

17/4

N-gram	count
your	883,614
rights	80,891
doorposts	21
your rights	378
your doorposts	0

corpus $520 \cdot 10^6$ words = Na) stima con MLE di $P(\text{your}) = \frac{883,614}{520 \cdot 10^6}$

$$Pr(\text{right} | \text{your}) = \frac{378}{883,614} \left(= \frac{C(w_i, w_{i+1})}{C(w_i, *)} \right)$$

b) stima $Pr(\text{doorposts} | \text{your})$; MLE con add- κ smoothing per $\kappa = 0.01$; $|V| = 1,254,193$

$$Pr = \frac{0 + 0.01}{883,614 + |V| \cdot \kappa} \Rightarrow C(\text{doorposts}, \text{your}) \text{ aumenta di } \kappa$$

\downarrow
 $C(w)$ aumenta di $\kappa \forall w \in \text{corpus}$

Argomento extra (non in syllabus)

SINTASSI: regole che governano struttura di frasi \rightarrow determinata da mente

CONSTITUENT/PHRASE: gruppo di parole che fungono da unità in struttura gerarchica

(es. "He saw the house on the hill" \rightarrow "He", "the house", "on the hill")Per identificarli, CONSTITUENCY TESTS \rightarrow più importanti: noun phrase, verb phraseParte fondamentale di NLP: AMBIGUITY \rightarrow ambiguità in syntactic parsing che influenza interpretazione: problema di PP ATTACHMENT

Si può rappresentare struttura sintattica come albero

