

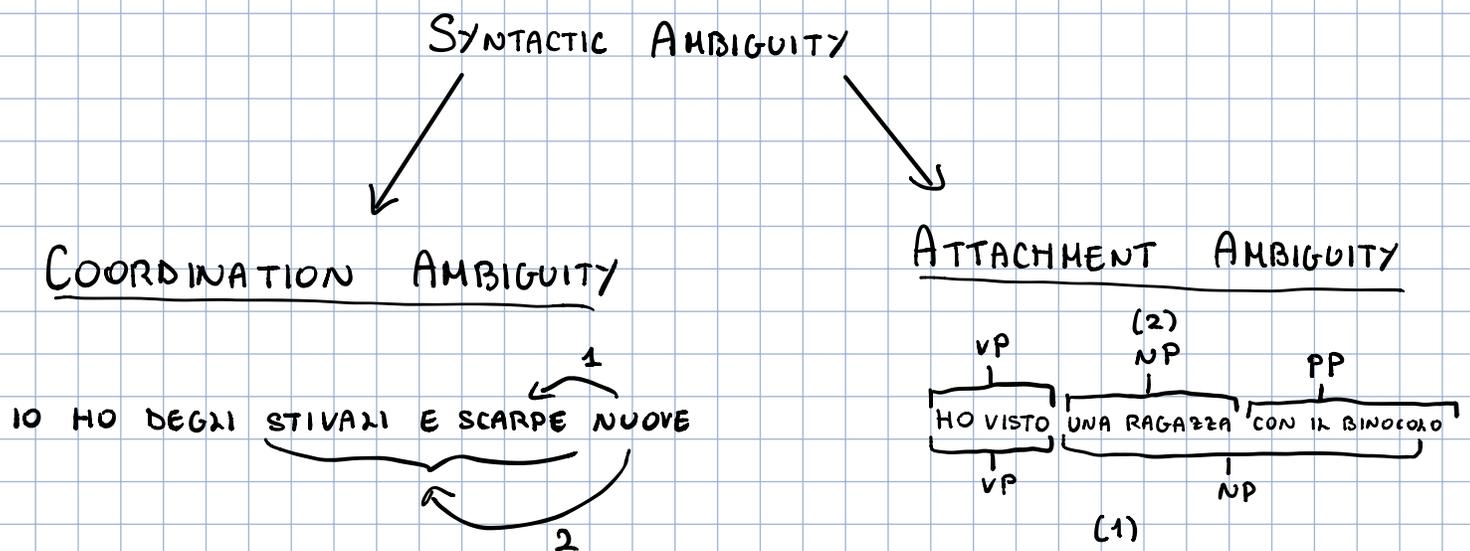
NLP - LEZIONE 11

DEL 07/11/2019

NCFG (MICHELE & SAMIR)

AMBIGUITÀ

Abbiamo due tipologie di AMBIGUITÀ SINTATTICHE, che sono



Non sappiamo se mi gli stivali e le scarpe sono nuove (1) oppure se solo le scarpe sono nuove.

Notiamo che secondo Chomsky l'ambiguità sintattica PREVALE su quella semantica in quanto la semantica è definita nel PARSE TREE.

PROBABILISTIC CF - GRAMMARS (PCFG)

Una PCFG è una quadrupla $\langle N, \Sigma, R, S \rangle$ simile a una tipica CF grammar ma con la differenza che le regole hanno la seguente forma

$$A \rightarrow \beta [P]$$

con,

- A non-terminale;
- β stringa di terminali e non-terminali;
- $P \in [0, 1]$ t.c. $P := P(\beta | A)$, ovvero la prob. che dato un non-terminale A , questo si scrive β .

Notiamo che vale la seguente EQUIVALENZA

$$\forall A \in N: \sum_{\beta: A \rightarrow \beta} P = 1$$

Le PCFG non sono state introdotte al fine di avere uno strumento di analisi utilizzare per gestire l'ambiguità sintattica: la **PROBABILITÀ**.

Utilizzando le PCFG e infatti non è possibile CALCOLARE la probabilità di ogni PARSE TREE. Una volta che viene in grado di fare questo, per rimuovere l'ambiguità sintattica ci basta considerare, data una frase S , solamente il parse tree che MASSIMIZZA la sua probabilità, o se vogliamo anche N parse tree più probabili.

Al fine di calcolare la prob. di un PARSE TREE, utilizziamo le seguenti IPOTESI DI INDIPENDENZA:

1) INDIPENDENZA DAL CONTESTO

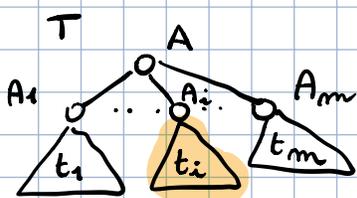
Dato un parse tree t con radice A , siano t_1, \dots, t_m i sub parse-tree e siano A_1, \dots, A_m i nodi radice corrispondenti.

L'ipotesi di INDIPENDENZA DAL CONTESTO ci dice che

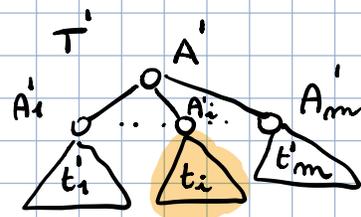
$$P(t_1 \mid A, A \rightarrow A_1 \dots A_m, t_2, \dots, t_m) = P(t_1)$$

ovvero non importa chi sono i FRATELLI di t_1 , la sua probabilità non cambia.

Graficamente abbiamo,



La prob. di t_1 è la stessa in entrambi gli alberi



2) INDIPENDENZA DALLA POSIZIONE

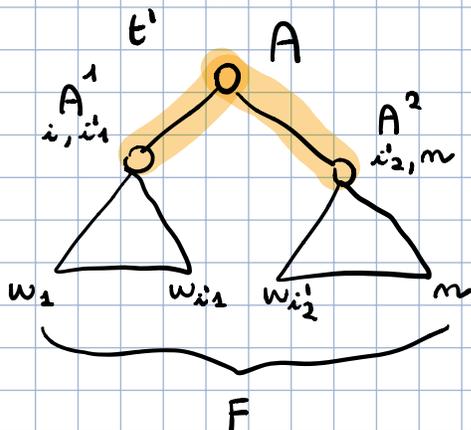
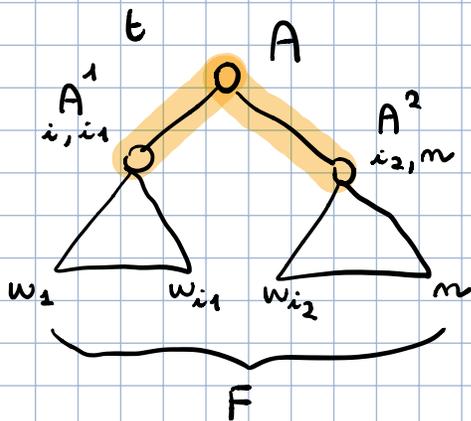
Dato un nave tree t con radice A , siano t_1, \dots, t_m i sub nave-tree e siano $A_{1,3_1}^1 \dots A_{3_m, m}^m$ i nodi radice corrispondenti. Con la notazione $A_{1,3_1}^1$ vogliamo indicare che il sotto-albero t_1 SPAZIA nella sottoparte $F[1 \dots 3_1]$.

L'insieme di INDIPENDENZA DALLA POSIZIONE ci dice che

$$P(A_{1,m} \rightarrow A_{1,3_1}^1 \dots A_{3_m, m}^m) = P(A_{1,m} \rightarrow A^1 \dots A^m)$$

ovvero che non ci interessa la posizione relative del non terminale nella radice.

Graficamente abbiamo,



Anche se li stessi non terminali A^1 e A^2 nascono la stringa F in due modi diversi in t e t' , la mb. di nascono da A e $A^1 A^2$ è la STESSA.

Supponiamo di avere due alberi t_1 e t_2 per una data SENTENZA S . Vogliamo trovare l'albero t^* tale che

$$t^* = \text{ARG MAX} \{ P(t_1 | S), P(t_2 | S) \}$$

Come facciamo?

Notiamo inizialmente che non è possibile STIMARE $P(T|S)$ in quanto, anche utilizzando BAYES, è ormai IMPROBABILE trovare la STESSA FRASE con tanti ALBERI DIFFERENTI all'interno di un corpus.

Fortunatamente è possibile dim. il seguente fatto: la probabilità di un PARSE TREE T è pari al PRODOTTO DELLE PROBABILITÀ ASSOCIATE ALLE m REGOLE UTILIZZATE PER ESPANDERE OGNUNO DEGLI n NODI NON-TERMINALI.

Formalmente,

$$P(T | G) = \prod_{(A \rightarrow \gamma) \in T} P(A \rightarrow \gamma) \quad (1)$$

Inoltre, dato che $P(S|S) = 1$, abbiamo

$$P(T, S) = P(S|T) \cdot P(T) = P(T) \quad (2)$$

Andiamo odess e dimostrare (1)...

DIM.: Procediamo per induzione su $m := \#$ di noduzioni utilizzate per ottenere t .

- CASO BASE, $m = 1$

Se t e formato da una sola noduzione

$A_{11} \rightarrow a$, allora abbiamo che

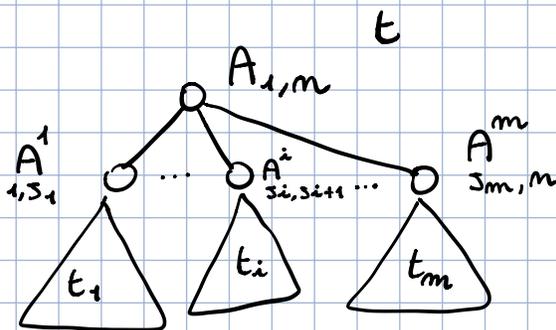
$$P(t|G) = P(A_{11} \rightarrow a) \stackrel{\text{(IND. DALLA POSIZIONE)}}{=} P(A \rightarrow a)$$

- PASSO INDUTTIVO, $m > 1$

Sia $A_{1,m}$ la radice di t , siano t_1, \dots, t_m i sottoalberi di t , e siano $A_{1,s_1}^1, \dots, A_{s_m, s_m}^m$ le radici dei rispettivi sottoalberi.

Con A_{1,s_1}^1 stiamo nuovamente indicando il fatto che il non-terminale A^1 inizia nella SOTTOFRASE $F[1 \dots s_1]$.

Graficamente,



La probabilità di t è quindi data da

$$P(t|G) = P(A_{1,m}, A_{1,m} \rightarrow A_{1,31}^1 \cdots A_{3m,m}^m, t_1, t_2, \dots, t_m)$$

Applicando la CHAIN RULE, che ci dice che $P(B, C) = P(B|C) \cdot P(C)$, otteniamo la seguente

$$P(A_{1,m}, A_{1,m} \rightarrow A_{1,31}^1 \cdots A_{3m,m}^m, t_1, t_2, \dots, t_m) =$$

$$P(t_1 | A_{1,m}, A_{1,m} \rightarrow A_{1,31}^1 \cdots A_{3m,m}^m, t_2, \dots, t_m) \cdot$$

$$P(t_2 | A_{1,m}, A_{1,m} \rightarrow A_{1,31}^1 \cdots A_{3m,m}^m, t_3, \dots, t_m) \cdot$$

⋮

$$P(t_m | A_{1,m}, A_{1,m} \rightarrow A_{1,31}^1 \cdots A_{3m,m}^m) \cdot$$

$$P(A_{1,m} \rightarrow A_{1,31}^1 \cdots A_{3m,m}^m | A_{1,m}) \cdot P(A_{1,m})$$

Notiamo che dalle ipotesi di indipendenza
assunte prima otteniamo che $\forall i=1, \dots, m$:

$$i) P(t_i | A_{1,m}, A_{1,m} \rightarrow A_{1,31}^1 \cdots A_{3m,m}^m, t_2, \dots, t_m) = P(t_i)$$

$$ii) P(A_{1,m}, A_{1,m} \rightarrow A_{1,31}^1 \cdots A_{3m,m}^m) = P(A, A \rightarrow A^1 \cdots A^m)$$

Infine notiamo che, dato che la nozione di t, A , è FISSA, abbiamo che



$$P(A, A \rightarrow A^1 \dots A^m) = P(A \rightarrow A^1 \dots A^m)$$

Mettendo tutto insieme e applicando l'IPOTESI INDUTTIVA, che ci dice che $P(t_i) = \prod_{A \rightarrow \alpha \in t_i} P(A \rightarrow \alpha)$, troviamo

$$\begin{aligned} P(t|G) &= P(A_{1,m}, A_{1,m} \rightarrow A_{1,1}^1 \dots A_{m,m}^m, t_1, t_2, \dots, t_m) \\ &= \prod_{i=1}^m P(t_i) \cdot P(A \rightarrow A^1 \dots A^m) \\ &= \prod_{i=1}^m \prod_{A \rightarrow \alpha \in t_i} P(A \rightarrow \alpha) \cdot P(A \rightarrow A^1 \dots A^m) \\ &= \prod_{A \rightarrow \alpha \in t} P(A \rightarrow \alpha) \end{aligned}$$



LANGUAGE MODELS

~ (1:06:00 min)

Il LANGUAGE MODELING viene utilizzato per calcolare la probabilità di una FRASE in una LINGUA, e viene utilizzato come tool di GENERAZIONE, come ad esempio durante la TRADUZIONE.

La probabilità di una frase non AMBIGUA è la probabilità del suo unico parse tree. Per le frasi AMBIGUE invece abbiamo

$$P(S) = \sum_{T \in \text{TREES}(S)} P(T)$$

l'insieme dei parse trees di S.

CKK PER PCFG

È possibile estendere l'algoritmo CKK per gestire PCFG messe in forma CNF.

Per fare questo si MEMORIZZA, per ogni non-term. nella matrice un valore di prob. Quando confrontiamo due non-term A, B e troviamo un match con la regola $C \rightarrow AB$ la prob. del non-term C sarà data da

$$P(C) = P(A) \cdot P(B) \cdot P(C \rightarrow AB)$$

dove $P(A)$ e $P(B)$ sono i valori di prob. associati ai simboli A e B.

Per diminuire la COMPLESSITÀ dell'algoritmo possiamo scegliere di mantenere solamente le N prob. più alte. Questo ci permette di confrontare

o il più N·N coppie di non-terminali ad ogni passo dell'algoritmo.

STIMARE PROBABILITÀ NELLE PCFGS ~ (1:24:00 min)

Per stimare le prob. in una PCFGS abbiamo due approcci:

1) Dato un TREEBANK possiamo utilizzare la MASSIMA VEROSIMIGLIANZA per stimare le probs

$$P(\alpha \rightarrow \beta | \alpha) = \frac{\text{COUNT}(\alpha \rightarrow \beta)}{\sum_Y \text{COUNT}(\alpha \rightarrow Y)} = \frac{\text{COUNT}(\alpha \rightarrow \beta)}{\text{COUNT}(\alpha)}$$

2) Utilizzando l'algoritmo INSIDE-OUTSIDE.

ALGORITHM INSIDE-OUTSIDE

L'idea è quella di ASSUMERE dei parametri iniziali Φ e utilizzare un CORPUS al fine di aggiornare 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 nuovo 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.

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

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

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

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

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

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

$\bullet C_{\phi}(A \rightarrow w, W_i) :=$ 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 può 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

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

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

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 la FRASE in modo da poter AGGIORNARE le stime dei parametri.