

NATURAL LANGUAGE PROCESSING
Guidelines for Assignment 2

Note: do not include $\langle s \rangle$ or $\langle /s \rangle$ in the vocabulary (V) or when counting tokens (N).

Laplace Unigram

$$P(\langle s \rangle \text{ I love Soho } \langle /s \rangle) = P(\text{I}) * P(\text{love}) * P(\text{Soho})$$

$$P(\text{I}) = \frac{\text{count}(\text{I}) + 1}{N + V}$$

Laplace Bigram

$$P(\langle s \rangle \text{ I love Soho } \langle /s \rangle) = P(\text{I}/\langle s \rangle) * P(\text{love}/\text{I}) * P(\text{Soho}/\text{love}) * P(\langle /s \rangle/\text{Soho})$$

$$P(\text{I}/\langle s \rangle) = \frac{\text{count}(\langle s \rangle, \text{I}) + 1}{\text{count}(\langle s \rangle) + V}$$

Stupid Backoff

$$P(\langle s \rangle \text{ I love Soho } \langle /s \rangle) = P(\text{I}/\langle s \rangle) * P(\text{love}/\text{I}) * P(\text{Soho}/\text{love}) * P(\langle /s \rangle/\text{Soho})$$

$$P(\text{I}/\langle s \rangle) = \begin{cases} \frac{\text{count}(\langle s \rangle, \text{I})}{\text{count}(\langle s \rangle)} & \text{if } \text{count}(\langle s \rangle, \text{I}) > 0 \\ 0.4 * \frac{\text{count}(\text{I}) + 1}{N + V} & \text{otherwise} \end{cases}$$

Kneser-Ney

$$P(\langle s \rangle \text{ I love Soho } \langle /s \rangle) = P(\text{I}/\langle s \rangle) * P(\text{love}/\text{I}) * P(\text{Soho}/\text{love}) * P(\langle /s \rangle/\text{Soho})$$

$$P(\text{I}/\langle s \rangle) = \begin{cases} \frac{\text{numerator}}{\text{count}(\langle s \rangle)} & \text{if } \text{numerator} > 0 \\ \frac{\text{count}(\text{I}) + 1}{N + V} & \text{otherwise} \end{cases}$$

$$\text{numerator} = \max(\text{count}(\langle s \rangle, \text{I}) - d, 0) + d * n\text{Next}(\langle s \rangle) * n\text{Prev}(\text{I}) / n\text{bigrams}$$

where $n\text{bigrams}$ is the total number of bigrams, $n\text{Next}(w)$ is the total number of continuations for w and $n\text{Prev}(w)$ is the total number of word types seen to precede w .