

L'idea è quella di ASSUMERE dei parametri iniziali ϕ e utilizzare un CORPUS di frasi di assistenza in modo ITERATIVO i valori delle varie stime. Ad esempio è possibile iniziare con una DISTRIBUZIONE UNIFORME su tutte le regole della grammatica con lo stesso non-terminale $A \rightarrow \cdot$.

La probabilità di una frase w è data da

$$P_{\phi}(w) = \sum_T P_{\phi}(w, T)$$

dove \sum_T è la sommatoria su tutti i parse tree della frase w , mentre ϕ è l'insieme dei PARAMETRI utilizzati per calcolare le probabilità.

Dato un CORPUS formato dalle frasi w_1, \dots, w_m , la LIKELIHOOD $L(\phi)$ del corpus è data da

$$L(\phi) = P_{\phi}(w_1) \cdot P_{\phi}(w_2) \dots P_{\phi}(w_m)$$

L'inside-outside algorithm parte da un insieme di parametri ϕ per OTTENERE un nuovo insieme di parametri ϕ' t.c. $L(\phi') \geq L(\phi)$. Questo processo di UPDATE dei parametri viene ripetuto fino a quando la LIKELIHOOD CONVERGE.

OSS: L'INSIDE-OUTSIDE ALGORITHM è un caso speciale dell'EXPECTATION MAXIMIZATION algorithm e viene utilizzato per MASSIMIZZARE LOCALMENTE la LIKELIHOOD dei dati di TRAINING.

Annunciamo SENZA PERDITA DI GENERALITÀ che la nostra PCFG sia in CNF. Fissati i parametri ϕ l'algoritmo procede nel seguente modo per aggiornarli:

$$\bullet \phi'(A \rightarrow BC) := \frac{\text{count}(A \rightarrow BC)}{\sum_d \text{count}(A \rightarrow d)}$$

$$\bullet \text{count}(A \rightarrow BC) := \sum_{i=1}^N C_\phi(A \rightarrow BC, w_i)$$

• $C_\phi(A \rightarrow BC, w_i) :=$ NUMERO ATTESO di volte in cui la regola $A \rightarrow BC$ viene utilizzata per generare la stringa w_i .

$$\bullet \phi'(A \rightarrow w) := \frac{\text{count}(A \rightarrow w)}{\sum_d \text{count}(A \rightarrow d)}$$

$$\bullet \text{count}(A \rightarrow w) := \sum_{i=1}^N C_\phi(A \rightarrow w, w_i)$$

• $C_\phi(A \rightarrow w, w_i) :=$ NUMERO ATTESO di volte in cui la regola $A \rightarrow w$ viene utilizzata per generare la stringa w_i .

Al fine di calcolare il valore di $P_\phi(A \rightarrow \alpha, w_i)$, ovvero il # medio di volte in cui viene utilizzata la regola $A \rightarrow \alpha$ per generare la frase w_i , introduciamo le seguenti probabilità

- INSIDE PROBABILITY

Denotate con $\alpha_{i,s}(A)$, rappresenta la probabilità che il non terminale A sia essere utilizzato per DERIVARE la SOTTOSTRINGA $w_i \dots w_s$ nella FRASE $W = w_1 \dots w_m$ dati i parametri specificati in ϕ .

In formule,

$$\alpha_{i,s}(A) := P_\phi(A \xrightarrow{*} w_i \dots w_s)$$

- OUTSIDE PROBABILITY

Denotate con $\beta_{i,s}(A)$, rappresenta la prob. di ottenere la stringa $w_1 \dots w_{i-1} A w_{s+1} \dots w_m$ partendo dal non-terminale S e utilizzando i parametri specificati in ϕ .

In formule,

$$\beta_{i,s}(A) := P_\phi(S \xrightarrow{*} w_1 \dots w_{i-1} A w_{s+1} \dots w_m)$$

Al fine di calcolare queste probabilità possiamo utilizzare le seguenti RELAZIONI RICORSIVE, che funzionano se la grammatica è in CNF

$$\alpha_{i,j}(A) = \sum_{B,C} \sum_{i \leq k \leq j} \phi(A \rightarrow BC) \cdot \alpha_{i,k}(B) \cdot \alpha_{k+1,j}(C)$$

per $i < j$. Se $i = j$ invece

$$\alpha_{i,i}(A) = \phi(A \rightarrow w_{ii})$$

$$\beta_{i,j}(A) = \sum_{B,C} \sum_{i \leq k} \phi(B \rightarrow CA) \alpha_{k,i-1}(C) \beta_{k,j}(B) + \sum_{B,C} \sum_{m \geq k > j} \phi(B \rightarrow AC) \cdot \alpha_{j+1,k}(C) \cdot \beta_{i,k}(B)$$

Utilizzando queste formule si è in grado di trovare ϕ e ϕ' . L'algoritmo GARANTISCE che la LOG-LIKELIHOOD non decresce, ovvero che

$$LL(\phi) := \sum_{i=1}^N \log P_{\phi}(w_i) \leq LL(\phi')$$

L'idea è quindi quella di SMETTERE quando la differenza $LL(\phi') - LL(\phi) > 0$ risulta essere ABBASTANZA PICCOLA.

L'implementazione di questi calcoli viene effettuata utilizzando l'algoritmo CYK: durante la fase di PARSING e di costruzione della CHART vengono calcolate le INSIDE PROBABILITIES con un approccio BOTTOM-UP; Successivamente vengono calcolate le OUTSIDE PROBABILITIES con un approccio TOP-DOWN.

In generale per l'algoritmo INSIDE-OUTSIDE si basa sul calcolo del NUMERO ATTESO di volte in cui applicare la data REGOLA al fine di generare una FRASE in modo da poter AGGIORNARE le stime dei parametri.