

# Linguistica computazionale

(A/A 2023-24)

Prova in itinere del 20 dicembre 2023

## Testo A

(studenti la cui ultima 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](mailto:alessandro.lenci@unipi.it)

---

1. Spiegate la nozione di n-fold crossvalidation (PUNTI: 4) (**la risposta deve essere lunga min. 5, max 10 righe**)

2. Illustrate il task di Named Entity Recognition e lo schema di annotazione usato (PUNTI: 4) (**la risposta deve essere lunga min. 5, max 10 righe**)

3. Illustrate come si valuta l'affidabilità di un'annotazione linguistica manuale (PUNTI: 4) (la risposta deve essere lunga min. 5, max 10 righe)

4. In un test corpus ci sono 820 nomi concreti (classe C) e 740 nomi astratti (classe A). Un sistema di annotazione semantica assegna la classe C a 900 nomi, di cui 400 sono corretti, mentre assegna la classe A a 660 nomi, di cui 350 sono corretti. Costruite la matrice di confusione e calcolate precision, recall e F-measure del sistema per entrambe le classi (PUNTI: 4):

$$TP\_C = 400$$

$$TP\_A = 350$$

$$FP\_C = 900 - 400 = 500$$

$$FP\_A = 660 - 350 = 310$$

$$FN\_C = 820 - 400 = 420$$

$$FN\_A = 740 - 350 = 390$$

$$P\_C = 400/(400+500) = 0.44$$

$$P\_A = 350/(350+310) = 0.53$$

$$R\_C = 400/(400+420) = 0.49$$

$$R\_A = 350/(350+390) = 0.47$$

$$F1\_C = 2*((0.44*0.49)/(0.44+0.49)) = 0.46$$

$$F1\_A = 2*((0.53*0.47)/(0.53+0.47)) = 0.50$$

5. In un corpus di 10.000 parole annotato con le Named Entity, i.) ci sono 1.500 nomi LOC e 1.800 nomi ORG; ii.) 1.700 ORG sono preceduti da un articolo, mentre solo 300 nomi LOC sono preceduti da un articolo. Calcolate qual è la classe più probabile da assegnare al nome *Roma* nella frase *Ho visto la Roma giocare*. (PUNTI: 5).

$$P(\text{LOC}|la Roma) = P(\text{LOC}) * P(\text{la Roma}|\text{LOC}) = 1500/10000 * 300/1500 = 0.15 * 0.2 = 0.03$$

$$P(\text{ORG}|\text{la Roma}) = P(\text{ORG}) * P(\text{la Roma}|\text{ORG}) = 1800/10000 * 1700/1800 = 0.18 * 0.94 = 0.17$$

La classe più probabile è ORG

6. Dato il contesto “C = *La macchina insegue la ...*”, un language model genera la parola seguente con questa distribuzione di probabilità  $m$ :  $P(\text{talpa}|\text{C}) = 0.05$ ;  $P(\text{dolce}|\text{C}) = 0.01$ ;  $P(\text{carta}|\text{C}) = 0.03$ ;  $P(\text{moto}|\text{C}) = 0.70$ ;  $P(\text{gatta}|\text{C}) = 0.20$ ;  $P(\text{felpa}|\text{C}) = 0.01$ . Calcolate l'entropia del language model e, assumendo  $m$  come la distribuzione reale, calcolate la sua cross-entropy rispetto a un modello che assegna a tutte le parole la stessa probabilità (PUNTI: 5).

$$H(\text{LM}) = - (P(\text{talpa}|\text{C}) * \log_2 P(\text{talpa}|\text{C}) + P(\text{dolce}|\text{C}) * \log_2 P(\text{dolce}|\text{C}) + P(\text{carta}|\text{C}) * \log_2 P(\text{carta}|\text{C}) + P(\text{moto}|\text{C}) * \log_2 P(\text{moto}|\text{C}) + P(\text{gatta}|\text{C}) * \log_2 P(\text{gatta}|\text{C}) + P(\text{felpa}|\text{C}) * \log_2 P(\text{felpa}|\text{C})) = -(0.05 * \log_2(0.05) + 0.01 * \log_2(0.01) + 0.03 * \log_2(0.03) + 0.70 * \log_2(0.70) + 0.20 * \log_2(0.20) + 0.01 * \log_2(0.01)) = 1.32$$

Nel modello equiprobabile, ogni parola riceve una probabilità  $1/6 = 0.17$

$$\text{Cross\_entropy}(\text{LM}, \text{equiProb}) = - (0.05 * \log_2(0.17) + 0.01 * \log_2(0.17) + 0.03 * \log_2(0.17) + 0.70 * \log_2(0.17) + 0.20 * \log_2(0.17) + 0.01 * \log_2(0.17)) = 2.56$$

7. Due language model A e B assegnano a un corpus lungo 1.000 token le seguenti probabilità:  $P_A = 0.045$ ,  $P_B = 0.0023$ . Misurate la bontà dei modelli e indicate qual è il modello migliore (PUNTI: 4):

Si misura la bontà del modello con la perplexity:

$$\text{Perplexity}(A) = 0.045^{-(1/1000)} = 1.003$$

$$\text{Perplexity}(B) = 0.0023^{-(1/1000)} = 1.006$$

Il modello migliore è A, perché ha la perplexity minore