

N-gram Language Modeling

$$P(\bar{w}) = P(w_1, \dots, w_m) = P(w_1) P(w_2 | w_1) P(w_3 | w_1 w_2) \dots$$

n-gram LM:

$$P(\bar{w}) = \prod_{i=1}^m P(w_i | \underbrace{w_{i-n+1}, \dots, w_{i-1}}_{\text{previous } n-1 \text{ words}}) \quad n \approx 3-7$$

2-gram LM: $P(w_1 | \langle s \rangle) P(w_2 | w_1) P(w_3 | w_2) \dots$

w_3 cond. indep. of w_1 | w_2

n-gram LM \Leftrightarrow n-1-order Markov model

3-gram LM: $P(w_1 | \langle s \rangle \langle s \rangle) P(w_2 | \langle s \rangle w_1) P(w_3 | w_1 w_2) \dots$

2-gram: multinomial distributions, $|V| \times |V|$ params

$$P(\text{w}|\text{the}) = \begin{array}{l} 0.001 \text{ house} \\ 0.0005 \text{ dog} \\ 0.0005 \text{ cat} \\ \vdots \end{array} \quad \text{very flat distribution!}$$

Parameter estimation MLE from a large corpus

$$P(\text{dog}|\text{the}) = \frac{\text{count}(\text{the}, \text{dog})}{\text{count}(\text{the})}$$

- Why
- ① Generation: machine translation
 - ② Grammatical error correction
 - ③ Way to build "word2vec++"