

Linguistica computazionale

(A/A 2022-23)

Prova in itinere del 8 novembre 2022

Testo A

(studenti la cui PENULTIMA cifra del numero di matricola è compresa tra 0 e 4)

Cognome e nome:

Matricola:

Istruzioni per la consegna:

- ridenominare il file del compito con *CognomeNomeMatricola* (es. RossiMario001293.doc)
- inviare il file a: alessandro.lenci@unipi.it

1. Illustrate la Legge di Zipf (riportando anche la formula matematica) e le sue conseguenze per l'analisi computazionale del linguaggio (PUNTI: 4) **(la risposta deve essere lunga min. 5, max 10 righe)**

2. Spiegate la differenza tra ASCII, Unicode e Iso-Latin1 (PUNTI: 4) **(la risposta deve essere lunga min. 5, max 10 righe)**

3. Spiegate cosa sono le concordanze di una parola (PUNTI: 4) (la risposta deve essere lunga min. 5, max 10 righe)

4. In un corpus lungo 15.000 token, il bigramma <centro, sociale > ricorre 80 volte. a). Costruite la tabella di contingenza del bigramma e calcolate la sua frequenza attesa; b.) calcolate la mutua informazione del bigramma sapendo che $P(\text{sociale}) = 0.02$, e $F(\text{centro}) = 120$ (NB: F = frequenza; P = probabilità) (PUNTI: 4).

	Y=sociale	Y!=sociale
X=centro	80	40
X!=centro	220	14660

Frequenza attesa= $(120 * 300) / 15000 = 2.4$

MI(centro, sociale) = $\log_2 \left(\frac{80 * 15000}{300 * 120} \right) = 5.058$

5. $|C| = 3000$ $|V| = 375$

$F(\text{lo}) = 50$, $F(\text{un}) = 35$, $F(\text{di}) = 65$, $F(\text{studente}) = 30$, $F(\text{linguistica}) = 22$, $F(\text{legge}) = 10$, $F(\text{libro}) = 14$, $F(\text{interessante}) = 6$

$F(\langle \text{lo}, \text{studente} \rangle) = 15$, $F(\langle \text{studente}, \text{di} \rangle) = 25$, $F(\langle \text{linguistica}, \text{legge} \rangle) = 5$, $F(\langle \text{legge}, \text{un} \rangle) = 9$, $F(\langle \text{un}, \text{libro} \rangle) = 10$

$F(\langle \text{lo}, \text{studente}, \text{di} \rangle) = 7$, $F(\langle \text{linguistica}, \text{legge}, \text{un} \rangle) = 4$, $F(\langle \text{legge}, \text{un}, \text{libro} \rangle) = 8$

Usando un modello markoviano del I ordine e lo add-one smoothing, calcolate la probabilità della frase *Lo studente di linguistica legge un libro interessante*. (PUNTI: 5).

$P(\text{lo}, \text{studente}, \text{di}, \text{linguistica}, \text{legge}, \text{un}, \text{libro}, \text{interessante}) = P(\text{lo}) * P(\text{studente}|\text{lo}) * P(\text{di}|\text{studente}) * P(\text{linguistica}|\text{di}) * P(\text{legge}|\text{linguistica}) * P(\text{un}|\text{legge}) * P(\text{libro}|\text{un}) * P(\text{interessante}|\text{libro}) =$

$$(51/(3000+375)) * (16/(50+375)) * (26/(30+375)) * (1/(65+375)) * (6/(22+375)) * (10/(10+375)) * (11/(35+375)) * (1/(14+375)) =$$

$$(51/3375) * (16/425) * (26/405) * (1/440) * (6/397) * (10/385) * (11/410) * (1/389) =$$

$$0.015 * 0.038 * 0.064 * 0.002 * 0.015 * 0.025 * 0.027 * 0.002 = 1.47744e-15$$

6.

In una caverna sotto terra viveva uno hobbit. Non era una caverna brutta, sporca, umida, piena di resti di vermi e di trasudo di fetido, e neanche una caverna arida, spoglia, sabbiosa, con dentro niente per sedersi o da mangiare: era una caverna hobbit, cioè comodissima.

Tokenizzate questo testo (normalizzando le maiuscole e includendo la punteggiatura) mettendo un token per riga; calcolate la frequenza cumulata relativa di parole token, calcolate la frequenza relativa della parola *di*, l'indice di ricchezza lessicale e costruite lo spettro di frequenza (PUNTI: 5).

in
una
caverna
sotto
terra
viveva
uno
hobbit
.
Non
era
una
caverna
brutta
,
sporca
,
umida
,
piena
di
resti
di
vermi
e
di
trasudo
di
fetido
,
e
neanche

una
caverna
arida
,
spoglia
,
sabbiosa
,
con
dentro
niente
per
sedersi
o
da
mangiare
:
era
una
caverna
hobbit
,
cioè
comodissima
.

$$\text{TTR} = 37/57 = 0.649$$

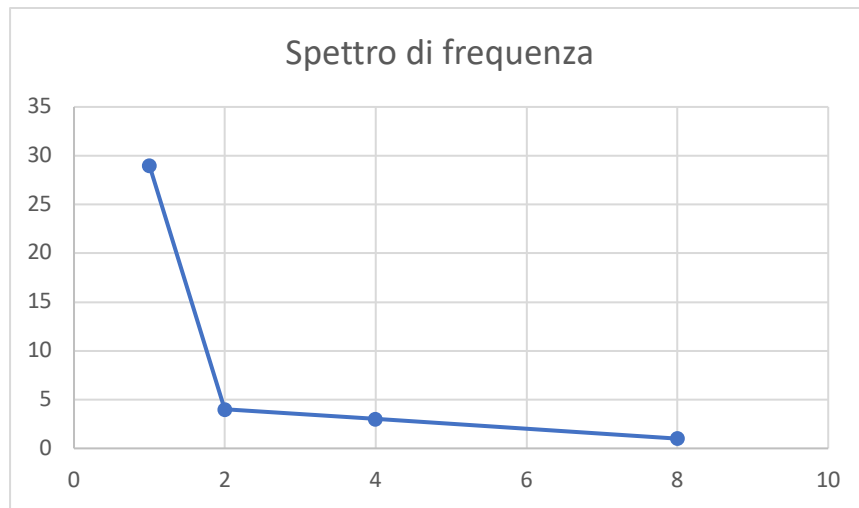
$$\text{Frequenza relativa di } di: 4/57 = 0.07$$

Classi di Frequenza:

V1	29
V2	4
V4	3
V8	1

Frequenza cumulata relativa di parole token:

$(V1*1)/57$	$= 0.50877193$
$((V1*1)+(V2*2))/57$	$= 0.649122807$
$((V1*1)+(V2*2)+(V4*3))/57$	$= 0.859649123$
$((V1*1)+(V2*2)+(V4*3)+(V8*8))/57$	$= 1$



7. Scrivete le espressioni regolari corrispondenti ai seguenti pattern (PUNTI: 4):

1) sostituite la consonante iniziale di una parola con *a*, soltanto se la parola non termina per *e*, ed è lunga almeno 15 caratteri;

`s/b([cdfghklmpqrstvz])(\w{13,})[^e]\b/a\2/`

2) inserite la stringa *ro* dopo la prima lettera delle parole che compaiono all'inizio di una riga di testo e terminano per *da* o *li*

`s/^(w)(w*)(da|li)\b\1ro\2\3/`

3) trovate le parole formate da un numero dispari di caratteri, che terminano per *r* oppure *l* e compaiono alla fine di una riga di testo;

`^b(w\w)*[rl]$/`