

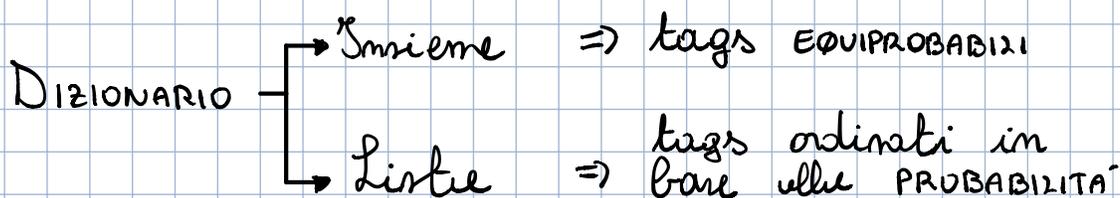
In questa lezione vedremo dei nomi di
metodi per effettuare il PART-OF-SPEECH TAGGING.

POS TAGGING (CONT.)

Date n parole $w_1 w_2 \dots w_m$ vogliamo
assegnare ad ogni parola w_i ed w_m
un tag $T_1 T_2, \dots, T_m$.

Notiamo che se vedo la parola "la", a
questa possiamo assegnare i TAG: NOUN, PRON, DET, ADV.
Se nei dati "la" c'è in NOUN, allora non vuol
dire che il "la" non è in ARTICOLO.

Una prima idea è quella di utilizzare un
DIZIONARIO V che ci dice per ogni parola quali
sono i possibili TAGS che possiamo assegnare
a quella parola.



RULE BASED POS TAGGING ~ (18:00 min)

Cominceremo definendo le coppie di tags validi

- ARTICOLO NOME
- NOME VERBO
- AUSILIARE VERBO
- NOME AGGETTIVO

Continuando con questo approccio non dovremmo costruire l'ALBERO ultraverso ma grammaticale che non abbiamo.

Quello che vogliamo è un metodo che ci permetta di ridurre il # di tags per ogni parola. Il POS-TAGGING è stato proprio utilizzato per ridurre lo SPAZIO DI RICERCA per lo step successivo del PARSING.

Un primo metodo del genere è stato introdotto da BRILL e funziona come segue:

- 1) Assegnare ad ogni parola il TAG con la MAGGIORE prior probability.

2) Scons ogni nave adiacente e se
know coppie di navi "strane" e
TRASFORMO utilizzando delle regole.

L'idea è quindi quella di utilizzare dei
CONTROLLORI DI STRANEZZE che trasformano i
tags associati alle navi. Dato che i
controllori NON INFLUENZARSI in modo
indiretto, l'ORDINE in cui applico questi
controllori è molto importante. Una volta
che ho applicato tutti i controllori ritorna
i TAGS originali.

Il metodo descritto prima è detto TRANSFORMATION
BASED POS-TAGGING. Gli aspetti critici di questo
metodo sono i seguenti:

- i) Definizione delle REGOLE DI TRASFORMAZIONE.
- ii) Definizione dell'ORDINE in cui applichiamo
i vari controllori.

In particolare nominerò accoppiare a
questo metodo un modulo che APPRENDE
l'ordine migliore in cui applicare i
controllori.

HIDDEN MARKOV MODEL POS TAGGING ~ (52:00 min)

Date la frase w_1, \dots, w_m vogliamo trovare dei tag che MASSIMIZZANO la probabilità

$$P(t_1, t_2, \dots, t_m \mid w_1, w_2, \dots, w_m)$$

Ovvero vogliamo trovare $\bar{t}_1, \bar{t}_2, \dots, \bar{t}_m$ e.c.

$$(\bar{t}_1, \dots, \bar{t}_m) = \underset{t_1, \dots, t_m}{\text{ARGMAX}} \left\{ P(t_1, \dots, t_m \mid w_1, \dots, w_m) \right\}$$

Un primo approccio potrebbe essere quello di stimare $P(t_1 \dots t_m \mid w_1 \dots w_m)$ contando quante volte la sequenza di tags (t_1, \dots, t_m) è associata alle frase $w_1 \dots w_m$ nel nostro DATASET. Notiamo però che questo è estremamente difficile, in quanto raramente in un qualsiasi CORPUS la sequenza di tags $t_1 \dots t_m$ è associata alle frase $w_1 \dots w_m$. Dovrò dunque trasformare le formule in un qualcosa che sia STIMABILE. Per fare questo è possibile utilizzare delle ipotesi di INDIPENDENZA tra i vari tags.

Coni precedenti otteniamo che

$$P(t_1 \dots t_m | w_1 \dots w_m) \approx \prod_{i=1}^m P(t_i | w_1 \dots w_m)$$

Qui ipotizziamo che
i TAGS non ind.
tra loro.

$$\approx \prod_{i=1}^m P(t_i | w_{i-1} w_i w_{i+1})$$

Qui ipotizziamo che il tag associato
alle i -esime parole dipende solo dalle
parole ADIACENTI alle i -esime parole.

Anche se adesso è più facile stimare la prob.,
il lavoro può comunque essere molto ONEROSO, in
quanto lavorando con un DIZIONARIO di 30K
parole dovremmo comunque stimare $(30K)^3$
distribuzioni di prob., ciascuna associata ad
una tripletta di parole.

Notando che non abbiamo solamente 15 tags,
possiamo ricominciare applicando BAYES come segue

$$P(t_1 \dots t_m | w_1 \dots w_m) = \frac{P(w_1 \dots w_m | t_1 \dots t_m) \cdot P(t_1 \dots t_m)}{P(w_1 \dots w_m)}$$

Dato che siamo interessati solamente all'ARGMAX, possiamo tranquillamente non considerare il DENOMINATORE. Vogliamo quindi massimizzare la seguente espressione

$$P(w_1 \dots w_m | t_1 \dots t_m) \cdot P(t_1 \dots t_m)$$

Se assumiamo poi che le PAROLE SONO IND., troviamo che

$$\begin{aligned} P(w_1 \dots w_m | t_1 \dots t_m) &\approx \prod_{i=1}^m P(w_i | t_1 \dots t_m) \\ &\approx \prod_{i=1}^m P(w_i | t_i) \end{aligned}$$

A questo punto notiamo che possiamo stimare $P(w_i | t_i)$ tranquillamente contando nel CORPUS quante volte il tag t_i è ASSOCIATO alle parole w_i .

Per quanto riguarda $P(t_1 \dots t_m)$ abbiamo che

$$\begin{aligned} P(t_1 \dots t_m) &= P(t_1 | t_2 \dots t_m) \cdot P(t_2 \dots t_m) \\ &= P(t_1 | t_2 \dots t_m) \cdot P(t_2 | t_3 \dots t_m) \dots P(t_{m-1} | t_m) \cdot P(t_m) \\ &= \prod_{i=1}^m P(t_i | t_{i+1} \dots t_m) \cdot P(t_m) \end{aligned}$$

Se noi assumiamo che OGNI TAG DIPENDE O
DAL PRECEDENTE O DAL SUCCESSIVO (o da uno dei due)
troviamo la seguente semplificazione

$$\prod_{i=1}^m P(t_i | t_{i+1} \dots t_m) \cdot P(t_m) \approx \prod_{i=1}^m P(t_i | t_{i-1}) \cdot P(t_m)$$

Ci rimangono quindi i seguenti problemi:

i) Come troviamo l'ARGMAX?

Siamo infatti invitati a dover
calcolare l'ARGMAX della seguente
espressione

$$\prod_{i=1}^m P(w_i | t_i) \cdot P(t_i | t_{i-1}) \cdot P(t_m)$$

ii) Come stimiamo le probabilità?

Vedere la sezione successiva.

MASSIMA VEROSIMIGLIANZA

~ (1:26:00 min)

Perché per stimare le prob. devo contare il # di casi favorevoli rispetto al # di eventi totali?

Per rispondere a tale domanda introduciamo il concetto di LIKELIHOOD.

$$\text{LIKELIHOOD} := P(\text{OSSERVAZIONE} \mid \text{PARAMETRI } \theta)$$

Il metodo di MASSIMA VEROSIMIGLIANZA ci permette di stimare i PARAMETRI di una distribuzione andando a massimizzare la VEROSIMIGLIANZA.

ESEMPIO: Consideriamo n lanci di una moneta

$$X_1, \dots, X_n \text{ i.i.d. con } X_i \sim \text{Bern}(\pi)$$

Volendo applicare il metodo di MASSIMA VEROSIMIGLIANZA per stimare π procediamo come segue

$$P(X_1 \dots X_n \mid \pi) = \pi^{x_1 + \dots + x_n} \cdot (1 - \pi)^{n - (x_1 + \dots + x_n)}$$

$$\log \downarrow \log P(X_1 \dots X_n \mid \pi) = (x_1 + \dots + x_n) \cdot \log \pi + (n - \sum_i x_i) \log(1 - \pi)$$

$$\frac{d}{d\pi} \downarrow \frac{d}{d\pi} \log P(X_1 \dots X_n \mid \pi) = \frac{\sum_i x_i}{\pi} - \frac{n - \sum_i x_i}{1 - \pi}$$

Troviamo quindi

$$\frac{\sum_i x_i}{n} - \frac{n - \sum_i x_i}{1-n} = 0 \Leftrightarrow \frac{\sum_i x_i - n \frac{\sum_i x_i}{n} - n \cdot n + n \frac{\sum_i x_i}{n}}{n(1-n)} = 0$$

$$\Leftrightarrow \sum_i x_i - n \cdot n = 0$$

$$\Leftrightarrow n = \frac{\sum_i x_i}{n}$$

Troviamo quindi che per stimare n dobbiamo contare casi favorevoli su casi continui.

Nel nostro caso particolare possiamo quindi stimare le prob. come segue

$$P(t_i | t_{i-1}) = \frac{c(t_i, t_{i-1})}{c(t_{i-1})}$$

di volte
in cui abbiamo
esattamente i
tags t_i e t_{i-1}
nel CORPUS.

$$P(w_i | t_i) = \frac{c(w_i, t_i)}{c(t_i)}$$

In questi ultimi casi può accadere che $c(w_i, t_i) = 0$ in quanto noi non avremo mai visto la parola w_i nel CORPUS. Per risolvere questo problema esistono varie tecniche, tra cui una che prende il nome di GOOD-TURING FREQUENCY ESTIMATION

LAVORARE SULL' ARGMAX ~ (1:48:00 min)

Vogliamo trovare l'ARGMAX della seguente espressione

$$\prod_{i=1}^m P(w_i | t_i) \cdot P(t_i | t_{i-1}) \cdot P(t_m)$$

Per fare questo possiamo partire scegliendo il TAG t_1 che massimizza la probabilità e continuare mantenendo di volta in volta il più K SEQUENZE DI TAGS DISTINTE e prendendo la sequenza con la prob. maggiore.

Per fare questo devo introdurre il tag $t_0 := \text{START}$ e devo stimare le prob $P(t | \text{START})$.

Questo metodo è chiamato LINK SEARCH.
Se noi $K=1$ allora otteniamo il MONTE CARLO
ALGORITHM.

Notiamo che questo metodo non sempre ci
permette di ottenere il MASSIMO ASSOLUTO.
Detto questo è FORSE possibile dimostrare
che con ALTA PROBABILITA' otteniamo un buon
MASSIMO.

NOTE FINALI ~ (2:02:00 min)

Breve riferimento ai concetti di
ACCURACY, PRECISION e RECALL.