

$N \rightarrow \textit{nose}$

$N \rightarrow \textit{glavi}$

$V \rightarrow \textit{primati}$

$V \rightarrow \textit{kape}$

$V \rightarrow \textit{nose}$

$P \rightarrow \textit{na}$

Chomsky – eliminacija unarnih pravila

$S \rightarrow NP VP$

$VP \rightarrow V NP$

$S \rightarrow V NP$

$VP \rightarrow V$

~~$S \rightarrow V$~~

$VP \rightarrow V NP PP$

$S \rightarrow V NP PP$

$VP \rightarrow V PP$

$S \rightarrow V PP$

$NP \rightarrow NP NP$

$NP \rightarrow NP$

$NP \rightarrow NP PP$

$NP \rightarrow PP$

$NP \rightarrow N$

$PP \rightarrow P NP$

$PP \rightarrow P$

$N \rightarrow \textit{primati}$

$N \rightarrow \textit{kape}$

$N \rightarrow \textit{nose}$

$N \rightarrow \textit{glavi}$

$V \rightarrow \textit{primati}$

$V \rightarrow \textit{kape}$

$V \rightarrow \textit{nose}$

$P \rightarrow \textit{na}$

$S \rightarrow \textit{primati}$

$S \rightarrow \textit{kape}$

$S \rightarrow \textit{nose}$

Chomsky – eliminacija unarnih pravila

$S \rightarrow NP VP$

$VP \rightarrow V NP$

$S \rightarrow V NP$

~~$VP \rightarrow V$~~

$VP \rightarrow V NP PP$

$S \rightarrow V NP PP$

$VP \rightarrow V PP$

$S \rightarrow V PP$

$NP \rightarrow NP NP$

$NP \rightarrow NP$

$NP \rightarrow NP PP$

$NP \rightarrow PP$

$NP \rightarrow N$

$PP \rightarrow P NP$

$PP \rightarrow P$

$VP \rightarrow primati$

$VP \rightarrow kape$

$VP \rightarrow nose$

$N \rightarrow primati$

$N \rightarrow kape$

$N \rightarrow nose$

$N \rightarrow glavi$

$V \rightarrow primati$

$S \rightarrow primati$

$V \rightarrow kape$

$S \rightarrow kape$

$V \rightarrow nose$

$S \rightarrow nose$

$P \rightarrow na$

Chomsky – eliminacija unarnih pravila

$S \rightarrow NP VP$

$VP \rightarrow V NP$

$S \rightarrow V NP$

$VP \rightarrow V NP PP$

$S \rightarrow V NP PP$

$VP \rightarrow V PP$

$S \rightarrow V PP$

$NP \rightarrow NP NP$

~~$NP \rightarrow NP$~~

$NP \rightarrow NP PP$

$NP \rightarrow PP$

$NP \rightarrow N$

$PP \rightarrow P NP$

$PP \rightarrow P$

$N \rightarrow \textit{primati}$

$N \rightarrow \textit{kape}$

$N \rightarrow \textit{nose}$

$N \rightarrow \textit{glavi}$

$V \rightarrow \textit{primati}$

$S \rightarrow \textit{primati}$

$VP \rightarrow \textit{primati}$

$V \rightarrow \textit{kape}$

$S \rightarrow \textit{kape}$

$VP \rightarrow \textit{kape}$

$V \rightarrow \textit{nose}$

$S \rightarrow \textit{nose}$

$VP \rightarrow \textit{nose}$

$P \rightarrow \textit{na}$

Chomsky – eliminacija unarnih pravila

$S \rightarrow NP VP$

$VP \rightarrow V NP$

$S \rightarrow V NP$

$VP \rightarrow V NP PP$

$S \rightarrow V NP PP$

$VP \rightarrow V PP$

$S \rightarrow V PP$

$NP \rightarrow NP NP$

$NP \rightarrow NP PP$

~~$NP \rightarrow PP$~~

$NP \rightarrow N$

$PP \rightarrow P NP$

$PP \rightarrow P$

$PP \rightarrow P NP$

$PP \rightarrow P$

$N \rightarrow \textit{primati}$

$N \rightarrow \textit{kape}$

$N \rightarrow \textit{nose}$

$N \rightarrow \textit{glavi}$

$V \rightarrow \textit{primati}$

$S \rightarrow \textit{primati}$

$VP \rightarrow \textit{primati}$

$V \rightarrow \textit{kape}$

$S \rightarrow \textit{kape}$

$VP \rightarrow \textit{kape}$

$V \rightarrow \textit{nose}$

$S \rightarrow \textit{nose}$

$VP \rightarrow \textit{nose}$

$P \rightarrow \textit{na}$

Chomsky – eliminacija unarnih pravila

$S \rightarrow NP VP$
 $VP \rightarrow V NP$
 $S \rightarrow V NP$
 $VP \rightarrow V NP PP$
 $S \rightarrow V NP PP$
 $VP \rightarrow V PP$
 $S \rightarrow V PP$
 $NP \rightarrow NP NP$
 $NP \rightarrow NP PP$
 ~~$NP \rightarrow N$~~
 $PP \rightarrow P NP$
 $NP \rightarrow P NP$
 $PP \rightarrow P$
 $NP \rightarrow P$

$NP \rightarrow \textit{primati}$
 $NP \rightarrow \textit{kape}$
 $NP \rightarrow \textit{nose}$
 $NP \rightarrow \textit{glavi}$

~~$N \rightarrow \textit{primati}$~~
 ~~$N \rightarrow \textit{kape}$~~
 ~~$N \rightarrow \textit{nose}$~~
 ~~$N \rightarrow \textit{glavi}$~~
 $V \rightarrow \textit{primati}$
 $S \rightarrow \textit{primati}$
 $VP \rightarrow \textit{primati}$
 $V \rightarrow \textit{kape}$
 $S \rightarrow \textit{kape}$
 $VP \rightarrow \textit{kape}$
 $V \rightarrow \textit{nose}$
 $S \rightarrow \textit{nose}$
 $VP \rightarrow \textit{nose}$
 $P \rightarrow \textit{na}$

Chomsky – eliminacija unarnih pravila

$S \rightarrow NP VP$

$VP \rightarrow V NP$

$S \rightarrow V NP$

$VP \rightarrow V NP PP$

$S \rightarrow V NP PP$

$VP \rightarrow V PP$

$S \rightarrow V PP$

$NP \rightarrow NP NP$

$NP \rightarrow NP PP$

$PP \rightarrow P NP$

$NP \rightarrow P NP$

~~$PP \rightarrow P$~~

$NP \rightarrow P$

$NP \rightarrow \textit{primati}$

$NP \rightarrow \textit{kape}$

$NP \rightarrow \textit{nose}$

$NP \rightarrow \textit{glavi}$

$V \rightarrow \textit{primati}$

$S \rightarrow \textit{primati}$

$VP \rightarrow \textit{primati}$

$V \rightarrow \textit{kape}$

$S \rightarrow \textit{kape}$

$VP \rightarrow \textit{kape}$


$V \rightarrow \textit{nose}$

$S \rightarrow \textit{nose}$

$VP \rightarrow \textit{nose}$

$P \rightarrow \textit{na}$

$P \rightarrow \textit{na}$
 $PP \rightarrow \textit{na}$



Chomsky – eliminacija unarnih pravila

$S \rightarrow NP VP$

$VP \rightarrow V NP$

$S \rightarrow V NP$

$VP \rightarrow V NP PP$

$S \rightarrow V NP PP$

$VP \rightarrow V PP$

$S \rightarrow V PP$

$NP \rightarrow NP NP$

$NP \rightarrow NP PP$

$PP \rightarrow P NP$

$NP \rightarrow P NP$

~~$NP \rightarrow P$~~

$NP \rightarrow \textit{primati}$

$NP \rightarrow \textit{kape}$

$NP \rightarrow \textit{nose}$

$NP \rightarrow \textit{glavi}$

$V \rightarrow \textit{primati}$

$S \rightarrow \textit{primati}$

$VP \rightarrow \textit{primati}$

$V \rightarrow \textit{kape}$

$S \rightarrow \textit{kape}$

$VP \rightarrow \textit{kape}$

$V \rightarrow \textit{nose}$

$S \rightarrow \textit{nose}$

$VP \rightarrow \textit{nose}$

$P \rightarrow \textit{na}$

$PP \rightarrow \textit{na}$

$NP \rightarrow \textit{na}$

Chomsky – binarizacija

$S \rightarrow NP VP$

$VP \rightarrow V NP$

$S \rightarrow V NP$

~~$VP \rightarrow V NP PP$~~

$VP \rightarrow V @VP-V$
 $@VP-V \rightarrow NP PP$

~~$S \rightarrow V NP PP$~~

$S \rightarrow V @S-V$
 $@S-V \rightarrow NP PP$

$VP \rightarrow V PP$

$S \rightarrow V PP$

$NP \rightarrow NP NP$

$NP \rightarrow NP PP$

$PP \rightarrow P NP$

$NP \rightarrow P NP$

$NP \rightarrow \textit{primati}$

$NP \rightarrow \textit{kape}$

$NP \rightarrow \textit{nose}$

$NP \rightarrow \textit{glavi}$

$V \rightarrow \textit{primati}$

$S \rightarrow \textit{primati}$

$VP \rightarrow \textit{primati}$

$V \rightarrow \textit{kape}$

$S \rightarrow \textit{kape}$

$VP \rightarrow \textit{kape}$

$V \rightarrow \textit{nose}$

$S \rightarrow \textit{nose}$

$VP \rightarrow \textit{nose}$

$P \rightarrow \textit{na}$

$PP \rightarrow \textit{na}$

$NP \rightarrow \textit{na}$

Chomsky normalna forma

$S \rightarrow NP VP$

$VP \rightarrow V NP$

$S \rightarrow V NP$

$VP \rightarrow V @VP-V$

$@VP-V \rightarrow NP PP$

$S \rightarrow V @S-V$

$@S-V \rightarrow NP PP$

$VP \rightarrow V PP$

$S \rightarrow V PP$

$NP \rightarrow NP NP$

$NP \rightarrow NP PP$

$PP \rightarrow P NP$

$NP \rightarrow P NP$

$NP \rightarrow \textit{primati}$

$NP \rightarrow \textit{kape}$

$NP \rightarrow \textit{nose}$

$NP \rightarrow \textit{glavi}$

$V \rightarrow \textit{primati}$

$S \rightarrow \textit{primati}$

$VP \rightarrow \textit{primati}$

$V \rightarrow \textit{kape}$

$S \rightarrow \textit{kape}$

$VP \rightarrow \textit{kape}$

$V \rightarrow \textit{nose}$

$S \rightarrow \textit{nose}$

$VP \rightarrow \textit{nose}$

$P \rightarrow \textit{na}$

$PP \rightarrow \textit{na}$

$NP \rightarrow \textit{na}$

CFG i Chomsky normalna forma

$S \rightarrow NP VP$	$N \rightarrow \textit{primati}$
$VP \rightarrow V NP$	$N \rightarrow \textit{kape}$
$VP \rightarrow V NP PP$	$N \rightarrow \textit{nose}$
$NP \rightarrow NP NP$	$N \rightarrow \textit{glavi}$
$NP \rightarrow NP PP$	$V \rightarrow \textit{primati}$
$NP \rightarrow N$	$V \rightarrow \textit{kape}$
$NP \rightarrow \varepsilon$	$V \rightarrow \textit{nose}$
$PP \rightarrow P NP$	$P \rightarrow \textit{na}$

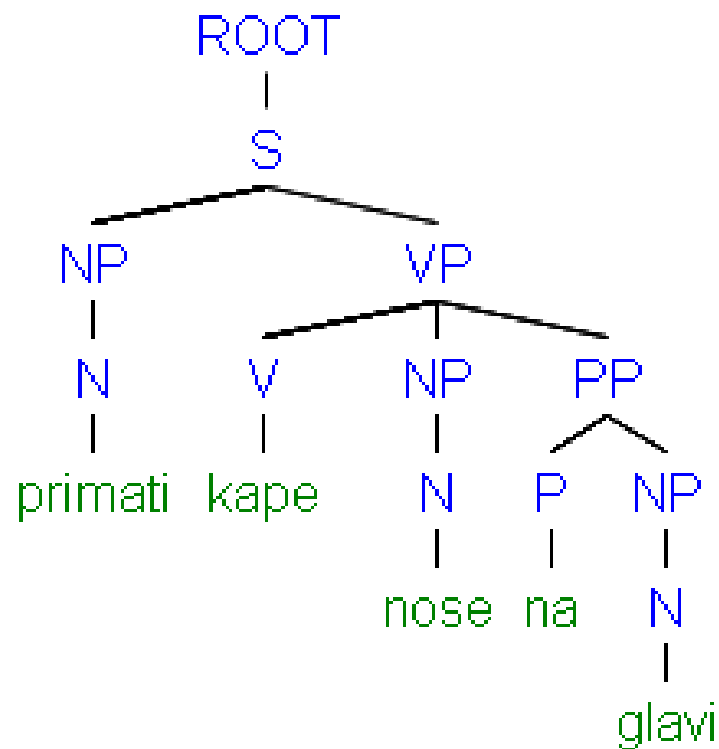
$S \rightarrow NP VP$	$NP \rightarrow \textit{primati}$
$VP \rightarrow V NP$	$NP \rightarrow \textit{kape}$
$S \rightarrow V NP$	$NP \rightarrow \textit{nose}$
$VP \rightarrow V @VP-V$	$NP \rightarrow \textit{glavi}$
$@VP-V \rightarrow NP PP$	$V \rightarrow \textit{primati}$
$S \rightarrow V @S-V$	$S \rightarrow \textit{primati}$
$@S-V \rightarrow NP PP$	$VP \rightarrow \textit{primati}$
$VP \rightarrow V PP$	$V \rightarrow \textit{kape}$
$S \rightarrow V PP$	$S \rightarrow \textit{kape}$
$NP \rightarrow NP NP$	$VP \rightarrow \textit{kape}$
$NP \rightarrow NP PP$	$V \rightarrow \textit{nose}$
$PP \rightarrow P NP$	$S \rightarrow \textit{nose}$
$NP \rightarrow P NP$	$VP \rightarrow \textit{nose}$
	$P \rightarrow \textit{na}$
	$PP \rightarrow \textit{na}$
	$NP \rightarrow \textit{na}$

Chomsky normalna forma

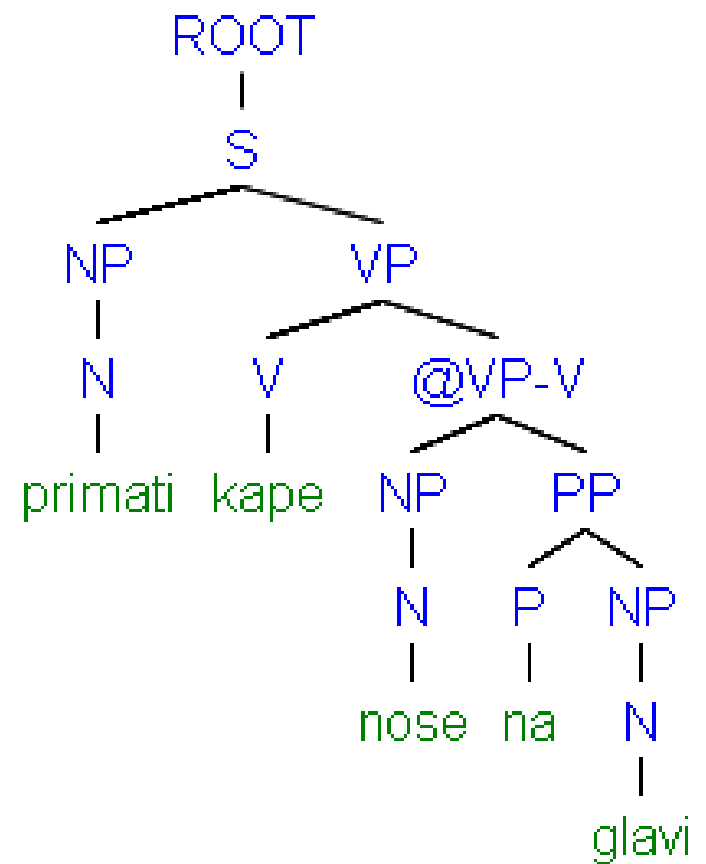
- Transformacija radi efikasnog parsiranja
- S pažljivim izborom neterminala moguće je rekonstruirati ista stabla detransformacijom
- Chomsky normalizacija nije lagana
 - rekonstrukcija n-arnih pravila je laka
 - rekonstrukcija unarnih i praznih pravila je složenija
- **Binarizacija** je ključna za $O(n^3)$ parsiranje CFG

Primjer: prije i poslije binarizacije

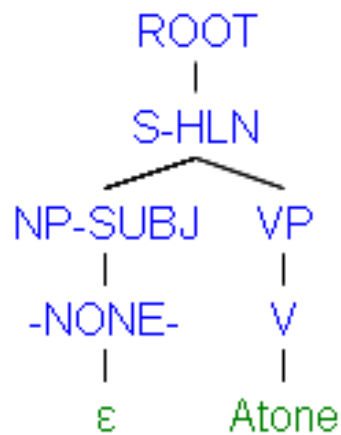
CFG



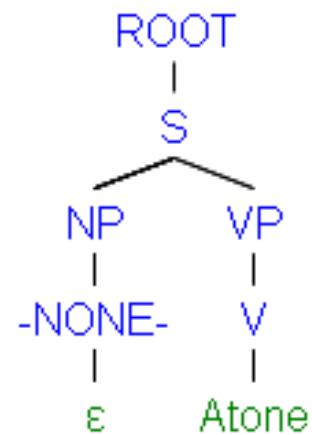
Normalizirani CFG



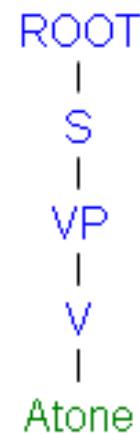
Banka stabala: unarna i prazna pravila



PTB Stablo



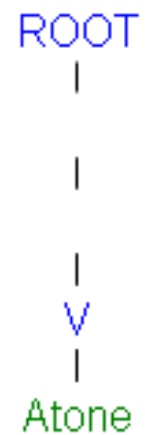
Bez
funkcijskih
oznaka



Bez
praznih
pravila



Visoko



Nisko

Bez
unarnih
pravila

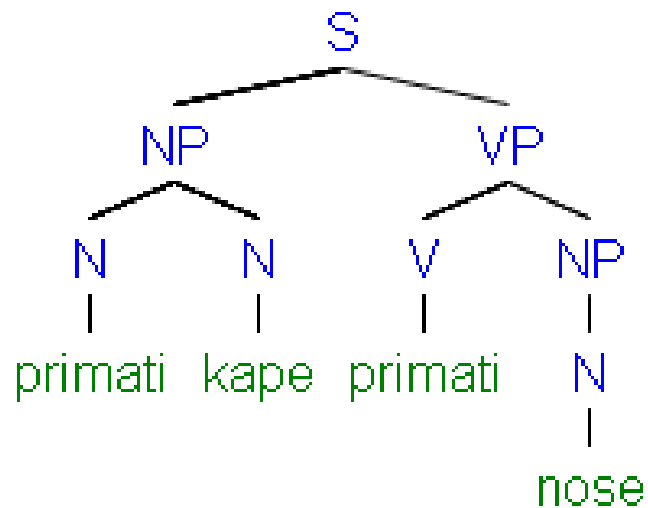
Uvod u obradu prirodnog jezika

14.3. CKY parsiranje

Branko Žitko

prevedeno od: Dan Jurafsky, Chris Manning

Strukturno parsiranje



PCFG

Vjerojatnost pravila θ_i

$S \rightarrow NP VP \quad \theta_1$

$VP \rightarrow V NP \quad \theta_2$

...

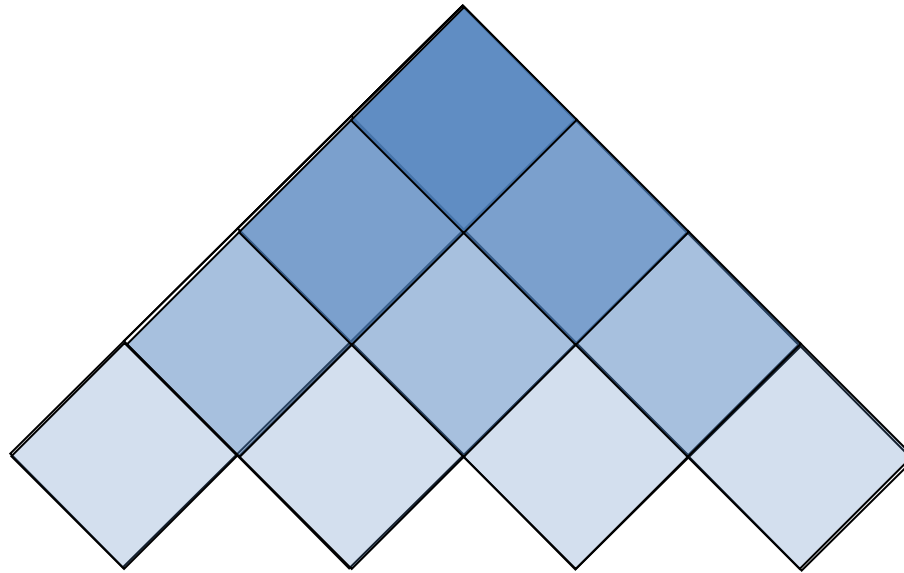
$N \rightarrow \textit{primati} \quad \theta_{42}$

$N \rightarrow \textit{kape} \quad \theta_{43}$

$V \rightarrow \textit{primati} \quad \theta_{44}$

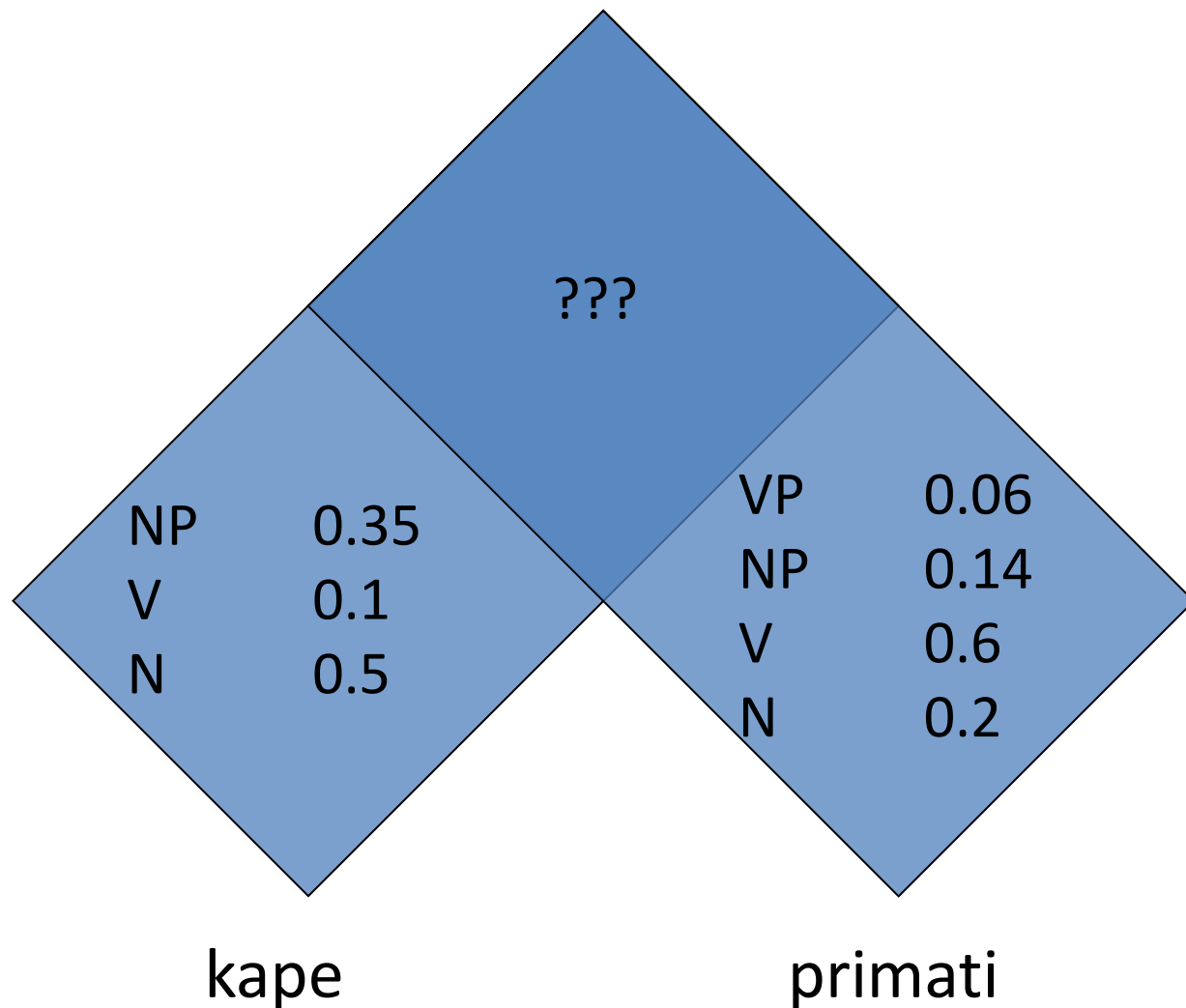
...

Cocke-Kasami-Younger (CKY) parsiranje



kape primati kape nose

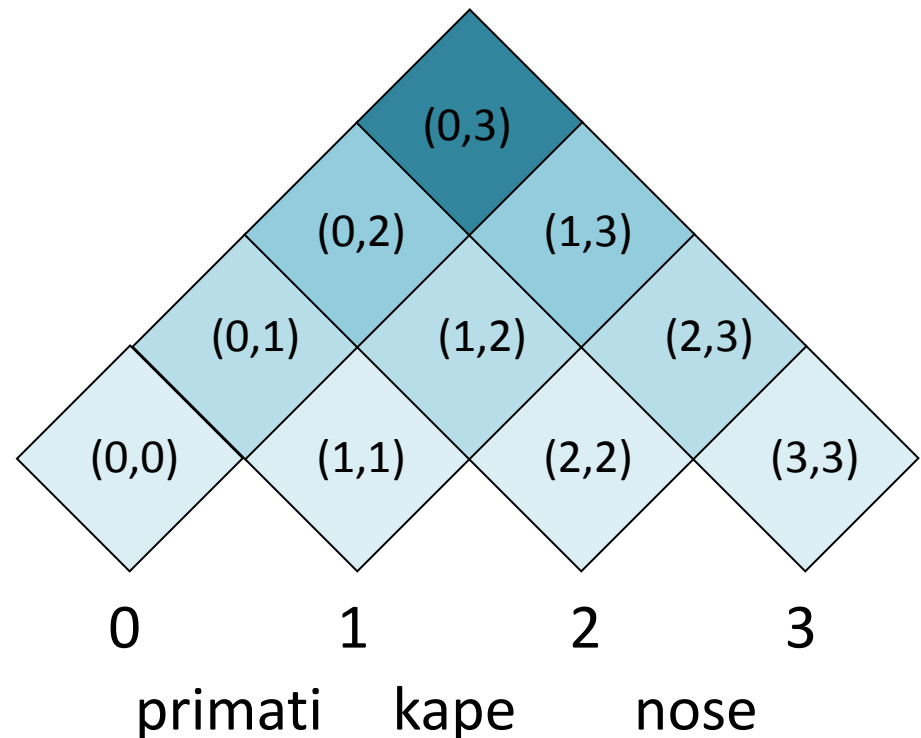
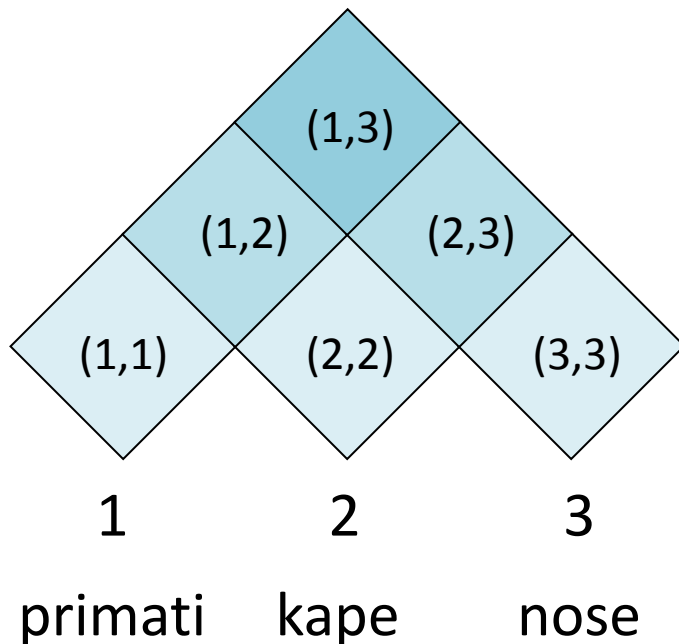
Viterbi (maksimalno bodovanje)



$S \rightarrow NP VP$	0.9
$S \rightarrow VP$	0.1
$VP \rightarrow V NP$	0.5
$VP \rightarrow V$	0.1
$VP \rightarrow V @VP_V$	0.3
$VP \rightarrow V PP$	0.1
$@VP_V \rightarrow NP PP$	1.0
$NP \rightarrow NP NP$	0.1
$NP \rightarrow NP PP$	0.2
$NP \rightarrow N$	0.7
$PP \rightarrow P NP$	1.0

Prošireno CKY parsiranje

- Unarna pravila se mogu uključiti u algoritam
 - malo neredno, ali ne povećava složenost algoritma
- Prazna pravila se mogu uključiti u algoritam
 - korištenje praznih ćelija
 - ne povećava složenost algoritma; slično kao unarna pravila



Prošireno CKY parsiranje

- Binarizacija je vitalna
 - bez binarizacije se ne dobiva kubično vrijeme parsiranja u odnosu na duljinu rečenice i broja neterminala u gramatici
 - Binarizacija može biti eksplicitna ili implicitna kod rada algoritma (kao Earley-ev algoritam), ali je uvijek prisutna

CKY algoritam

```
function CKY(rijeci, gramatika)
    # inicijalizacija
    bod = realna matrica dimenzije |rijeci|+1 x |rijeci|+1 x |neterminali|
    nazad = matrica parova dimenzije |rijeci|+1 x |rijeci|+1 x |neterminali|

    # prvi red
    for pocetak = 0 to |rijeci|-1 do
        kraj = pocetak + 1
        for A -> rijeci[pocetak] in gramatika do
            bod[pocetak][kraj][A] = P(A -> rijeci[pocetak])

    # unarna pravila za prvi red
    dodan = True
    while dodan do
        dodan = False
        for A -> B in gramatika do
            if bod[pocetak][kraj][B] > 0 then
                prob = P(A->B) * bod[pocetak][kraj][B]
                if prob > bod[pocetak][kraj][A] then
                    bod[pocetak][kraj][A] = prob
                    nazad[pocetak][kraj][A] = B
                    dodan = True
```

CKY algoritam

```
# ostali redovi
for raspon = 2 to |rijeci| do
    for pocetak = 0 to |rijeci|-raspon do
        kraj = pocetak + raspon
        for podjela = pocetak+1 to kraj-1 do
            for A -> B C in gramatika do
                prob = bod[pocetak][podjela][B] * bod[podjela][kraj][C] *
                    P(A -> B C)
                if prob > bod[pocetak][kraj][A] then
                    bod[pocetak][kraj][A] = prob
                    nazad[pocetak][kraj][A] = (podjela, B, C)

# unarna pravila za ostale redove
dodan = True
while dodan do
    dodan = False
    for A -> B in gramatika do
        prob = P(A -> B) * bod[pocetak][kraj][B]
        if prob > bod[pocetak][kraj][A] then
            bod[pocetak][kraj][A] = prob
            nazad[pocetak][kraj][A] = B
            dodan = True

return NapraviStablo(bod, nazad)
```

Uvod u obradu prirodnog jezika

14.4. CKY parsiranje: primjer

Branko Žitko

prevedeno od: Dan Jurafsky, Chris Manning

Binarna gramatika bez praznih pravila

$S \rightarrow NP VP$	0.9	$N \rightarrow \textit{primati}$	0.5
$S \rightarrow VP$	0.1	$N \rightarrow \textit{kape}$	0.2
$VP \rightarrow V NP$	0.5	$N \rightarrow \textit{nose}$	0.2
$VP \rightarrow V$	0.1	$N \rightarrow \textit{glavi}$	0.1
$VP \rightarrow V @VP_V$	0.3	$V \rightarrow \textit{primati}$	0.1
$VP \rightarrow V PP$	0.1	$V \rightarrow \textit{kape}$	0.6
$@VP_V \rightarrow NP PP$	1.0	$V \rightarrow \textit{nose}$	0.3
$NP \rightarrow NP NP$	0.1	$P \rightarrow \textit{na}$	1.0
$NP \rightarrow NP PP$	0.2		
$NP \rightarrow N$	0.7		
$PP \rightarrow P NP$	1.0		

CKY

S \rightarrow NP VP	0.9
S \rightarrow VP	0.1
VP \rightarrow V NP	0.5
VP \rightarrow V	0.1
VP \rightarrow V @VP_V	0.3
VP \rightarrow V PP	0.1
@VP_V \rightarrow NP PP	1.0
NP \rightarrow NP NP	0.1
NP \rightarrow NP PP	0.2
NP \rightarrow N	0.7
PP \rightarrow P NP	1.0
N \rightarrow <i>primati</i>	0.5
N \rightarrow <i>kape</i>	0.2
N \rightarrow <i>nose</i>	0.2
N \rightarrow <i>glavi</i>	0.1
V \rightarrow <i>primati</i>	0.1
V \rightarrow <i>kape</i>	0.6
V \rightarrow <i>nose</i>	0.3
P \rightarrow <i>na</i>	1.0

0	kape	1	primati	2	kape	3	nose	4
	bod[0][1]		bod[0][2]		bod[0][3]		bod[0][4]	
1			bod[1][2]		bod[1][3]		bod[1][4]	
2					bod[2][3]		bod[2][4]	
3							bod[3][4]	
4								

CKY – leksička pravila

$S \rightarrow NP VP$	0.9
$S \rightarrow VP$	0.1
$VP \rightarrow V NP$	0.5
$VP \rightarrow V$	0.1
$VP \rightarrow V @VP_V$	0.3
$VP \rightarrow V PP$	0.1
$@VP_V \rightarrow NP PP$	1.0
$NP \rightarrow NP NP$	0.1
$NP \rightarrow NP PP$	0.2
$NP \rightarrow N$	0.7
$PP \rightarrow P NP$	1.0
$N \rightarrow \textit{primati}$	0.5
$N \rightarrow \textit{kape}$	0.2
$N \rightarrow \textit{nose}$	0.2
$N \rightarrow \textit{glavi}$	0.1
$V \rightarrow \textit{primati}$	0.1
$V \rightarrow \textit{kape}$	0.6
$V \rightarrow \textit{nose}$	0.3
$P \rightarrow \textit{na}$	1.0

	0	1	2	3	4
0	kape				
1	$N \rightarrow \textit{kape} 0.2$ $V \rightarrow \textit{kape} 0.6$				
2		$N \rightarrow \textit{primati} 0.5$ $V \rightarrow \textit{primati} 0.1$			
3			$N \rightarrow \textit{kape} 0.2$ $V \rightarrow \textit{kape} 0.6$		
				$N \rightarrow \textit{nose} 0.1$ $V \rightarrow \textit{nose} 0.3$	

prvi red
for pocetak = 0 to |rijeci|-1 do
 kraj = pocetak + 1
 for A -> rijeci[pocetak] do
 bod[pocetak][kraj][A] = P(A -> rijeci[pocetak])

CKY – unarna pravila

$S \rightarrow NP VP$	0.9
$S \rightarrow VP$	0.1
$VP \rightarrow V NP$	0.5
$VP \rightarrow V$	0.1
$VP \rightarrow V @VP_V$	0.3
$VP \rightarrow V PP$	0.1
$@VP_V \rightarrow NP PP$	1.0
$NP \rightarrow NP NP$	0.1
$NP \rightarrow NP PP$	0.2
$NP \rightarrow N$	0.7
$PP \rightarrow P NP$	1.0
$N \rightarrow \textit{primati}$	0.5
$N \rightarrow \textit{kape}$	0.2
$N \rightarrow \textit{nose}$	0.2
$N \rightarrow \textit{glavi}$	0.1
$V \rightarrow \textit{primati}$	0.1
$V \rightarrow \textit{kape}$	0.6
$V \rightarrow \textit{nose}$	0.3
$P \rightarrow \textit{na}$	1.0

	0	1	2	3	4
kape					
primati					
kape					
nose					
	0	1	2	3	4
0	$N \rightarrow \textit{kape}$ 0.2 $V \rightarrow \textit{kape}$ 0.6 $NP \rightarrow N$ 0.14 $VP \rightarrow V$ 0.06 $S \rightarrow VP$ 0.006				
1		$N \rightarrow \textit{primati}$ 0.5 $V \rightarrow \textit{primati}$ 0.1 $NP \rightarrow N$ 0.35 $VP \rightarrow V$ 0.01 $S \rightarrow VP$ 0.001			
2			$N \rightarrow \textit{kape}$ 0.2 $V \rightarrow \textit{kape}$ 0.6 $NP \rightarrow N$ 0.14 $VP \rightarrow V$ 0.06 $S \rightarrow VP$ 0.006		
				$N \rightarrow \textit{nose}$ 0.1 $V \rightarrow \textit{nose}$ 0.3 $NP \rightarrow N$ 0.14 $VP \rightarrow V$ 0.03 $S \rightarrow VP$ 0.003	

```
# unarna pravila za prvi red
dodan = True
while dodan do
  dodan = False
  for A -> B in gramatika do
    if bod[pocetak][kraj][B] > 0 then
      prob = P(A->B) * bod[pocetak][kraj][B]
      if prob > bod[pocetak][kraj][A] then
        bod[pocetak][kraj][A] = prob
        nazad[pocetak][kraj][A] = B
        dodan = True
```

CKY – binarna pravila

		0	1	2	3	4
S → NP VP	0.9					
S → VP	0.1					
VP → V NP	0.5					
VP → V	0.1					
VP → V @VP_V	0.3					
VP → V PP	0.1					
@VP_V → NP PP	1.0					
NP → NP NP	0.1					
NP → NP PP	0.2					
NP → N	0.7					
PP → P NP	1.0					
N → <i>primati</i>	0.5					
N → <i>kape</i>	0.2					
N → <i>nose</i>	0.2					
N → <i>glavi</i>	0.1					
V → <i>primati</i>	0.1					
V → <i>kape</i>	0.6					
V → <i>nose</i>	0.3					
P → <i>na</i>	1.0					

CKY – unarna pravila

$S \rightarrow NP VP$	0.9
$S \rightarrow VP$	0.1
$VP \rightarrow V NP$	0.5
$VP \rightarrow V$	0.1
$VP \rightarrow V @VP_V$	0.3
$VP \rightarrow V PP$	0.1
$@VP_V \rightarrow NP PP$	1.0
$NP \rightarrow NP NP$	0.1
$NP \rightarrow NP PP$	0.2
$NP \rightarrow N$	0.7
$PP \rightarrow P NP$	1.0
$N \rightarrow \textit{primati}$	0.5
$N \rightarrow \textit{kape}$	0.2
$N \rightarrow \textit{nose}$	0.2
$N \rightarrow \textit{glavi}$	0.1
$V \rightarrow \textit{primati}$	0.1
$V \rightarrow \textit{kape}$	0.6
$V \rightarrow \textit{nose}$	0.3
$P \rightarrow \textit{na}$	1.0

	0	1	2	3	4
kape					
primati					
kape					
nose					
0	$N \rightarrow \textit{kape}$ 0.2 $V \rightarrow \textit{kape}$ 0.6 $NP \rightarrow N$ 0.14 $VP \rightarrow V$ 0.06 $S \rightarrow VP$ 0.006	$NP \rightarrow NP NP$ 0.0049 $VP \rightarrow V NP$ 0.105 $S \rightarrow VP$ 0.0105			
1		$N \rightarrow \textit{primati}$ 0.5 $V \rightarrow \textit{primati}$ 0.1 $NP \rightarrow N$ 0.35 $VP \rightarrow V$ 0.01 $S \rightarrow VP$ 0.001	$NP \rightarrow NP NP$ 0.0049 $VP \rightarrow V NP$ 0.007 $S \rightarrow NP VP$ 0.0189		
2			$N \rightarrow \textit{kape}$ 0.2 $V \rightarrow \textit{kape}$ 0.6 $NP \rightarrow N$ 0.14 $VP \rightarrow V$ 0.06 $S \rightarrow VP$ 0.006	$NP \rightarrow NP NP$ 0.00196 $VP \rightarrow V NP$ 0.042 $S \rightarrow VP$ 0.0042	
				$N \rightarrow \textit{nose}$ 0.1 $V \rightarrow \textit{nose}$ 0.3 $NP \rightarrow N$ 0.14 $VP \rightarrow V$ 0.03 $S \rightarrow VP$ 0.003	

```
# unarna pravila za ostale redove
dodan = True
while dodan do
  dodan = False
  for A -> B in gramatika do
    prob = P(A -> B) * bod[pocetak][kraj][B]
    if prob > bod[pocetak][kraj][A] then
      bod[pocetak][kraj][A] = prob
      nazad[pocetak][kraj][A] = B
  dodan = True
```

CKY – binarna i unarna pravila

$S \rightarrow NP VP$	0.9
$S \rightarrow VP$	0.1
$VP \rightarrow V NP$	0.5
$VP \rightarrow V$	0.1
$VP \rightarrow V @VP_V$	0.3
$VP \rightarrow V PP$	0.1
$@VP_V \rightarrow NP PP$	1.0
$NP \rightarrow NP NP$	0.1
$NP \rightarrow NP PP$	0.2
$NP \rightarrow N$	0.7
$PP \rightarrow P NP$	1.0
$N \rightarrow \textit{primati}$	0.5
$N \rightarrow \textit{kape}$	0.2
$N \rightarrow \textit{nose}$	0.2
$N \rightarrow \textit{glavi}$	0.1
$V \rightarrow \textit{primati}$	0.1
$V \rightarrow \textit{kape}$	0.6
$V \rightarrow \textit{nose}$	0.3
$P \rightarrow \textit{na}$	1.0

	0	1	2	3	4
kape					
primati					
kape					
nose					
0					
1	$N \rightarrow \textit{kape}$ 0.2 $V \rightarrow \textit{kape}$ 0.6 $NP \rightarrow N$ 0.14 $VP \rightarrow V$ 0.06 $S \rightarrow VP$ 0.006	$NP \rightarrow NP NP$ 0.0049 $VP \rightarrow V NP$ 0.105 $S \rightarrow VP$ 0.0105	$NP \rightarrow NP NP$ 0.0000686 $VP \rightarrow V NP$ 0.00147 $S \rightarrow NP VP$ 0.000882		
2		$N \rightarrow \textit{primati}$ 0.5 $V \rightarrow \textit{primati}$ 0.1 $NP \rightarrow N$ 0.35 $VP \rightarrow V$ 0.01 $S \rightarrow VP$ 0.001	$NP \rightarrow NP NP$ 0.0049 $VP \rightarrow V NP$ 0.007 $S \rightarrow NP VP$ 0.0189		
			$N \rightarrow \textit{kape}$ 0.2 $V \rightarrow \textit{kape}$ 0.6 $NP \rightarrow N$ 0.14 0.06 0.006	$NP \rightarrow NP NP$ 0.00196 $VP \rightarrow V NP$ 0.042 $S \rightarrow VP$ 0.0042	
				$N \rightarrow \textit{nose}$ 0.1 $V \rightarrow \textit{nose}$ 0.3 $NP \rightarrow N$ 0.14 $VP \rightarrow V$ 0.03 $S \rightarrow VP$ 0.003	

```
# ostali redovi
for raspon = 2 to |rijeci| do
  for pocetak = 0 to |rijeci| - raspon do
    kraj = pocetak + raspon
    for podjela = pocetak + 1 to kraj - 1 do
      for A -> B C in gramatika do
        prob = bod[pocetak][podjela][B] * bod[podjela][kraj][C]
        * P(A -> B C)
        if prob > bod[pocetak][kraj][A] then
          bod[pocetak][kraj][A] = prob
          nazad[pocetak][kraj][A] = (podjela, B, C)
```

CKY – binarna i unarna pravila

$S \rightarrow NP VP$	0.9
$S \rightarrow VP$	0.1
$VP \rightarrow V NP$	0.5
$VP \rightarrow V$	0.1
$VP \rightarrow V @VP_V$	0.3
$VP \rightarrow V PP$	0.1
$@VP_V \rightarrow NP PP$	1.0
$NP \rightarrow NP NP$	0.1
$NP \rightarrow NP PP$	0.2
$NP \rightarrow N$	0.7
$PP \rightarrow P NP$	1.0
$N \rightarrow \textit{primati}$	0.5
$N \rightarrow \textit{kape}$	0.2
$N \rightarrow \textit{nose}$	0.2
$N \rightarrow \textit{glavi}$	0.1
$V \rightarrow \textit{primati}$	0.1
$V \rightarrow \textit{kape}$	0.6
$V \rightarrow \textit{nose}$	0.3
$P \rightarrow \textit{na}$	1.0

	0	1	2	3	4
kape					
primati					
kape					
nose					
	<div>0</div> <div>$N \rightarrow \textit{kape}$ 0.2 $V \rightarrow \textit{kape}$ 0.6 $NP \rightarrow N$ 0.14 $VP \rightarrow V$ 0.06 $S \rightarrow VP$ 0.006</div>	<div>1</div> <div>$NP \rightarrow NP NP$ 0.0049 $VP \rightarrow V NP$ 0.105 $S \rightarrow VP$ 0.0105</div>	<div>2</div> <div>$NP \rightarrow NP NP$ 0.0000686 $VP \rightarrow V NP$ 0.00147 $S \rightarrow NP VP$ 0.000882</div>		
	<div>1</div>	<div>2</div> <div>$N \rightarrow \textit{primati}$ 0.5 $V \rightarrow \textit{primati}$ 0.1 $NP \rightarrow N$ 0.35 $VP \rightarrow V$ 0.01 $S \rightarrow VP$ 0.001</div>	<div>3</div> <div>$NP \rightarrow NP NP$ 0.0049 $VP \rightarrow V NP$ 0.007 $S \rightarrow NP VP$ 0.0189</div>	<div>4</div> <div>$NP \rightarrow NP NP$ 0.0000686 $VP \rightarrow V NP$ 0.000098 $S \rightarrow NP VP$ 0.01323</div>	
			<div>3</div> <div>$N \rightarrow \textit{kape}$ 0.2 $V \rightarrow \textit{kape}$ 0.6 $NP \rightarrow N$ 0.14 0.06 0.006</div>	<div>4</div> <div>$NP \rightarrow NP NP$ 0.00196 $VP \rightarrow V NP$ 0.042 $S \rightarrow VP$ 0.0042</div>	
				<div>4</div> <div>$N \rightarrow \textit{nose}$ 0.1 $V \rightarrow \textit{nose}$ 0.3 $NP \rightarrow N$ 0.14 $VP \rightarrow V$ 0.03 $S \rightarrow VP$ 0.003</div>	

```
# ostali redovi
for raspon = 2 to |rijeci| do
  for pocetak = 0 to |rijeci| - raspon do
    kraj = pocetak + raspon
    for podjela = pocetak + 1 to kraj - 1 do
      for A -> B C in gramatika do
        prob = bod[pocetak][podjela][B] * bod[podjela][kraj][C]
              * P(A -> B C)
        if prob > bod[pocetak][kraj][A] then
          bod[pocetak][kraj][A] = prob
          nazad[pocetak][kraj][A] = (podjela, B, C)
```


CKY – vraćanje unatrag

$S \rightarrow NP VP$	0.9
$S \rightarrow VP$	0.1
$VP \rightarrow V NP$	0.5
$VP \rightarrow V$	0.1
$VP \rightarrow V @VP_V$	0.3
$VP \rightarrow V PP$	0.1
$@VP_V \rightarrow NP PP$	1.0
$NP \rightarrow NP NP$	0.1
$NP \rightarrow NP PP$	0.2
$NP \rightarrow N$	0.7
$PP \rightarrow P NP$	1.0
$N \rightarrow \textit{primati}$	0.5
$N \rightarrow \textit{kape}$	0.2
$N \rightarrow \textit{nose}$	0.2
$N \rightarrow \textit{glavi}$	0.1
$V \rightarrow \textit{primati}$	0.1
$V \rightarrow \textit{kape}$	0.6
$V \rightarrow \textit{nose}$	0.3
$P \rightarrow \textit{na}$	1.0

	kape	1	primati	2	kape	3	nose	4
0	$N \rightarrow \textit{kape}$ $V \rightarrow \textit{kape}$ $NP \rightarrow N$ $VP \rightarrow V$ $S \rightarrow VP$	$NP \rightarrow NP NP$ 1 $VP \rightarrow V NP$ $S \rightarrow VP$		$NP \rightarrow NP NP$ $VP \rightarrow V NP$ $S \rightarrow NP VP$		$NP \rightarrow NP NP$ $VP \rightarrow V NP$ $S \rightarrow NP VP$ 2		
1		$N \rightarrow \textit{primati}$ $V \rightarrow \textit{primati}$ $NP \rightarrow N$ $VP \rightarrow V$ $S \rightarrow VP$		$NP \rightarrow NP NP$ $VP \rightarrow V NP$ $S \rightarrow NP VP$		$NP \rightarrow NP NP$ $VP \rightarrow V NP$ $S \rightarrow NP VP$		
2				$N \rightarrow \textit{kape}$ $V \rightarrow \textit{kape}$ $NP \rightarrow N$ $VP \rightarrow V$ $S \rightarrow VP$		$NP \rightarrow NP NP$ $VP \rightarrow V NP$ $S \rightarrow VP$ 3		
3						$N \rightarrow \textit{nose}$ $V \rightarrow \textit{nose}$ $NP \rightarrow N$ $VP \rightarrow V$ $S \rightarrow VP$		
4								

Uvod u obradu prirodnog jezika

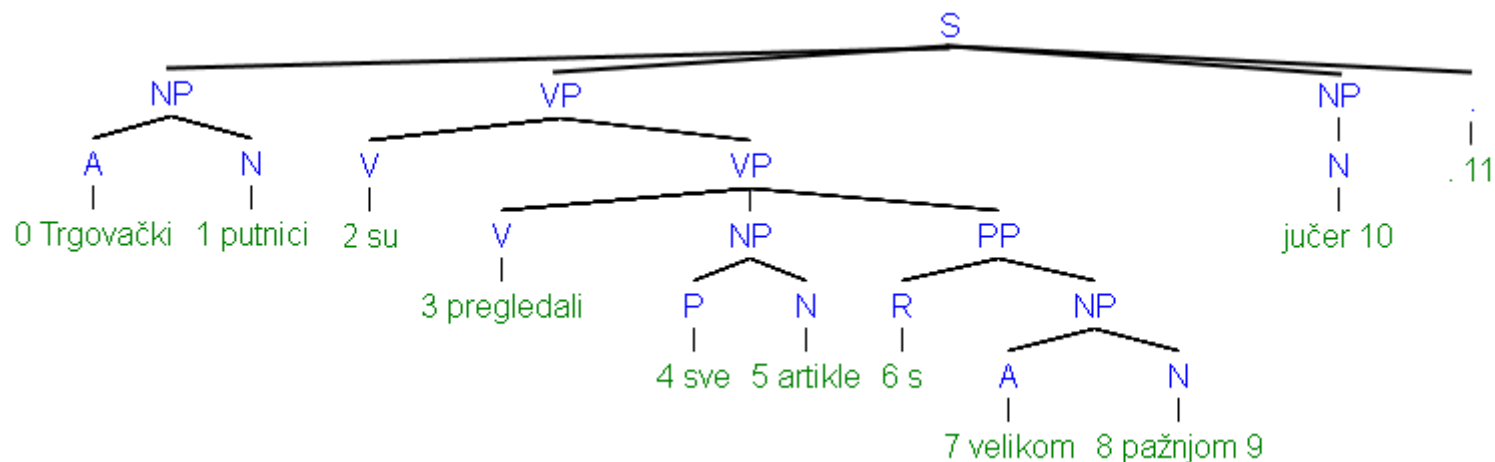
14.5. Evaluacija strukturnog parsiranja

Branko Žitko

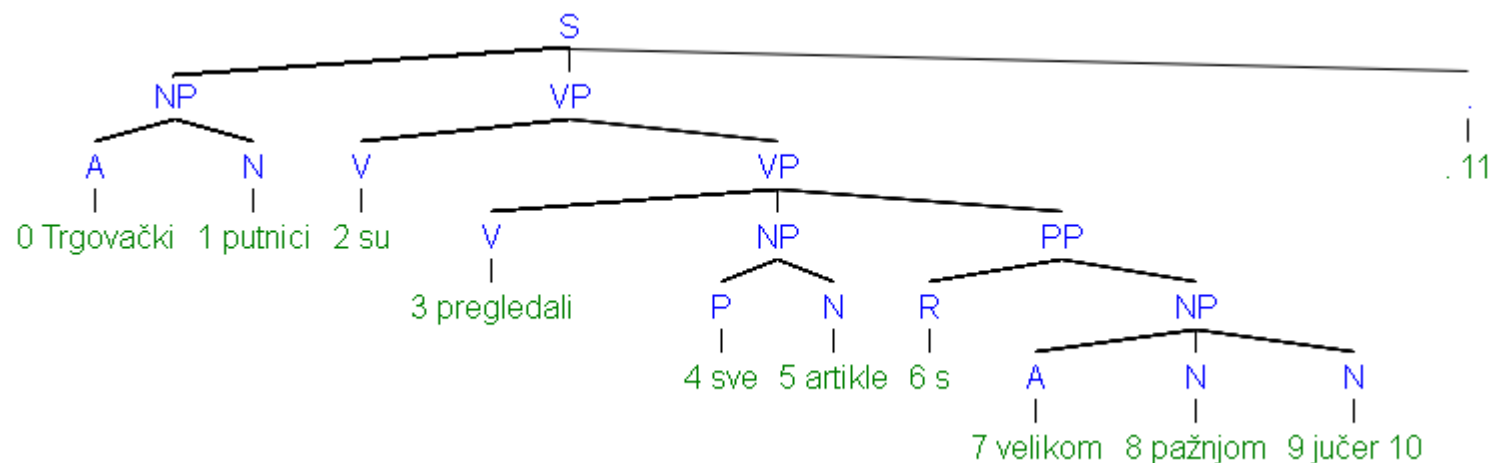
prevedeno od: Dan Jurafsky, Chris Manning

Evaluacija strukturnog parsiranja

Zlatni standard: **S(0:11)** NP(0:2) VP(2:9) VP(3:9) **NP(4:6)** PP(6:9) NP(7:9) NP(9:10)



Kandidat: **S(0:11)** NP(0:2) VP(2:10) VP(3:10) **NP(4:6)** PP(6:10) NP(7:10)



Evaluacija strukturnog parsiranja

Zlatni standard:

S(0:11) NP(0:2) VP(2:9) VP(3:9) NP(4:6) PP(6:9) NP(7:9) NP(9:10)

Kandidat:

S(0:11) NP(0:2) VP(2:10) VP(3:10) NP(4:6) PP(6:10) NP(7:10)

Preciznost oznake (PO): $3/7 = 42.9\%$

Odziv oznake (OO): $3/8 = 37.5\%$

F1 oznake: 40%

POS točnost: $11/11 = 100\%$

Koliko su dobre PCFG?

- Točnost parsiranja Penn WSJ: oko 73% F1
- Robusno
 - Obično prihvaća sve, ali s malom vjerojatnošću
- Parcijalno rješenje za višeznačnost gramatike
 - PCFG daje neke ideje vjerojatnosti parsiranja
 - ali ne toliko dobre jer su pretpostavke nezavisnosti previše jake
- Daju probabilistički model jezika
 - ali kod jednostavnih slučajeva daje lošije rezultate od trigram modela
- Izgleda da je problem PCFG-a u nedostatku leksikalizacije trigram modela