

content

- Viterbi for PCFG
- Inside-Outside-Algorithmus
- Berkeley-Parser

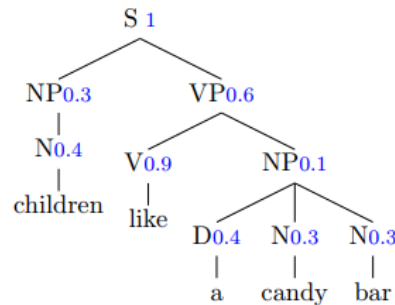
Viterbi for PCFG

What does the Viterbi algorithm compute and where and when do we use it?

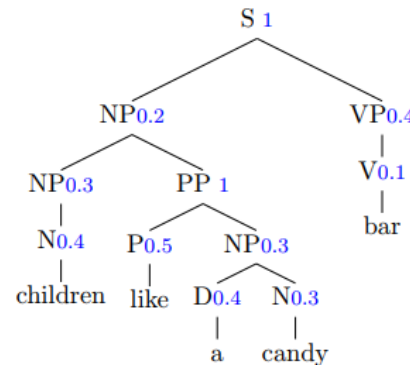
- we use it to find the best analysis (best parse tree) and the probability of it
- when the probability of each grammar rule is already given/computed.

Beispiele für Parsebaum-Wahrscheinlichkeiten

S	→ NP VP	1
NP	→ NP PP	0.2
NP	→ D N	0.3
NP	→ N	0.3
NP	→ D N N	0.1
NP	→ N N	0.1
VP	→ V NP	0.6
VP	→ V	0.4
PP	→ P NP	1
D	→ the	0.6
D	→ a	0.4
N	→ children	0.4
N	→ candy	0.3
N	→ bar	0.3
V	→ bar	0.1
V	→ like	0.9
P	→ like	0.5
P	→ for	0.5



$$p_1 = 0.0002333$$



$$p_2 = 0.0000173$$

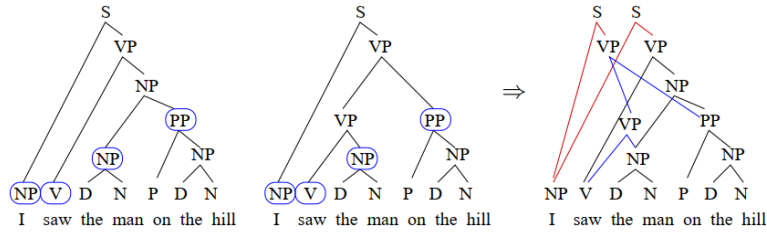
recap

- The initial problem that we want to solve is to find the best syntactic analysis of a sentence that can be parsed in more than one way by a symbolic syntactic parser (symbolische Parser)
- a symbolic parser (in this lecture) is a parser that uses some context-free grammar (without probability) to compute syntactic analyses of a sentence.
- In the example (left), we already have a PCFG (grammar rules with probabilities), so we can compute the probability of each tree for the sentence "children like a candy bar". The best tree is the one with the highest probability.
- we have learned how to estimate the prob for each rule (using Baumbank training or EM-training) and also how to compute the prob of a tree (by multiplying the prob of all nodes in the tree together)
- another way to determine the best tree is to use the **Viterbi algorithm**. To do that, we first have to merge both trees together to create a **Parsewald**.

Parsewald

Beispiel

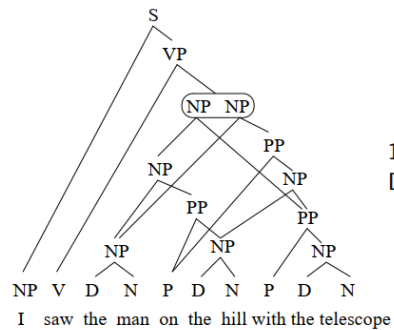
Zusammenfassen (maximaler) gemeinsamer Teilbäume



Die blauen Knoten der beiden Parsebäume werden jeweils zu einem Knoten zusammengefasst.

Parsewälder als spezialisierte Grammatiken

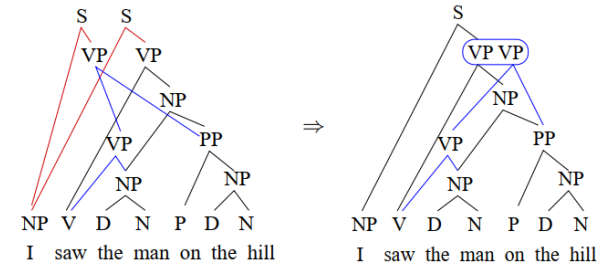
Wir definieren nun für jeden Parsewald eine Grammatik, die genau die darin enthaltenen Parsebäume generiert.



1. Schritt:
Durchnumerierung der Knoten

Beispiel

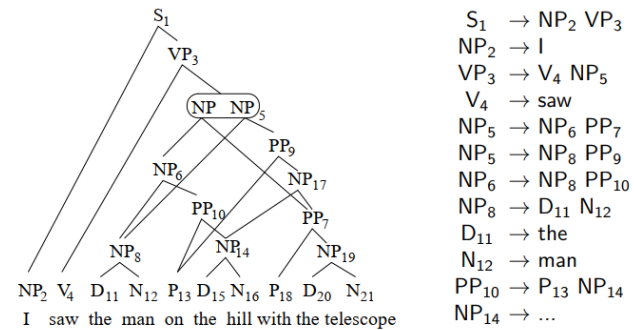
Zusammenfassen von Parsebäumen, die sich nur in einem Teilbaum unterscheiden



Hier werden die beiden roten Teilbäume zusammengefasst.
Die beiden VP-Teilbäume werden in einem ambigen VP-Knoten vereinigt.
Die blauen Kanten gehören zum 2. Parsebaum.

Parsewald als spezialisierte Grammatik

2. Schritt: Extraktion der Grammatikregeln



Diese Grammatik generiert genau die Parsebäume, die zu dem Parsewald zusammengefasst wurden (wenn man ignoriert, dass die Symbole noch Indizes tragen).

See the lecture video (Vorlesung8) to see how to use Viterbi

- Der **Viterbi**-Algorithmus berechnet die Wahrscheinlichkeit der **besten Analyse** jedes Parsewaldknotens

Berechnung der Viterbi-Wahrscheinlichkeiten

Viterbi-Algorithmus

$$\delta(a) = 1 \quad \text{für jedes Terminalsymbol } a$$

$$\delta(A \rightarrow X_1 \dots X_n) = p(A \rightarrow X_1 \dots X_n) \prod_{i=1}^n \delta(X_i) \quad \text{für Parsewald-Regeln}$$

$$\delta(A) = \max_{\alpha} \delta(A \rightarrow \alpha) \quad \text{für Nichtterminale } A$$

$$\psi(A) = \arg \max_{\alpha} \delta(A \rightarrow \alpha) \quad \text{beste Analyse von } A$$

$p(S_1 \rightarrow NP_2 VP_3)$ ist die Wahrscheinlichkeit $p(S \rightarrow NP VP)$ der entsprechenden PCFG-Regel.

Die ψ -Variable speichert die beste Parsewald-Regel für jedes Nichtterminal.

Viterbi-Beispiel

$$\delta(I) = 1$$

$$\delta(saw) = 1$$

...

$$\delta(NP_2) = \delta(NP_2 \rightarrow I)$$

$$= p(NP \rightarrow I) \delta(I)$$

$$\delta(V_4) = p(V \rightarrow saw) \delta(saw)$$

...

$$\delta(NP_8) = p(NP \rightarrow D N) \delta(D_{11}) \delta(N_{12})$$

...

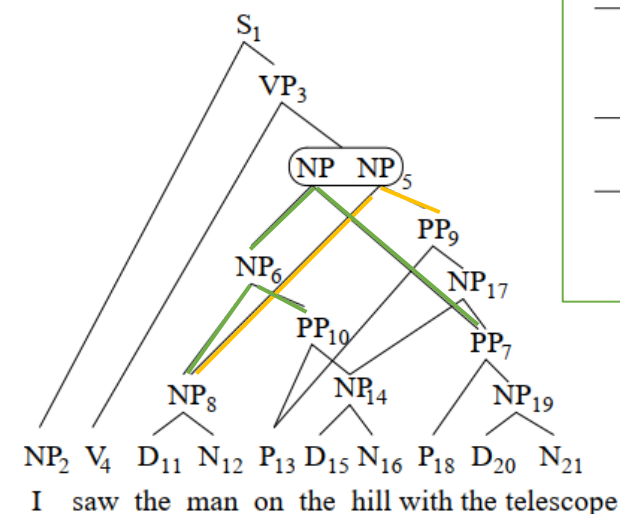
$$\delta(NP_6) = p(NP \rightarrow NP PP) \delta(NP_8) \delta(PP_{10})$$

...

$$\delta(NP_5) = \max(p(NP \rightarrow NP PP) \delta(NP_6) \delta(PP_7), p(NP \rightarrow NP PP) \delta(NP_8) \delta(PP_9))$$

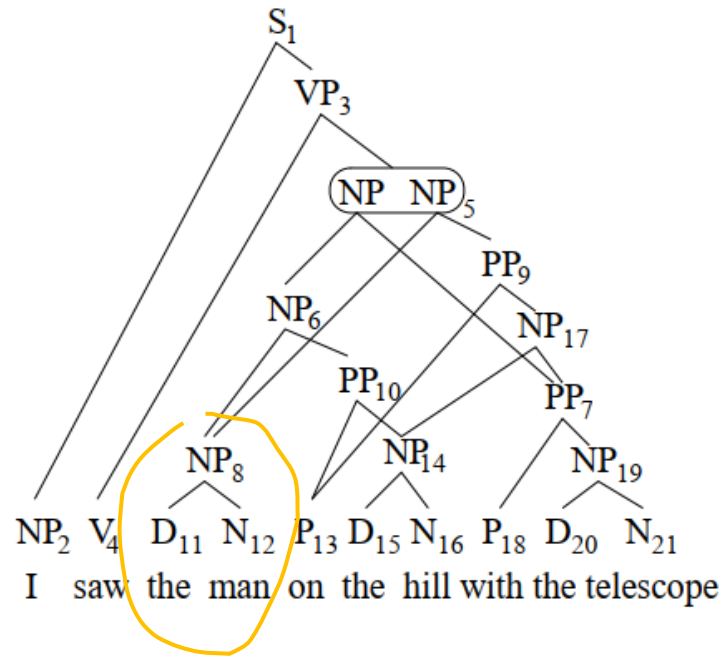
...

$$\delta(S_1) = p(S \rightarrow NP VP) \delta(NP_2) \delta(VP_3)$$



S	→ NP VP	1
NP	→ NP PP	0.2
NP	→ D N	0.3
NP	→ N	0.3
NP	→ D N N	0.1
NP	→ N N	0.1
VP	→ V NP	0.6
VP	→ V	0.4
PP	→ P NP	1
D	→ the	0.6
D	→ a	0.4
N	→ children	0.4
N	→ candy	0.3
N	→ bar	0.3
V	→ bar	0.1
V	→ like	0.9
P	→ like	0.5
P	→ for	0.5

Exercise1 : compute vit(NP8)



grammar for this
Parsewald

$S_1 \rightarrow NP_2 VP_3$
 $NP_2 \rightarrow I$
 $VP_3 \rightarrow V_4 NP_5$
 $V_4 \rightarrow \text{saw}$
 $NP_5 \rightarrow NP_6 PP_7$
 $NP_5 \rightarrow NP_8 PP_9$
 $NP_6 \rightarrow NP_8 PP_{10}$
 $NP_8 \rightarrow D_{11} N_{12}$
 $D_{11} \rightarrow \text{the}$
 $N_{12} \rightarrow \text{man}$
 $PP_{10} \rightarrow P_{13} NP_{14}$
 $NP_{14} \rightarrow \dots$

explain

1. compute vit(the), vit(man) using (1)

vit(the) = 1

vit(man) = 1

2. then vit(D11), vit(N12) using (2)

for D11, look at the grammar for Parsewald and find every rule that starts with D11, here we found only one (D11 → the), so we have

vit(D11) = max [vit (D11 → the)]

= vit (D11 → the) # then use (3)

= p(D → the) vit(the)

look at the global grammar for p(D → the)

= 0.6 * 1

= 0.6

Viterbi-Algorithmus

(1) $\delta(a) = 1$ für jedes Terminalsymbol a

(3) $\delta(A \rightarrow X_1 \dots X_n) = p(A \rightarrow X_1 \dots X_n) \prod_{i=1}^n \delta(X_i)$ für Parsewald-Regeln

(2) $\delta(A) = \max_{\alpha} \delta(A \rightarrow \alpha)$ für Nichtterminale A

$\psi(A) = \arg \max_{\alpha} \delta(A \rightarrow \alpha)$ beste Analyse von A

global grammar

S	→ NP VP	1
NP	→ NP PP	0.2
NP	→ D N	0.3
NP	→ N	0.3
NP	→ D N N	0.1
NP	→ N N	0.1
VP	→ V NP	0.6
VP	→ V	0.4
PP	→ P NP	1
D	→ the	0.6
D	→ a	0.4
N	→ children	0.4
N	→ candv	0.3
N	→ man	0.3
V	→ bar	0.1
V	→ like	0.9
P	→ like	0.5
P	→ for	0.5

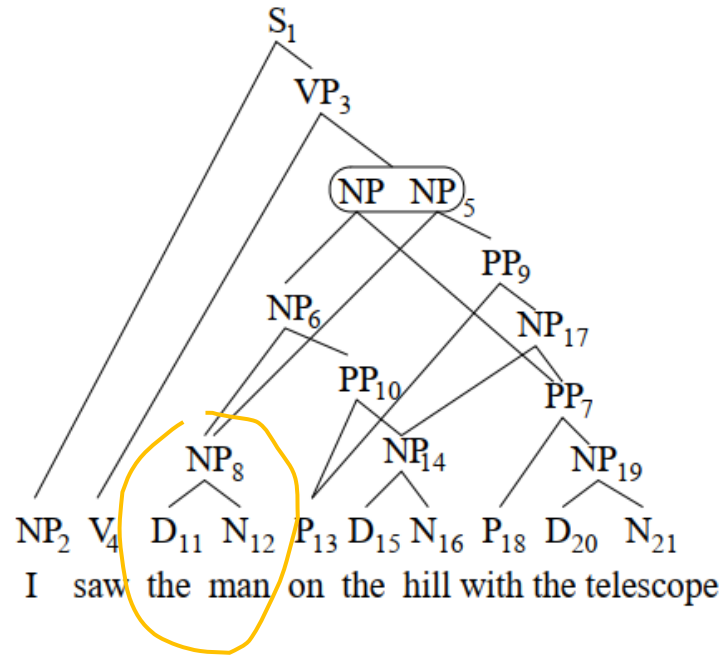
steps

1. compute vit(the), vit(man)

2. then vit(D11), vit(N12)

3. then vit(NP8)

Exercise1 : compute vit(NP8)



grammar for this
Parsewald

$S_1 \rightarrow NP_2 VP_3$
 $NP_2 \rightarrow I$
 $VP_3 \rightarrow V_4 NP_5$
 $V_4 \rightarrow \text{saw}$
 $NP_5 \rightarrow NP_6 PP_7$
 $NP_5 \rightarrow NP_8 PP_9$
 $NP_6 \rightarrow NP_8 PP_{10}$
 $NP_8 \rightarrow D_{11} N_{12}$
 $D_{11} \rightarrow \text{the}$
 $N_{12} \rightarrow \text{man}$
 $PP_{10} \rightarrow P_{13} NP_{14}$
 $NP_{14} \rightarrow \dots$

explain (continue)

2. then vit(D11), vit(N12) using (2)

$\text{vit}(N_{12}) = \max (\text{vit}(N_{12} \rightarrow \text{man}))$
 $= \text{vit}(N_{12} \rightarrow \text{man})$
 $= p(N \rightarrow \text{man}) \text{vit}(\text{man})$
 $= 0.3 * 1$
 $= 0.3$

steps

1. compute vit(the), vit(man)
2. then vit(D11), vit(N12)
3. then vit(NP8)

Viterbi-Algorithmus

(1) $\delta(a) = 1$ für jedes Terminalsymbol a

(3) $\delta(A \rightarrow X_1 \dots X_n) = p(A \rightarrow X_1 \dots X_n) \prod_{i=1}^n \delta(X_i)$ für Parsewald-Regeln

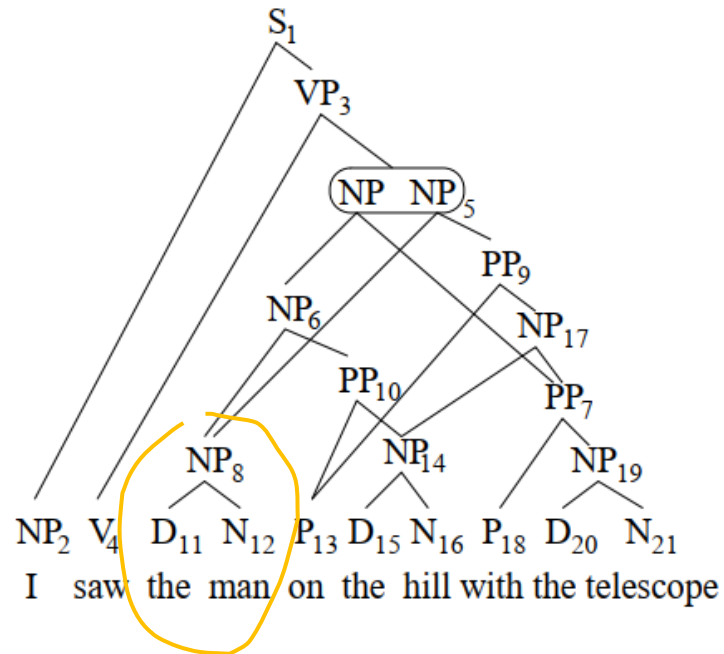
(2) $\delta(A) = \max_{\alpha} \delta(A \rightarrow \alpha)$ für Nichtterminale A

$\psi(A) = \arg \max_{\alpha} \delta(A \rightarrow \alpha)$ beste Analyse von A

global grammar

S	→ NP VP	1
NP	→ NP PP	0.2
NP	→ D N	0.3
NP	→ N	0.3
NP	→ D N N	0.1
NP	→ N N	0.1
VP	→ V NP	0.6
VP	→ V	0.4
PP	→ P NP	1
D	→ the	0.6
D	→ a	0.4
N	→ children	0.4
N	→ candv	0.3
N	→ man	0.3
V	→ bar	0.1
V	→ like	0.9
P	→ like	0.5
P	→ for	0.5

Exercise1 : compute vit(NP8)



grammar for this
Parsewald

$S_1 \rightarrow NP_2 VP_3$
 $NP_2 \rightarrow I$
 $VP_3 \rightarrow V_4 NP_5$
 $V_4 \rightarrow \text{saw}$
 $NP_5 \rightarrow NP_6 PP_7$
 $NP_5 \rightarrow NP_8 PP_9$
 $NP_6 \rightarrow NP_8 PP_{10}$
 $NP_8 \rightarrow D_{11} N_{12}$
 $D_{11} \rightarrow \text{the}$
 $N_{12} \rightarrow \text{man}$
 $PP_{10} \rightarrow P_{13} NP_{14}$
 $NP_{14} \rightarrow \dots$

explain (continue)

3. compute vit(NP8)

$\text{vit}(NP8) = \max [\text{vit}(NP8 \rightarrow D_{11} N_{12})]$
 $= \text{vit}(NP8 \rightarrow D_{11} N_{12})$
 $= p(NP \rightarrow D N) \text{vit}(D_{11}) \text{vit}(N_{12})$
 $= 0.3 * 0.6 * 0.3$

finished

steps

1. compute vit(the), vit(man)
2. then vit(D11), vit(N12)
3. then vit(NP8)

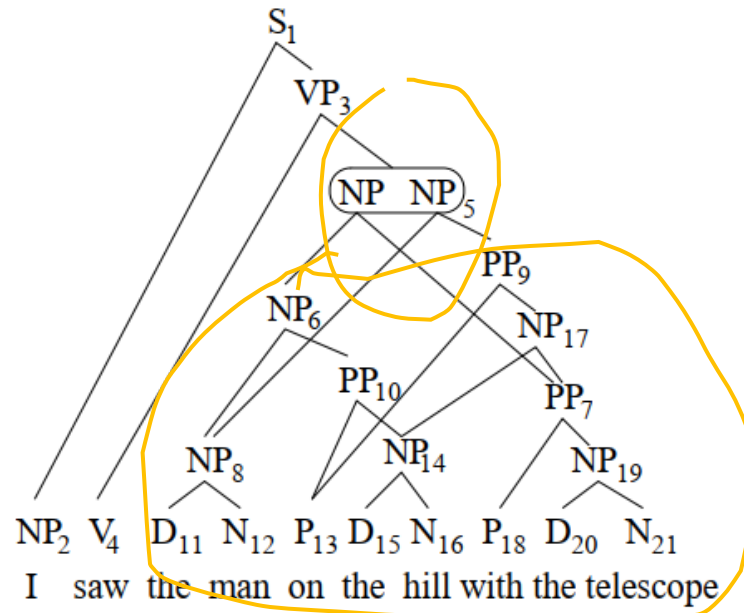
Viterbi-Algorithmus

- (1) $\delta(a) = 1$ für jedes Terminalsymbol a
- (3) $\delta(A \rightarrow X_1 \dots X_n) = p(A \rightarrow X_1 \dots X_n) \prod_{i=1}^n \delta(X_i)$ für Parsewald-Regeln
- (2) $\delta(A) = \max_{\alpha} \delta(A \rightarrow \alpha)$ für Nichtterminale A
- $\psi(A) = \arg \max_{\alpha} \delta(A \rightarrow \alpha)$ beste Analyse von A

global grammar

S	→ NP VP	1
NP	→ NP PP	0.2
NP	→ D N	0.3
NP	→ N	0.3
NP	→ D N N	0.1
NP	→ N N	0.1
VP	→ V NP	0.6
VP	→ V	0.4
PP	→ P NP	1
D	→ the	0.6
D	→ a	0.4
N	→ children	0.4
N	→ candv	0.3
N	→ man	0.3
V	→ bar	0.1
V	→ like	0.9
P	→ like	0.5
P	→ for	0.5

Exercise2 : compute vit(NP5)



$S_1 \rightarrow NP_2 VP_3$
 $NP_2 \rightarrow I$
 $VP_3 \rightarrow V_4 NP_5$
 $V_4 \rightarrow \text{saw}$
 $NP_5 \rightarrow NP_6 PP_9$
 $NP_6 \rightarrow NP_8 PP_{10}$
 $NP_8 \rightarrow D_{11} N_{12}$
 $D_{11} \rightarrow \text{the}$
 $N_{12} \rightarrow \text{man}$
 $PP_{10} \rightarrow P_{13} NP_{14}$
 $NP_{14} \rightarrow \dots$

Viterbi-Algorithmus

- (1) $\delta(a) = 1$ für jedes Terminalsymbol a
- (2) $\delta(A) = \max_{\alpha} \delta(A \rightarrow \alpha)$ für Nichtterminale A
- (3) $\delta(A \rightarrow X_1 \dots X_n) = p(A \rightarrow X_1 \dots X_n) \prod_{i=1}^n \delta(X_i)$ für Parsewald-Regeln
- $\psi(A) = \arg \max_{\alpha} \delta(A \rightarrow \alpha)$ beste Analyse von A

steps (in the actual computation)

- compute vit prob of every node below NP5
- then compute vit(NP5)

In the exam, you may only have to write down the computation of vit(NP5) and do not need to actually compute it for other nodes below NP5 (see the answer example)

Answer

$\text{vit}(NP_5) = \max [\text{vit}(NP_5 \rightarrow NP_6 PP_7), \text{vit}(NP_5 \rightarrow NP_8 PP_9)]$ (2)

where

$\text{vit}(NP_5 \rightarrow NP_6 PP_7) = p(NP \rightarrow NP PP) \text{vit}(NP_6) \text{vit}(PP_7)$ (3)

$\text{vit}(NP_5 \rightarrow NP_8 PP_9) = p(NP \rightarrow NP PP) \text{vit}(NP_8) \text{vit}(PP_9)$

finished

$S \rightarrow NP VP$	1
$NP \rightarrow NP PP$	0.2
$NP \rightarrow D N$	0.3
$NP \rightarrow N$	0.3
$NP \rightarrow D N N$	0.1
$NP \rightarrow N N$	0.1
$VP \rightarrow V NP$	0.6
$VP \rightarrow V$	0.4
$PP \rightarrow P NP$	1
$D \rightarrow \text{the}$	0.6
$D \rightarrow \text{a}$	0.4
$N \rightarrow \text{children}$	0.4
$N \rightarrow \text{candv}$	0.3
$N \rightarrow \text{man}$	0.3
$V \rightarrow \text{bar}$	0.1
$V \rightarrow \text{like}$	0.9
$P \rightarrow \text{like}$	0.5
$P \rightarrow \text{for}$	0.5

Inside Outside

Inside-Outside-Algorithmus

Der Inside-Outside-Algorithmus

Anwendung

- berechnet effizient die **erwarteten Regelhäufigkeiten** beim EM-Training
- und entspricht damit dem Forward-Backward-Algorithmus bei den HMMs.
- Er berechnet bottom-up **Inside**-Wahrscheinlichkeiten (ähnlich dem Viterbi-Algorithmus)
- und top-down **Outside**-Wahrscheinlichkeiten.
- Aus den Inside- und Outside-Wahrscheinlichkeiten berechnet er die **erwarteten Häufigkeiten**.

To use the Inside-Outside-Algorithm to compute "die erwartete Häufigkeit" $\gamma(A \rightarrow B)$, we first compute the Inside-WK for every node (1) and rule (2) of the given Parsewald and grammar rule with prob. Then compute the Outside-WK. Then we can compute $\gamma(A \rightarrow B)$ for every rule $A \rightarrow B$ in the Parsewald. Do it for all sentences in the training set, then from $\gamma(A \rightarrow B)$, we can compute "die erwartete Regelhäufigkeiten" $f(A \rightarrow B)$ that we can use in EM-Training.

Inside-Wahrscheinlichkeiten

- Der **Viterbi**-Algorithmus berechnet die Wahrscheinlichkeit der **besten Analyse** jedes Parsewaldknotens
- Der **Inside**-Algorithmus berechnet die Gesamtwahrscheinlichkeit **aller Analysen** für jeden Parsewaldknoten.

Inside-Algorithmus

$$\begin{aligned} (1) \quad \alpha(a) &= 1 \quad \text{für Terminalsymbol } a \\ (2) \quad \alpha(A \rightarrow X_1 \dots X_n) &= p(A \rightarrow X_1 \dots X_n) \prod_{i=1}^n \alpha(X_i) \quad \text{für Parsewaldregeln} \\ (1) \quad \alpha(A) &= \sum_{A \rightarrow \gamma} \alpha(A \rightarrow \gamma) \quad \text{für Nichtterminale } A \end{aligned}$$

Der Unterschied zum Viterbi-Algorithmus ist die Summen-Operation statt der max-Operation.

Old method (inefficient)

Training auf unannotierten Daten

Ansatz 3: (Forts.)

EM-Training

- 1 Initialisierung der Regelwahrscheinlichkeiten
- 2 für alle Sätze (E-Schritt)
 - Berechnung der Parsebäume für den Satz
 - Berechnung der Parsebaumgewichte $p(t|s)$
 - Extraktion der gewichteten Regelhäufigkeiten
- 3 Neuschätzung der Regelwahrscheinlichkeiten (M-Schritt)
- 4 Weiter mit Schritt 2 bis ein Endekriterium erfüllt ist

compute f

compute p

Regel	p_0	f_1	p_1	f_2	p_2	f_3	p_3
S → NP VP	1.00	2.00	1.00	2.00	1.00	2.00	1.00
NP → D N	0.25	0.50	0.11	0.10	0.02	0.00	0.00
NP → D N N	0.25	0.50	0.11	0.90	0.22	1.00	0.25
NP → N	0.25	3.00	0.67	3.00	0.73	3.00	0.75
NP → NP PP	0.25	0.50	0.11	0.10	0.02	0.00	0.00
VP → V	0.50	0.50	0.25	0.10	0.05	0.00	0.00
VP → V NP	0.50	1.50	0.75	1.90	0.95	2.00	1.00
PP → P NP	1.00	0.50	1.00	0.10	1.00	0.00	1.00
D → a	0.50	1.00	1.00	1.00	1.00	1.00	1.00
D → the	0.50	0.00	0.00	0.00	0.00	0.00	0.00
N → bar	0.25	0.50	0.11	0.90	0.18	1.00	0.20
N → candy	0.25	1.00	0.22	1.00	0.20	1.00	0.20
N → children	0.25	2.00	0.44	2.00	0.41	2.00	0.40
N → chocolate	0.25	1.00	0.22	1.00	0.20	1.00	0.20
V → bar	0.50	0.50	0.25	0.10	0.05	0.00	0.00
V → like	0.50	1.50	0.75	1.90	0.95	2.00	1.00
P → like	1.00	0.50	1.00	0.10	1.00	0.00	1.00
-logprob		15.2		11.3		9.3	

Old method (inefficient)

Training auf unannotierten Daten

Ansatz 3: (Forts.)

EM-Training

- ➊ Initialisierung der Regelwahrscheinlichkeiten
- ➋ für alle Sätze (E-Schritt)
 - ▶ Berechnung der Parsebäume für den Satz
 - ▶ Berechnung der Parsebaumgewichte $p(t|s)$
 - ▶ Extraktion der gewichteten Regelhäufigkeiten
- ➌ Neuschätzung der Regelwahrscheinlichkeiten (M-Schritt)
- ➍ Weiter mit Schritt 2 bis ein Endekriterium erfüllt ist

compute f

compute p

More efficient with Inside-outside-Algorithm

Neuschätzung der Regel-Wahrscheinlichkeiten

EM-Algorithmus: wiederholte Ausführung der beiden Schritte

- ➊ E-Schritt compute f
 - ▶ Jeder Satz des Trainingskorpus wird geparkt und ein Parsewald erstellt.
 - ▶ Die erwartete Häufigkeit $\gamma(A \rightarrow \delta)$ jeder Parsewaldregel $A \rightarrow \delta$ wird berechnet.
 - ▶ Die erwarteten Häufigkeiten $\gamma(A \rightarrow \delta)$ werden für jede CFG-Regel über alle Trainingssätze zu $f(A' \rightarrow \delta')$ summiert, wobei sich A' und δ' aus A und δ durch Entfernen der Knotennummern ergeben.
- ➋ M-Schritt
Die Regel-Wahrscheinlichkeiten werden aus den erwarteten Häufigkeiten neu geschätzt:

compute p

$$p(A \rightarrow \delta) = \frac{f(A \rightarrow \delta)}{\sum_{\delta'} f(A \rightarrow \delta')}$$

Neuschätzung der Regel-Wahrscheinlichkeiten

EM-Algorithmus: wiederholte Ausführung der beiden Schritte

① E-Schritt compute f + initialize p uniformly (1)

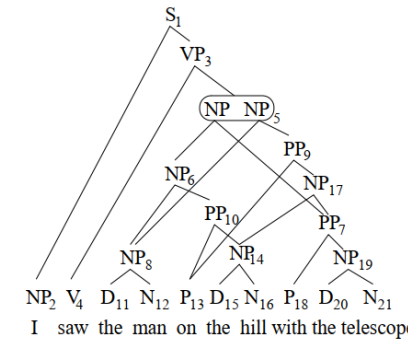
- ▶ Jeder Satz des Trainingskorpus wird geparkt und ein Parsewald erstellt.
- ▶ Die erwartete Häufigkeit $\gamma(A \rightarrow \delta)$ jeder Parsewaldregel $A \rightarrow \delta$ wird berechnet. (4)
- ▶ Die erwarteten Häufigkeiten $\gamma(A \rightarrow \delta)$ werden für jede CFG-Regel über alle Trainingssätze zu $f(A' \rightarrow \delta')$ summiert, wobei sich A' und δ' aus A und δ durch Entfernen der Knotennummern ergeben.

② M-Schritt
Die Regel-Wahrscheinlichkeiten werden aus den erwarteten Häufigkeiten neu geschätzt:

$$(3) \quad p(A \rightarrow \delta) = \frac{f(A \rightarrow \delta)}{\sum_{\delta'} f(A \rightarrow \delta')} \quad (2)$$

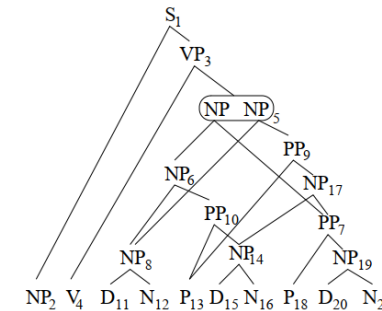
Suppose we have 2 sentences in our training set, want to compute $f(A \rightarrow \delta)$ (2) so that we can compute $p(A \rightarrow \delta)$ (3) and perform the EM-Training. Create a Parsewald for every sentence to merge all possible analyses of the sentence together. Then compute the Inside and Outside WK for all Parsewald. Using the Inside and Outside WK, we can compute **gamma**($A \rightarrow \delta$) (4) for every rule in each Parsewald. Then we sum up $\gamma(A \rightarrow \delta)$ from all Parsewald that has the same underlying rule together. The result is $f(A \rightarrow \delta)$ (2).
For example, to compute $f(\text{NP} \rightarrow \text{NP PP})$ we look in Parsewald1 for all rules of this form, sum the gamma value of them. Then look at the Parsewald2, do the same. Sum the results together.

$$f(\text{NP} \rightarrow \text{NP PP}) = \gamma(\text{NP}_5 \rightarrow \text{NP}_6 \text{ PP}_7) + \gamma(\text{NP}_5 \rightarrow \text{NP}_8 \text{ PP}_9) + \gamma(\text{NP}_6 \rightarrow \text{NP}_8 \text{ PP}_9) + \text{from Parsewald1} \\ \gamma(\text{NP}_8 \rightarrow \text{NP}_1 \text{ PP}_2) + \gamma(\text{NP}_6 \rightarrow \text{NP}_2 \text{ PP}_2) \quad \text{from Parsewald2}$$



sentence1

$S_1 \rightarrow \text{NP}_2 \text{ VP}_3$
 $\text{NP}_2 \rightarrow \text{I}$
 $\text{VP}_3 \rightarrow \text{V}_4 \text{ NP}_5$
 $\text{V}_4 \rightarrow \text{saw}$
 $\text{NP}_5 \rightarrow \text{NP}_6 \text{ PP}_7$ gamma= 0.3
 $\text{NP}_5 \rightarrow \text{NP}_8 \text{ PP}_9$ gamma= 0.5
 $\text{NP}_6 \rightarrow \text{NP}_8 \text{ PP}_{10}$ gamma= 0.5
 $\text{NP}_8 \rightarrow \text{D}_{11} \text{ N}_{12}$
 $\text{D}_{11} \rightarrow \text{the}$
 $\text{N}_{12} \rightarrow \text{man}$
 $\text{PP}_{10} \rightarrow \text{P}_{13} \text{ NP}_{14}$
 $\text{NP}_{14} \rightarrow \dots$



sentence2

$S_1 \rightarrow \text{NP}_2 \text{ VP}_3$
 $\text{NP}_2 \rightarrow \text{I}$
 $\text{VP}_3 \rightarrow \text{V}_4 \text{ NP}_5$
 $\text{NP}_8 \rightarrow \text{NP}_1 \text{ PP}_2$ gamma= 0.12
 $\text{NP}_6 \rightarrow \text{NP}_2 \text{ PP}_2$ gamma= 0.4
 $\text{NP}_8 \rightarrow \text{D}_{11} \text{ N}_{12}$
 $\text{D}_{11} \rightarrow \text{the}$
 $\text{N}_{12} \rightarrow \text{man}$
 $\text{PP}_{10} \rightarrow \text{P}_{13} \text{ NP}_{14}$
 $\text{NP}_{14} \rightarrow \dots$

(1) (2) (3)

Regel	p_0	f_1	p_1
$S \rightarrow \text{NP VP}$	1.00	2.00	1.00
$\text{NP} \rightarrow \text{D N}$	0.25	0.50	0.11
$\text{NP} \rightarrow \text{D N N}$	0.25	0.50	0.11
$\text{NP} \rightarrow \text{N}$	0.25	3.00	0.67
$\text{NP} \rightarrow \text{NP PP}$	0.25	0.50	0.11
$\text{VP} \rightarrow \text{V}$	0.50	0.50	0.25
$\text{VP} \rightarrow \text{V NP}$	0.50	1.50	0.75
$\text{PP} \rightarrow \text{P NP}$	1.00	0.50	1.00
$\text{D} \rightarrow \text{a}$	0.50	1.00	1.00
$\text{D} \rightarrow \text{the}$	0.50	0.00	0.00

Inside-Wahrscheinlichkeiten

- Der **Viterbi**-Algorithmus berechnet die Wahrscheinlichkeit der **besten Analyse** jedes Parsewaldknotens
- Der **Inside**-Algorithmus berechnet die Gesamtwahrscheinlichkeit **aller Analysen** für jeden Parsewaldknoten.

Inside-Algorithmus

$$\alpha(a) = 1 \quad \text{für Terminalsymbol } a$$

$$\alpha(A \rightarrow X_1 \dots X_n) = p(A \rightarrow X_1 \dots X_n) \prod_{i=1}^n \alpha(X_i) \quad \text{für Parsewaldregeln}$$

$$\alpha(A) = \sum_{A \rightarrow \gamma} \alpha(A \rightarrow \gamma) \quad \text{für Nichtterminale } A$$

Der Unterschied zum Viterbi-Algorithmus ist die Summen-Operation statt der max-Operation.

Berechnung der Outside-Wahrscheinlichkeiten

Outside-Algorithmus

$$\beta(S) = 1 \quad \text{für Startsymbol } S$$

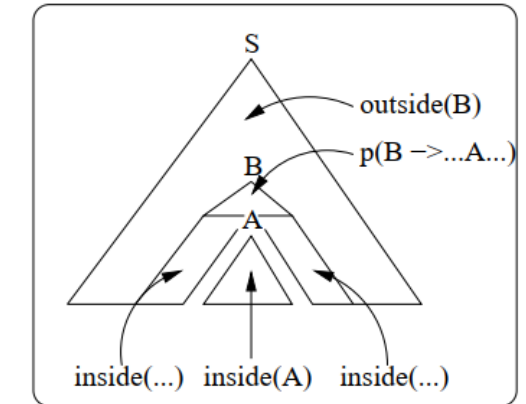
$$\beta(A) = \sum_{B \rightarrow \gamma A \delta} \beta(B \rightarrow \gamma A \delta)$$

$$\beta(B \rightarrow X_1 \dots X_m \underline{A} X_{m+1} \dots X_n) = \beta(B) p(B \rightarrow X_1 \dots X_m \underline{A} X_{m+1} \dots X_n) \prod_{i=1}^n \alpha(X_i)$$

$$\beta(B \rightarrow \gamma \underline{A} \delta) = \beta(B) \frac{\alpha(B \rightarrow \gamma A \delta)}{\alpha(A)} \quad \text{mit } \gamma = X_1 \dots X_m \text{ und } \delta = X_{m+1}, \dots, X_n$$

$\beta(B \rightarrow \gamma \underline{A} \delta)$: Beitrag der Regel $B \rightarrow \gamma A \delta$ zur Outside-Wahrscheinlichkeit von A

Beachte, dass $p(B \rightarrow X_1 \dots X_m \underline{A} X_{m+1} \dots X_n) \prod_{i=1}^n \alpha(X_i) = \alpha(B \rightarrow X_1 \dots X_m \underline{A} X_{m+1} \dots X_n) / \alpha(A)$

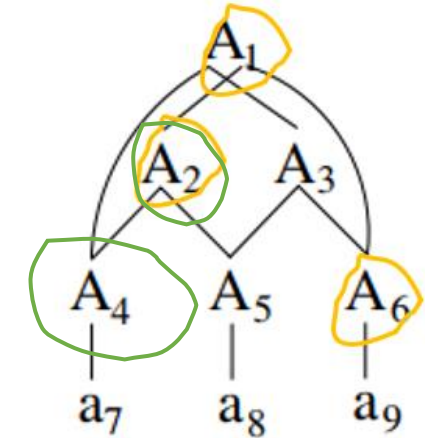


Outside-Algorithmus

$$b(A_2) = b(A_1 \rightarrow A_2 A_3) \\ = b(A_1) p(A \rightarrow A A) a(A_3)$$

$$b(A_4) = \text{Sum} (b(A_1 \rightarrow A_4 A_3), b(A_2 \rightarrow A_4 A_5))$$

where
 $b(A_1 \rightarrow A_4 A_3) = b(A_1) p(A \rightarrow A A) a(A_3)$ # method1
 $b(A_2 \rightarrow A_4 A_5) = b(A_2) p(A \rightarrow A A) a(A_5)$



$$\beta(S) = 1 \quad \text{für Startsymbol } S$$

$$\beta(A) = \sum_{B \rightarrow \gamma A \delta} \beta(B \rightarrow \gamma A \delta)$$

p

S → NP VP	1
NP → NP PP	0.2
NP → D N	0.3
NP → N	0.3
NP → D N N	0.1
NP → N N	0.1
VP → V NP	0.6
VP → V	0.4
PP → P NP	1
D → the	0.6
D → a	0.4
N → children	0.4
N → candy	0.3
N → bar	0.3
V → bar	0.1
V → like	0.9
P → like	0.5
P → for	0.5

Methode 1)

$$\beta(B \rightarrow \gamma A \delta) = \beta(B) \frac{\alpha(B \rightarrow \gamma A \delta)}{\alpha(A)} \quad \text{mit } \gamma = X_1 \dots X_m \text{ und } \delta = X_{m+1}, \dots, X_n$$

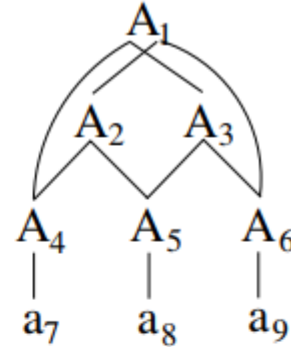
Methode 2)

$$\beta(B \rightarrow X_1 \dots X_m A X_{m+1} \dots X_n) = \beta(B) p(B \rightarrow X_1 \dots X_m A X_{m+1} \dots X_n) \prod_{i=1}^n \alpha(X_i)$$

inside WK

Inside-Outside-Beispiel

PCFG: $A \rightarrow A A$ 0.6
 $A \rightarrow a$ 0.4



$$\alpha(a_7) = 1 = \alpha(a_8) = \alpha(a_9)$$

$$\alpha(A_4) = p(A \rightarrow a) \alpha(a_7) = 0.4 * 1 = \mathbf{0.4} = \alpha(A_5) = \alpha(A_6)$$

$$\alpha(A_2) = p(A \rightarrow A A) \alpha(A_4) \alpha(A_5) = 0.6 * 0.4 * 0.4 = \mathbf{0.096}$$

$$\alpha(A_3) = \alpha(A_2)$$

$$\begin{aligned} \alpha(A_1) &= p(A \rightarrow A A) \alpha(A_4) \alpha(A_3) + p(A \rightarrow A A) \alpha(A_2) \alpha(A_6) \\ &= 0.6 * 0.4 * 0.096 * 2 = \mathbf{0.04608} \end{aligned}$$

$$\beta(A_1) = 1$$

$$\beta(A_2) = \beta(A_1) p(A \rightarrow A A) \alpha(A_6) = 1 * 0.6 * 0.4 = \mathbf{0.24} = \beta(A_3)$$

$$\begin{aligned} \beta(A_4) &= \beta(A_1) p(A \rightarrow A A) \alpha(A_3) + \beta(A_2) p(A \rightarrow A A) \alpha(A_5) \\ &= 1 * 0.6 * 0.096 + 0.24 * 0.6 * 0.4 = \mathbf{0.1152} = \beta(A_6) \end{aligned}$$

$$\begin{aligned} \beta(A_5) &= \beta(A_2) p(A \rightarrow A A) \alpha(A_4) + \beta(A_3) p(A \rightarrow A A) \alpha(A_6) \\ &= 0.24 * 0.6 * 0.4 * 2 = \mathbf{0.1152} \end{aligned}$$

$$\beta(a_7) = \beta(A_4) p(A \rightarrow a) = 0.1152 * 0.4 = \mathbf{0.04608}$$

$$\begin{aligned} \gamma(A_2 \rightarrow A_4 A_5) &= \beta(A_2) p(A \rightarrow A A) \alpha(A_4) \alpha(A_5) / \alpha(A_1) \\ &= 0.24 * 0.6 * 0.4 * 0.4 / 0.04608 = 0.5 \end{aligned}$$

- erwartete Regelhäufigkeit

$$\gamma(A \rightarrow \delta) = \alpha(A \rightarrow \delta) \beta(A) / \alpha(S)$$

Aufgabe 7) Erläutern Sie den EM-Algorithmus am Beispiel des unüberwachten Trainings von PCFGs (also Training auf Roh⁷texten). Welche Daten benötigen Sie? Welche Berechnungsschritte führt der EM-Algorithmus aus? (4 Punkte)

- ▶ Training auf automatisch geparsten Texten (von einem Parser erzeugte Parsewälder)

EM-Algorithmus: wiederholte Ausführung der beiden Schritte

① E-Schritt

- ▶ Jeder Satz des Trainingskorpus wird geparkt und ein Parsewald ausgegeben.
- ▶ Die erwartete Häufigkeit jeder Parsewaldregel wird berechnet.
- ▶ Die erwarteten Häufigkeiten werden für jede CFG-Regel über alle Trainingssätze summiert.

② M-Schritt

Die Regel-Wahrscheinlichkeiten werden aus den erwarteten Häufigkeiten neu geschätzt:

$$p(A \rightarrow \delta) = \frac{\gamma(A \rightarrow \delta)}{\sum_{\delta'} \gamma(A \rightarrow \delta')}$$

7

Parser-Evaluierung

Precision

Es werden korrekte **Konstituenten** statt korrekter **Parsebäume** gezählt

Eine Konstituente ist **korrekt**, wenn der Goldstandard-Parse eine Konstituente mit derselben Start- und Endposition und derselben Kategorie enthält.

$$\text{Precision} = \frac{\text{Anzahl korrekte Konstituenten}}{\text{Anzahl ausgegebene Konstituenten}}$$

Berechnung von Precision und Recall

TP (True Positives): Zahl der ausgegebenen Konstituenten, die korrekt sind

FP (False Positives): Zahl der ausgegebenen Konstituenten, die falsch sind

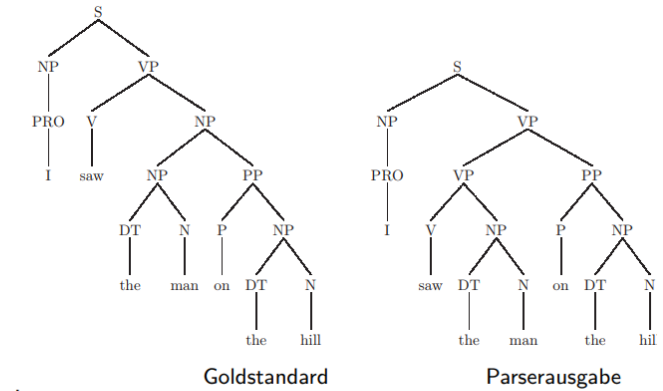
FN (False Negatives): Zahl der Goldstandard-Konstituenten, die fehlten

$$\text{Precision} = \frac{\text{TP}}{\text{TP} + \text{FP}} \quad \text{Recall} = \frac{\text{TP}}{\text{TP} + \text{FN}}$$

Mit dem **Recall** messen wir zusätzlich, wieviel Prozent der Goldstandard-Konstituenten im Parse enthalten waren.

Zusammen ergeben Precision und Recall ein genaues Bild der Parsing-Genauigkeit.

Beispiel



Goldstandard-Konstituenten:

(S,0,7) (NP,0,1) (VP,1,7) (NP,2,7) (NP,2,4) (PP,4,7) (NP,5,7)

Konstituenten in Parserausgabe:

(S,0,7) (NP,0,1) (VP,1,7) (VP,1,4) (NP,2,4) (PP,4,7) (NP,5,7)

True Positives = 6

False Positives = 1

False Negatives = 1

Precision = $\text{TP} / (\text{TP} + \text{FP}) = 6/7$

Recall = $\text{TP} / (\text{TP} + \text{FN}) = 6/7$

F-Score = $2 \text{ P R} / (\text{P} + \text{R}) = 2 * 6/7 * 6/7 / (6/7 + 6/7)$

F-Score

F-Score: harmonisches Mittel aus Precision und Recall

$$F_1 = \frac{1}{\frac{1}{P} + \frac{1}{R}} = \frac{2}{\frac{R}{PR} + \frac{P}{PR}} = \frac{2PR}{P + R}$$

Der gewichtete F-Score gibt Precision β -mal mehr Gewicht als Recall:

$$F_\beta = \frac{(1 + \beta^2)PR}{\beta^2 P + R}$$

See the lecture video for an explanation

Markowisierung

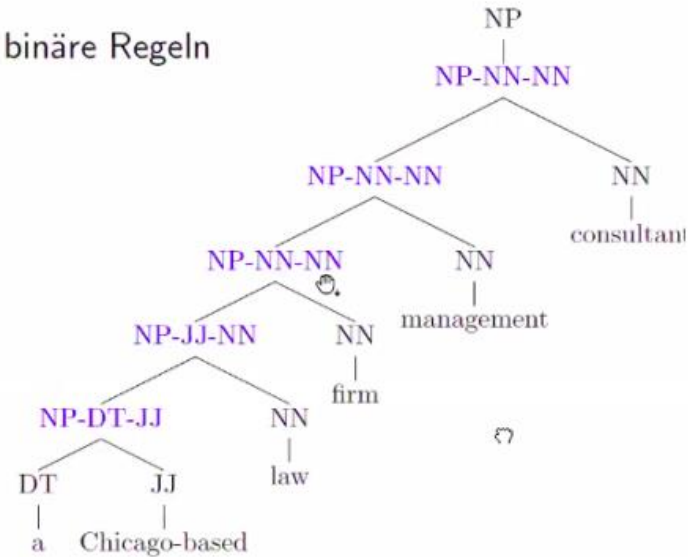
- Aufspaltung von langen Regeln in binäre Regeln
- Neue **Hilfskategorien** mit z.B.

- ▶ der Elternkategorie und
- ▶ den Kategorien der beiden letzten Tochterknoten

- Das Entfernen der Hilfsknoten liefert wieder den Originalparse

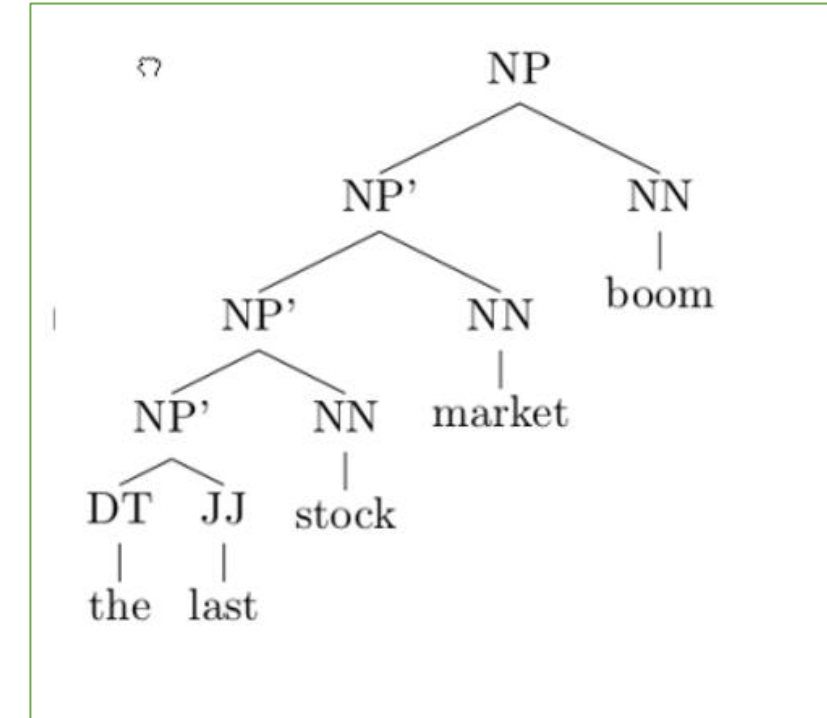
⇒ Reduktion der Grammatikgröße

⇒ Verbesserung ihrer Abdeckung



Die gezeigte Markowisierung ist äquivalent dazu, dass wir die rechten Seiten (und Wahrscheinlichkeiten) von bspw. NP-Regeln mit einem Markow-Modell 2. Ordnung erzeugen.

Binarisierung



See the lecture video for an explanation

Berkeley-Parser

Berkeley Parser : Grundidee

Wir wollen die Parse-Kategorien (die nicht-terminalen Symbole wie S, NP, PRO,...) in SubKategorien verfeinern.

Die Methode von Petrov und Klein soll die Subkategorisierung automatisch lernt.

Jede Kategorie wird in zwei neue Kategorien aufgespaltet. Z.B. NP wird NP/0 und NP/1.

Damit bekommen wir auch neue Regeln z.B. NP/0 → PRO/0 oder NP/0 → PRO/1 und viele Parsebäume aus der neuen Regeln. Die neuen Parsebäume sollen mit EM-Training trainiert werden (EM schätzt die WK der neuen Regeln).
Nach dem Training: Wenn wir die Regeln nach dem Kategorie auf der linken Seite gruppieren und die Regeln nach der Regel-WK sortieren, sehen wir die gelernte Subkategorisierung.

Beispiel: wahrscheinlichste Expansionen der DT-(Unter-)Kategorien

DT	the (0.50)	a (0.24)	The (0.08)	...
DT/0	that (0.15)	this (0.14)	some (0.11)	...
DT/1	the (0.54)	a (0.25)	The (0.09)	...

DT/0 → that (0.15)
DT/0 → this (0.14)
DT/0 → some (0.11)

DT/1 → the (0.54)
DT/1 → a (0.25)
DT/1 → The (0.09)

Der originale Parsebaum wird dann mit den neuen Subkategorien annotiert. Dadurch dass die Kategorien spezifischer werden, wird die Performanz der Parser auch erhöht.

Synthetische Merkmale

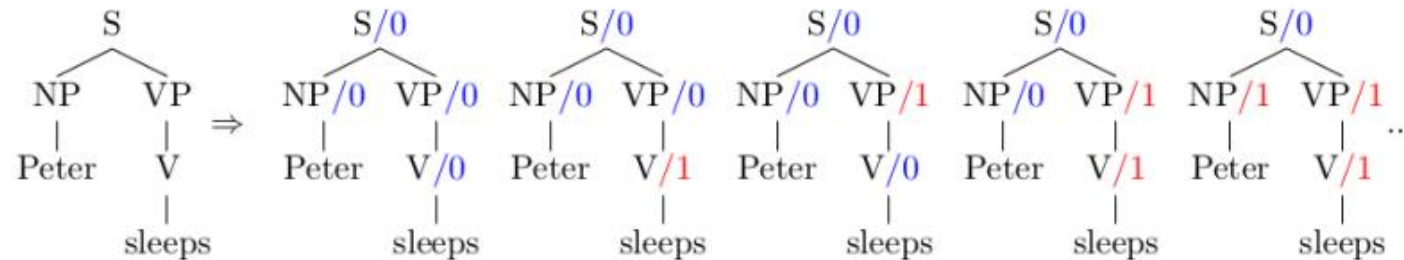
Grundidee (von Petrov/Klein)

- Alle Kategorien werden durch ein synthetisches Merkmal mit den Werten 0 bzw. 1 **aufgespalten**.
- Jeder Parse der Baumbank kann von der neuen Grammatik auf viele unterschiedliche Arten generiert werden.
- Durch **EM-Training** wird die neue Grammatik an die Baumbank angepasst.

...

$NP'/0 \rightarrow NP'/0 \text{ } NN/0$	0.24
$NP'/0 \rightarrow NP'/0 \text{ } NN/1$	0.26
$NP'/0 \rightarrow NP'/1 \text{ } NN/0$	0.25

...

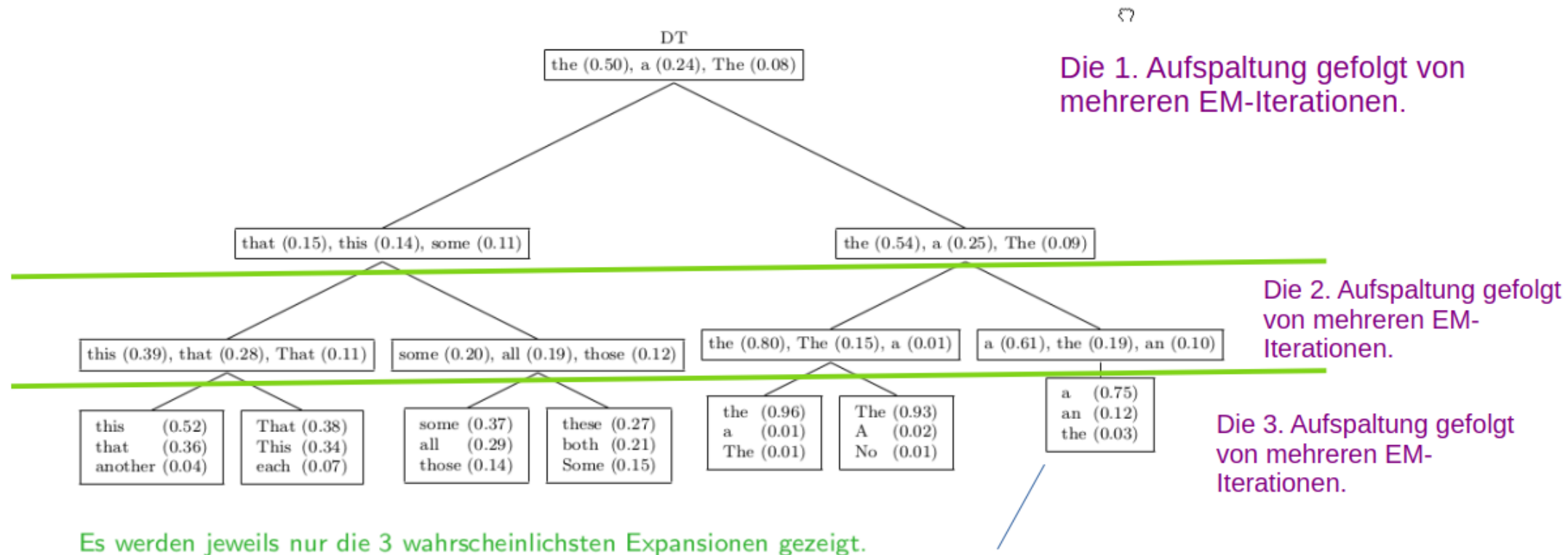


Die modifizierte Grammatik liefert für den alten Parsebaum $2^4 = 16$ neue Parsebäume.

Die Subkategorisieren können noch weiter verfeinert werden.

Rekursive Verfeinerung

Rekursives Aufspalten der Kategorien gefolgt von mehreren EM-Iterationen verfeinert die Kategorien immer mehr.



Wenn die Aufspaltung nicht mehr vorteilhaft ist, werden die Subkategorien wieder zusammengefasst

Beispiel für andere Regeln : hier sehen wir, nach mehrmaligen Aufspaltungen und EM-Trainings wurde ADVP in 9 Subkagegorien unterteilt.

ADVP			
ADVP-0	RB-13 NP-2	RB-13 PP-3	IN-15 NP-2
ADVP-1	NP-3 RB-10	NP-3 RBR-2	NP-3 IN-14
ADVP-2	IN-5 JJS-1	RB-8 RB-6	RB-6 RBR-1
ADVP-3	RBR-0	RB-12 PP-0	RP-0
ADVP-4	RB-3 RB-6	ADVP-2 SBAR-8	ADVP-2 PP-5
ADVP-5	RB-5	NP-3 RB-10	RB-0
ADVP-6	RB-4	RB-0	RB-3 RB-6
ADVP-7	RB-7	IN-5 JJS-1	RB-6
ADVP-8	RB-0	RBS-0	RBR-1 IN-14
ADVP-9	RB-1	IN-15	RBR-0

Original paper: Learning Accurate, Compact, and Interpretable Tree Annotation
<https://dl.acm.org/doi/pdf/10.3115/1220175.1220230>

Vereinigung

- Ab einem bestimmten Punkt ist eine weitere Aufspaltung der Kategorien nicht mehr vorteilhaft, weil sie zu **Sparse-Data-Problemen** führt.

Die Penn-Treebank verwendet bspw. eine eigene Kategorie „**,**“ für Kommas, deren Aufspaltung nichts nützt.

⇒ Unterkategorien, die sich zu ähnlich sind, werden daher nach dem EM-Training wieder **zusammengefasst**.

Strategie:

- Nach jedem EM-Training und vor dem Aufspalten der Kategorien, werden **50%** der vorherigen Aufspaltungen rückgängig gemacht.
- Für jede Aufspaltung wird berechnet, wie stark sich die **Wahrscheinlichkeit der Trainingsdaten** verringert, wenn die Aufspaltung rückgängig gemacht wird. (Details im Originalartikel)
- Die Aufspaltungen mit der **kleinsten Differenz** werden rückgängig gemacht.

Experimentelle Ergebnisse

- “Berkeley”-Parser von Slav Petrov und Dan Klein
- Der Parser liefert für viele verschiedene Sprachen und Baumbanken gute Ergebnisse
- Deutsch (Precision=80.1%, Recall=80.1%)
- Englisch (Precision=90.2%, Recall=89.9%)
- Chinesisch (Precision=84.8%, Recall=81.9%)

Aber: Neuere Parser auf Basis neuronaler Netze sind noch besser.

Kurzzusammenfassung: Berkeley-Parser

- ① gegeben: eine **Baumbank**
- ② **Binarisierung** der Parsebäume durch Einfügen von Hilfsknoten
- ③ Extraktion einer Grammatik mit Regelhäufigkeiten
- ④ für mehrere Iterationen
 - ① Aufspaltung jeder Kategorie X in zwei neue Kategorien X1 und X2
 - ② annähernd uniforme Verteilung der Häufigkeit einer Regel auf die neuen Regeln
 - ③ für mehrere Iterationen
 - ① Schätzung der Regelwahrscheinlichkeiten aus den Regelhäufigkeiten
 - ② Für jeden Parsebaum der Baumbank wird der Parsewald mit Analysen entsprechend der verfeinerten Grammatik berechnet.
 - ③ Mit dem Inside-Outside-Algorithmus werden erwartete Häufigkeiten für die Parsewaldregeln berechnet.
 - ④ Die erwarteten Regelhäufigkeiten werden über alle Sätze summiert.
 - ④ Die schlechtere Hälfte der Aufspaltungen wird rückgängig gemacht.