

Linguistica Computazionale

Un corpus, può essere rappresentativo per:

- Contesto: ossia contenente token specifici per un determinato ambito
- Lingua: ossia contenente token di diversi campi affinché copra il più possibile una lingua

Generalmente un corpus non copre interamente la lingua

Per definizione: un'opera di selezione per tipo

Corpus

Un corpus può essere anche classificato per token contenuti.
"The larger, the better"

Anche l'intero web può essere utilizzato come corpus, ma non totalmente affidabile per via di dati possibilmente errati (*rumore*)

Corpus principali conosciuti

Brown Corpus, 1961
US English
Circa un milione di token

British National Corpus (BNC), misto
British English
Circa 100 milioni di token, 90% scritto e 10% orale

La Repubblica
Italiano
Circa 326 milioni di token
Tema generale, specialistico nella tipologia testuale: scritto, monolingua e annotato

Parametri di un corpus

Generalità
generali o specialistici

Modalità:
scritto, audio o misto

Cronologia:
diacronico o sincronico

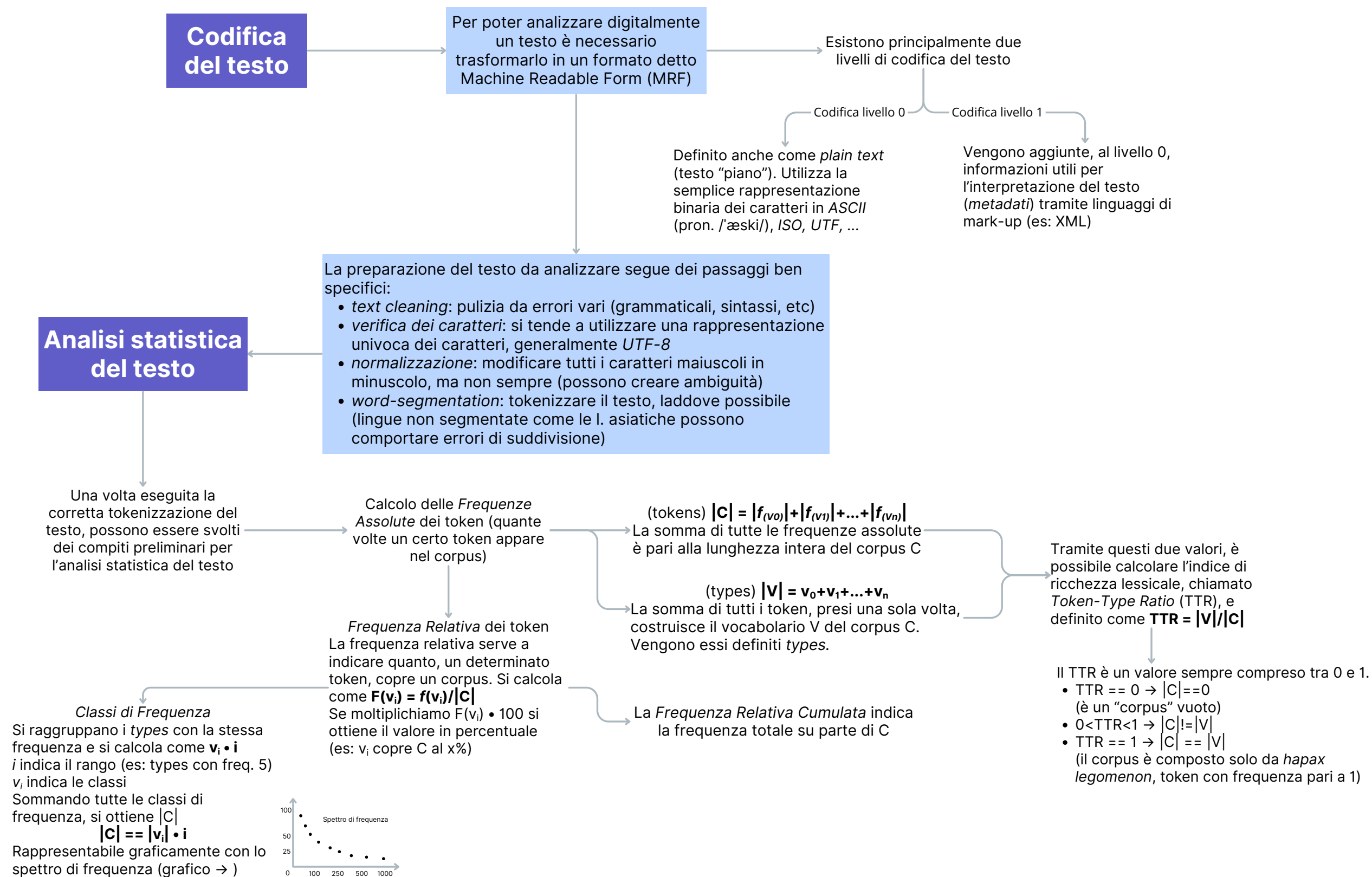
Lingua:

- monolingua
- paralleli (più lingue)
- comparabili

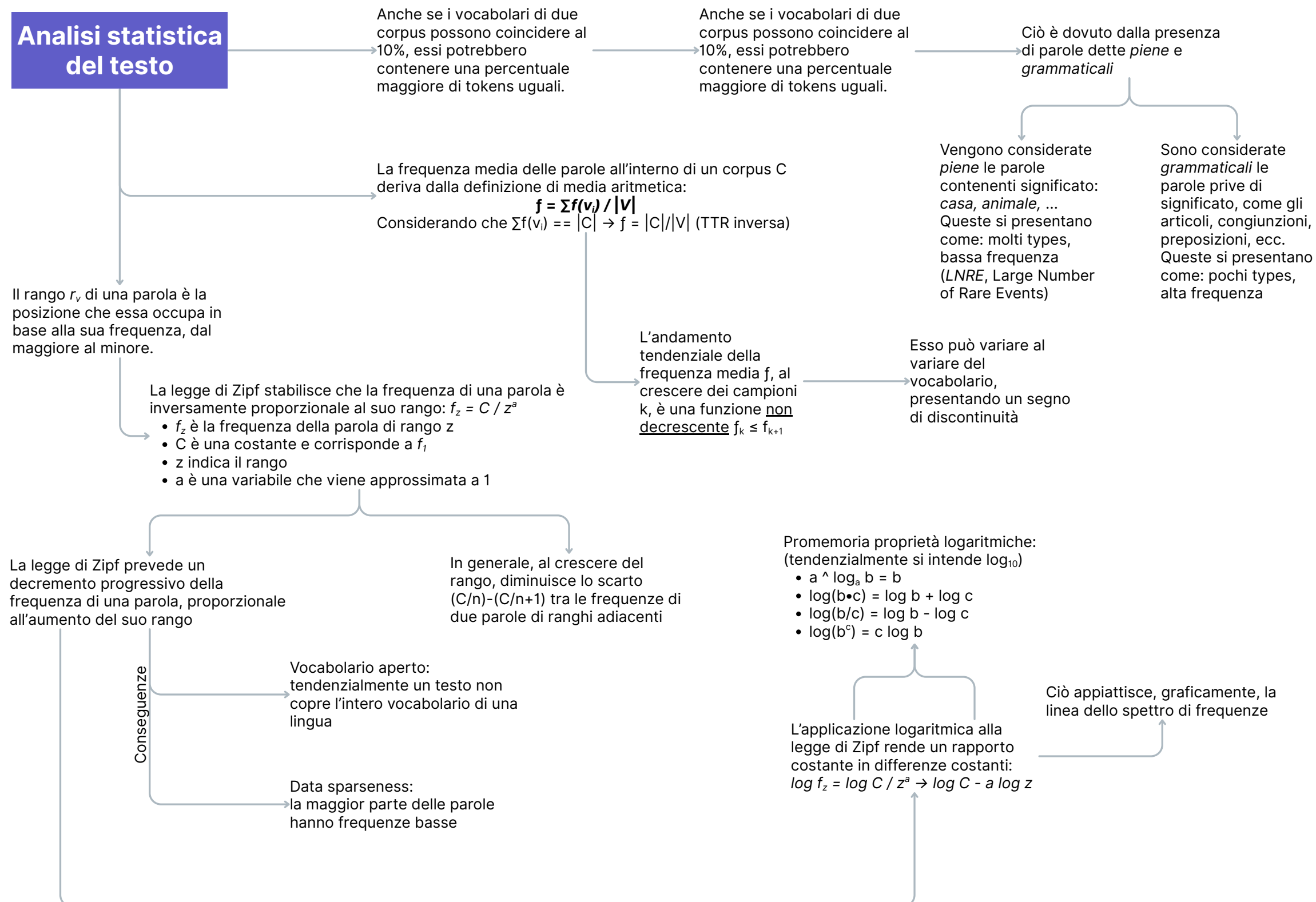
Integrità dei testi

Livello di codifica
digitale dei testi

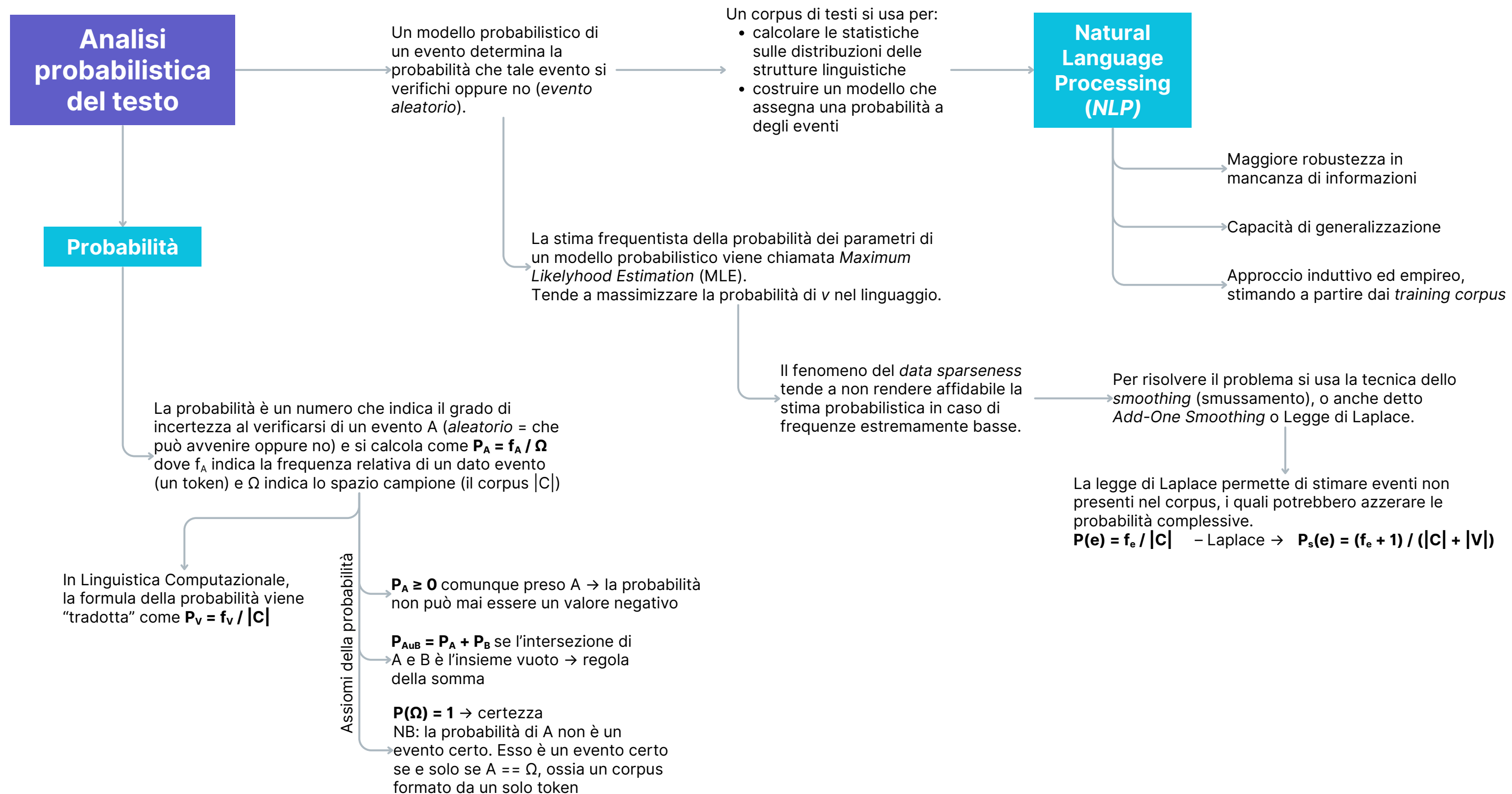
Linguistica Computazionale



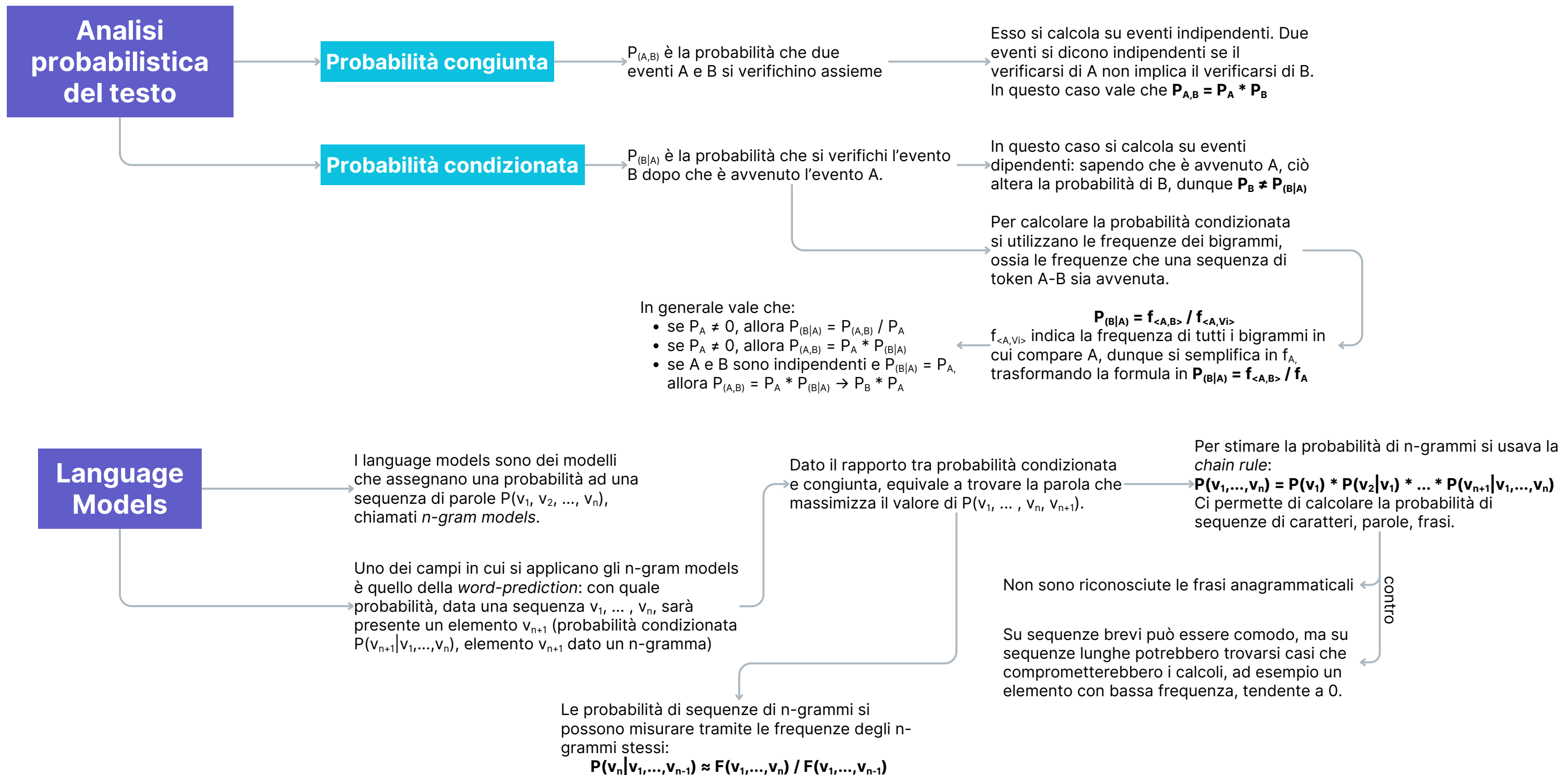
Linguistica Computazionale



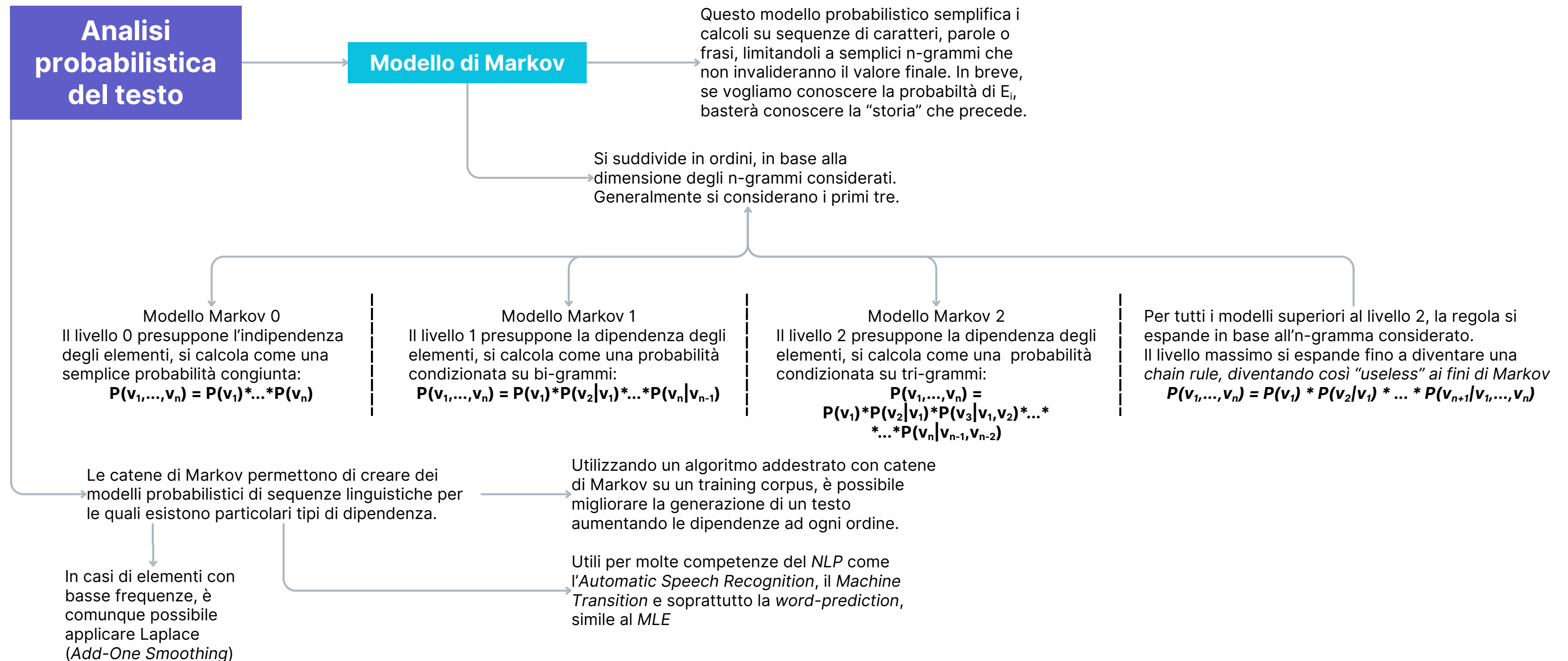
Linguistica Computazionale



Linguistica Computazionale



Linguistica Computazionale



Linguistica Computazionale

Combinazioni

Alcune parole sono legate ad altre tramite forti lessicazioni morfologiche e semantiche. Esse possono essere sostituite con altre parole per ottenere altre frasi grammaticali.

Esistono altri tipi di combinazioni che si basano su legami non riconducibili a classi linguistiche generali. Sono difficili da sostituire perché produrrebbero strani risultati.

Combinazioni di due o più parole caratterizzate da un elevato grado di associazione, determinata dalla tendenza di co-occorrere:

- argomenti o modificatori "tipici"
- argomenti o modificatori "idiosincratici"
- costruzioni idiomatiche
- nomi propri composti
- ...

Collocazioni

Distinguibili in due categorie

Sono parole con un alto grado di associazione reciproca. Le misure di associazione lessicale:

Quantificano la forza del legame tra due o più parole sul piano sintagmatico

La nozione intuitiva di associazione lessicale viene trasformata in un indice quantitativo e misurabile

Empirico, o senso ampio: combinazioni ricorrenti e predicibili di parole, osservate nell'uso linguistico (*corpora*)

Teorico, o senso stretto: espressioni polirematiche fortemente lessicalizzate, idiomatiche e idiosincratiche (*multiword expressions*)

Elevata convenzionalità
Sono tendenzialmente espressioni di uso convenzionali, tipici di varietà linguistica

Ridotta composizionalità semantica
Non immediatamente ricavabile dalla composizione delle parole che lo formano, ad esempio:
{topolino grigio} ≠ {topolino} + {grigio}
{gatta morta} ≠ {gatta} + {morta}

Forte rigidità strutturale
Sono spesso resistenti a modificazioni aggettivali o avverbiali, oppure occorrono solo in particolari forme flesse e contesti sintattici.

Linguistica Computazionale

Collocazioni

- 1 → Analisi linguistica del corpus: il testo deve essere tokenizzato e possibilmente annotato con PoS (*Part of Speech*) tagging, lemmatizzazione, ecc.
- 2 → Selezione dei bigrammi: il tipo dei bigrammi che vengono selezionati dipende dal livello di annotazione.
- 3 → Costruzione della tabella di contingenza: dal totale dei bigrammi costruisco la tabella.
- 4 → Applicazione di una misura di associazione
- 5 → Ordinamento delle coppie in base alla forza di associazione.

Due parole si dicono fortemente associate quando si presentano insieme più spesso rispetto alle singole frequenze.

È necessario confrontare la frequenza osservata O di $\langle x, y \rangle$ con la sua frequenza attesa E , ossia la frequenza che ci aspettiamo qualora gli elementi fossero statisticamente indipendenti.

Un'altra versione con cui si calcola la MI in termini di probabilità:
 $MI_{\langle x, y \rangle} = \log_2 P_{x,y} / (P_x P_y) == \log_2 (f_{\langle x, y \rangle} / N) / ((f_x / N) (f_y / N)) \rightarrow$
 $\rightarrow (f_{\langle x, y \rangle} \cdot N) / (f_x f_y)$

Se x e y ricorrono sempre insieme:
 $MI_{\langle x, y \rangle} = f \cdot N / f^2 = N / f$

Estremamente sensibile agli eventi rari.
 In qualsiasi corpus i bigrammi formati da hapax avranno valori massimi.

Preso il bigramma $\langle x, y \rangle$ avremo

	x	$\neg x$	
y	O_{11}	O_{12}	R_1
$\neg y$	O_{21}	O_{22}	R_2
	C_1	C_2	

- O_{11} : bigramma dove occorrono x e y
- O_{12} : occorre y , ma non x
- O_{21} : occorre x , ma non y
- O_{22} : non occorrono né x né y
- $N = R_1 + R_2 == C_1 + C_2$

Partendo dalla tabella di contingenza:

	x	$\neg x$	
y	E_{11}	E_{12}	
$\neg y$	E_{21}	E_{22}	

- E_{11} : $(R_1 \cdot C_1) / N$
- E_{12} : $(R_1 \cdot C_2) / N$
- E_{21} : $(R_2 \cdot C_1) / N$
- E_{22} : $(R_2 \cdot C_2) / N$

Con entrambe le tabelle, è possibile calcolare la *Mutual Information* (MI):

$$MI_{\langle x, y \rangle} = \log_2 O_{\langle x, y \rangle} / E_{\langle x, y \rangle}$$

Se la $MI \leq 0$, c'è assenza di associazione tra le parole, altrimenti potrebbe essere una forte associazione.

La *Local Mutual Information* (LMI) privilegia i bigrammi più frequenti ed è il termine fondamentale nel calcolo del *Log Likelihood Ratio*
 $LMI_{\langle x, y \rangle} = f_{\langle x, y \rangle} \cdot MI_{\langle x, y \rangle}$

Linguistica Computazionale

Legge di Bayes

Viene usata per decidere qual è l'ipotesi I , o classe, più probabile all'interno di IP che spiega un certo tipo di osservazioni O

$\text{argmax}_{I \in \text{IP}} P_{(I|O)} = \text{argmax}_{I \in \text{IP}} (P_I * P_{(O|I)}) / P_O$
Poiché stiamo cercando l'ipotesi più probabile (argmax), data la stessa osservazione O , possiamo ignorare P_O :

$$\text{argmax}_{I \in \text{IP}} P_{(I|O)} = \text{argmax}_{I \in \text{IP}} P_I * P_{(O|I)}$$

Chiamato *teorema della probabilità delle cause*, viene impiegato per calcolare la probabilità di una causa che ha provocato l'evento verificato.

Molti fenomeni linguistici possono essere modellati come processo di inferenza bayesiana, come:

Riconoscimento del parlato come l'*Automatic Speech Recognition (ASR)*

Problema di decidere a quale classe appartiene una certa osservazione linguistica

Deriva della definizione di probabilità condizionata

$$P_{(A|B)} = P_{\langle A, B \rangle} / P_B \rightarrow P_{\langle A, B \rangle} = P_B * P_{(B|A)}$$

$$P_{(A|B)} = P_{\langle A, B \rangle} / P_B = P_{(B|A)} P_A / P_B \text{ dove:}$$

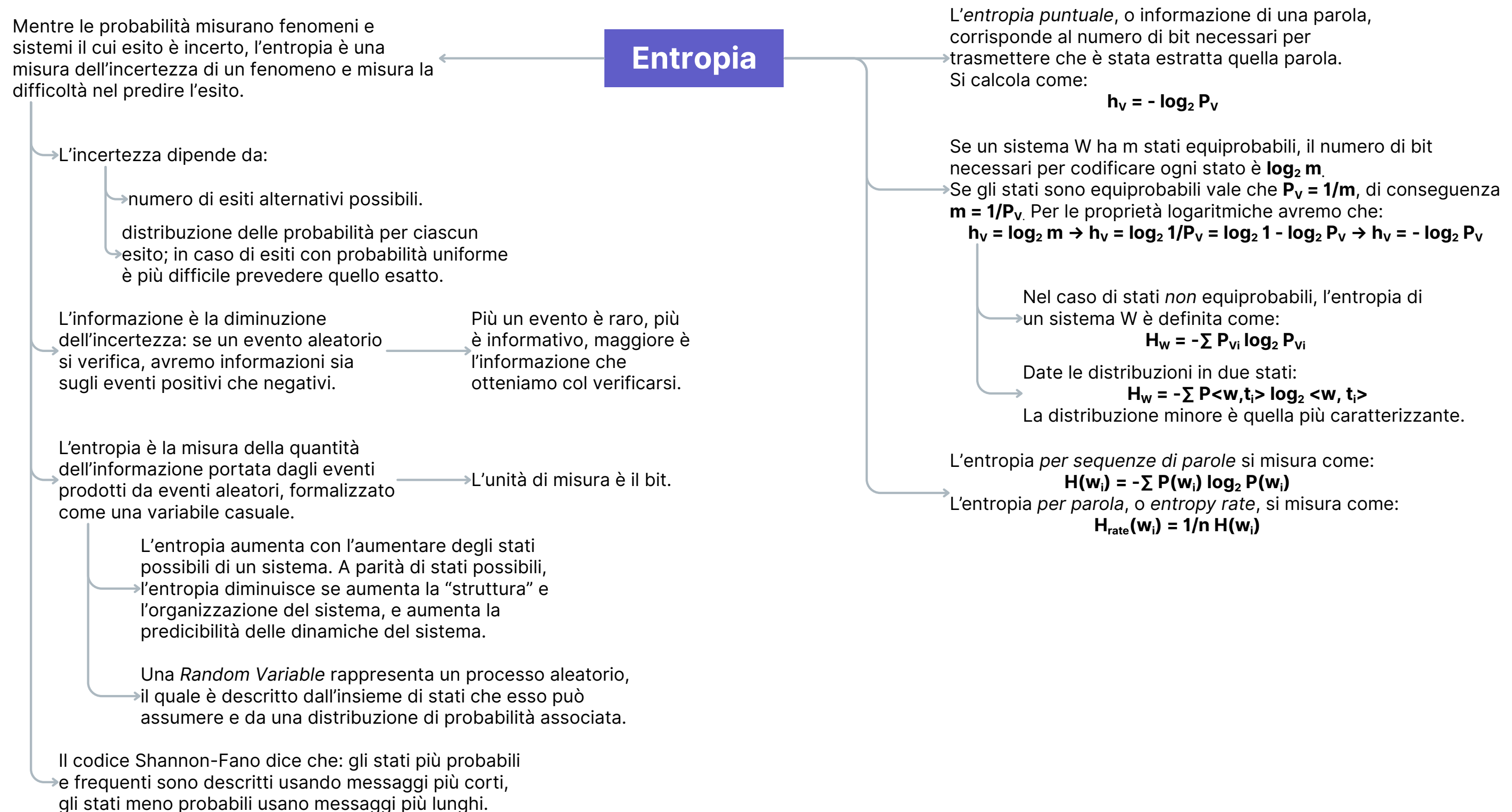
- $P_A \rightarrow$ probabilità a priori
- $P_{(B|A)} \rightarrow$ likelihood
- $P_{(A|B)} \rightarrow$ probabilità a posteriori

Ci permette di tradurre $P_{(I|O)}$ in un prodotto di probabilità più facile da stimare.

Noisy Channel Mode

Un sistema dove si introduce un input in un canale "rumore", alternandone l'output. Il canale rumoroso contiene gli elementi dell'inferenza bayesiana, che produce in output un'ipotesi con una certa probabilità.

Linguistica Computazionale



Linguistica Computazionale

Entropia

Shannon-McMillan-Breiman: se un linguaggio L è generato da un processo *stocastico ergodico* e *stazionario*, allora vale che:

$$H_{\text{rate}} L = \lim_{n \rightarrow \infty} -1/n * \log_2 P(w_1, \dots, w_n)$$

Cross Entropy: $H(w_1, m) = -\sum P(w_1) \log_2 P(w_1)$

- w_1 è una variabile casuale con distribuzione reale P
- m è un modello stocastico di w_1 che cerca di approssimare la sua distribuzione reale.

La cross entropy ci consente di misurare quanto bene un modello probabilistico approssima un certo processo stocastico $H(w) \leq H(w, m)$
È il costo in bits per usare m come modello descrittivo di un processo con distribuzione P .

Per Shannon-McMillan-Breiman è possibile approssimare la cross entropy prendendo un campione sufficientemente grande di testi del linguaggio come unica sequenza di parole

$$H(p, m) = \lim_{n \rightarrow \infty} -1/n * \log_2 m(w_1, \dots, w_n)$$

Il modello più accurato sarà quello con la cross entropy minore.

Un linguaggio è *stazionario* se la probabilità che assegna a sequenze di parole sono invarianti rispetto al tempo.

Un linguaggio è *ergodico* se, aumentando la lunghezza della sequenza di parole generate, possiamo ottenere un campione perfettamente rappresentativo del processo.

In realtà, però, il linguaggio non è né stazionario né ergodico.

Natural Language Processing

Il *Natural Language Processing* è un sistema in grado di accedere al contenuto di informazioni attraverso l'elaborazione del linguaggio.

Vengono effettuati delle task tipiche della preparazione per l'analisi statistica.

Esistono diversi problemi che possono presentarsi nell'addestramento degli algoritmi per via delle ambiguità fonologiche, morfologiche, sintattiche, semantiche, etc.

Il *Machine Learning* comprende gli algoritmi che permettono alla macchina di imparare a svolgere un compito X , partendo da degli esempi su come svolgere quel determinato compito.

Linguistica Computazionale

Machine Learning

Usa modelli statistici dei dati nel corpus al fine di costruire un modello per svolgere il lavoro. I componenti sono:

- training corpus
- metodologia
- testing corpus

Se si aggiunge il ML al NLP, viene migliorato drasticamente il suo funzionamento, ma rimane legato ai dati di partenza. Per l'addestramento esistono due tipi di algoritmi:

Non supervisionato: vengono utilizzati dei corpus non annotati per creare modelli. Usati per compiti come il ranking dei dati in base a qualche funzione o il clustering in base a similitudini (*raw corpora*)

Supervisionato: si basa su un training "annotato a mano" tramite linguaggi di mark-up (XML, ...). Sono adatti per le classificazioni come, per esempio, trovare le corrette classi di appartenenza (*machine readable*)

L'annotazione BIOES-style viene utilizzata per espressioni polirematiche, espressioni che sono composte da più entità, ma valgono come un'unica ("carta di credito", "scala mobile").

L'acronimo BIO sta per: **B**egin, **I**n (parte interna), **O**ther

L'annotazione può essere fatta:

- manualmente, ma è molto lenta e costosa, a volte incoerente
- semi-automatica, tramite *parsers* (corretti al 96%) o *taggers* (corretti al 98%)

L'annotazione segue un processo di sviluppo di annotazione denominato M.A.T.T.E.R.:

- **M**odel, descrizione del fenomeno linguistico
- **A**nnote, annotazione in base alle features
- **T**rainig, addestramento di un algoritmo su un corpus
- **T**esting, viene testato l'algoritmo su un altro corpus
- **E**valuate, viene valutato l'esito del testing
- **R**evise, viene eventualmente revisionato l'algoritmo

I testi possono essere annotati morfo-sintatticamente tramite il PoS (*Part of Speech*) tagging. Il procedimento del PoS tagging genera delle coppie di valori dove vengono associate le parole e le parti del discorso corrispondenti, ad esempio: <cani, Nome> <mangia, Verbo>

Le annotazioni sintattiche vengono svolte tramite i *TreeBanks*, i quali danno:

- *bracketing*, strutture a costituenti
- *relazioni grammaticali*
- *struttura predicativa*

Gli schemi di annotazione possono essere:

- *labelled bracketing* (es: [S[NP Gianni][V legge] ...])
- *dependency structure*

Le annotazioni semantiche vengono svolte tramite la *Named Entity Recognition*, basata sulle categorie semantiche, la quale assegna ai nomi propri le "entità" della loro categoria. Ad esempio: ("Luigi", "Person") ("Università di Pisa", "Organization")

Linguistica Computazionale

Annotazione

I dati annotati “a mano” devono sempre essere valutati per determinare il grado di affidabilità. Per essere tale, l’annotazione deve essere replicata in maniera coerente da più annotatori, anche in momenti differenti.

Il grado di accordo si basa sui dati di *almeno* due annotatori e viene chiamato *interannotator agreement*. Esistono diverse misure di accordo, il più comune è chiamato *k di Cohen*, che si basa su soli due annotatori. Per 3+ annotatori avremo il *multi-k*, etc.

Il *k di Cohen* misura l’accordo al netto dell’accordo che si verificherebbe casualmente. Per poterlo calcolare si crea la cosiddetta *tabella di confusione* e, sulla diagonale alto-sx basso-dx, troveremo le volte che A_1 e A_2 sono concordi.

il *PoS tagging* è usato per la disambiguazione morfosintattica. Il compito è quello di assegnare, ad ogni token, la categoria grammaticale adatta. Esso si complica nella specifica dei tratti morfologici

Può utilizzare la regola *Pattern-Action*:

- *<action>*: seleziona e/o rimuove i tag
- *IF*: blocca *<action>* fino alla conferma
- *<pattern>*: se verificato, esegue *<action>*

Nella fase di addestramento troviamo:

- *training set*, un corpus di testi annotati
- *feature*, le caratteristiche estratte dall’input
- *language model*: *<feature, peso>*

Negli algoritmi supervisionati, il numero di classi y di un certo input x ($\langle x; y \rangle$, coppia del *training set*) deve essere finito, mentre gli input non sono quasi mai finiti. L’algoritmo deve quindi imparare le generalizzazioni di eventi mai osservati.

Le features sono le caratteristiche che devono essere analizzate e sono indicate come coppia *<attr., valore>*, dove il valore è un numero binario 0 (assente) o 1 (presente).

L’estrazione delle features è uno dei compiti più complessi, poiché potrebbe falsificare il modello che viene generato.

Le features si distinguono in:

- *locali*, estratte dal token stesso (*forma, lemma, ecc*)
- *contestuali*, estratte dal contesto in cui si trova il token (analisi della sua “storia”)
- *globali*, più ampie di quelle contestuali ma meno usate (dominio del documento)

Tabella di confusione, esempio:

A2 \ A1	X	Y	
	C1	C2	
X	50	20	R1
Y	10	20	R2

Si calcola come: $k = (P_A - P_E) / (1 - P_E)$

- P_A è il totale delle volte che A_1 e A_2 concordano, diviso il numero dei casi totali ($P_A = (O_{11} + O_{22}) / N$)
- P_E è il totale delle volte che vengono fatte le scelte per caso ($P_E = P_X + P_Y$ dove $P_X = P_{R1} * P_{C1}$ e $P_Y = P_{R2} * P_{C2}$)

Per un buon accordo, si dovrebbe ottenere $k \geq 0.80$. Per valori inferiori, il sistema potrebbe necessitare una revisione.

Linguistica Computazionale

Gold Standard

Il Gold Standard (GS) è una porzione del test corpus annotato manualmente, che rappresenta l'output di riferimento.

Basandoci sugli output attestati nel GS, abbiamo:

- *True Positive* (TP), sono gli output confermati con il GS
- *False Positive* (FP), sono gli output errati nel GS
- *False Negative* (FN), sono gli output mancati ed errati ($FN = |N| - TP$)

La *N-fold crossvalidation* è la metodologia di valutazione ideale per quando i dati annotati sono limitati:

1. Si divide il corpus in N parti
2. Si compiono N cicli di *training-evaluation* usando a ciclo una parte come test e le n-1 come training
3. Si effettua la media delle prestazioni

La *baseline* è il limite inferiore che ci si attende dal sistema. L'*accuracy* è invece data dal rapporto degli output corretti e la lunghezza del corpus.

Information Retrieval: ci dà delle misure riguardanti:

- *precisione*, la correttezza delle risposte del sistema
- *richiamo*, la copertura del sistema

Dai valori ottenuti con il GS possiamo calcolare tre parametri:

- *Precision* (P) = $TP / (TP+FP)$
- *Recall* (R) = $TP / (TP+FN)$
- *F-Measure* = $2PR/(P+R)$

Precision e Recall sono due valori che si penalizzano a vicenda: crescendo uno, diminuisce l'altro

Natural Language Processing

È il processo di analisi di una frase per determinare la sua struttura sintattica. Possono essere a costituenti o a dipendenze.

Vengono individuate le relazioni sintattiche tra i token della frase. Producono alberi a dipendenze, dove i *nodi* sono i token e la *radice*, mentre gli *archi* sono le relazioni tra le parole. Si usano classificatori addestrati su delle *TreeBanks* allo scopo di predire l'azione in base alle features.

L'algoritmo termina dopo aver analizzato tutto e comincia l'analisi in base alle features estratte.

Le azioni effettuate per costruire:

- *shift*: non c'è relazione tra i due token, va avanti
- *right*: c'è relazione, il nodo sx dipende dal nodo dx
- *left*: c'è relazione, il nodo dx dipende dal nodo sx

Linguistica Computazionale

