

n-gram models

- The Markov assumption: $P(\mathbf{x}^{(t)} | \mathbf{x}^{(1)}, \dots, \mathbf{x}^{(t-1)}) = P(\mathbf{x}^{(t)} | \mathbf{x}^{(t-n)}, \dots, \mathbf{x}^{(t-1)})$
- *n*-gram models are language models that use the Markov assumption with a specific selection of *n*.

$P(dog | the, quick, brown, fox, jumped, over, the, lazy) = P(dog) \quad \leftarrow \quad (n = 1) \quad \text{Unigram}$

$P(dog | the, quick, brown, fox, jumped, over, the, lazy) = P(dog | lazy) \quad \leftarrow \quad (n = 2) \quad \text{Bigram}$

$P(dog | the, quick, brown, fox, jumped, over, the, lazy) = P(dog | the, lazy) \quad \leftarrow \quad (n = 3) \quad \text{Trigram}$

$P(dog | the, quick, brown, fox, jumped, over, the, lazy) = P(dog | over, the, lazy) \quad \leftarrow \quad (n = 4) \quad \text{4-gram}$

- This should look familiar to Lecture-03 (Word2Vec)! Skip-gram modeling generalizes *n*-grams by not considering the sequence order of the *context* words.

Estimating n -gram probabilities

- The maximum likelihood estimate of an n -gram model:

$$P(\mathbf{x}^{(t)} \mid \mathbf{x}^{(t-n)}, \dots, \mathbf{x}^{(t-1)}) \stackrel{MLE}{=} \frac{\text{count}(\mathbf{x}^{(t-n)}, \dots, \mathbf{x}^{(t)})}{\sum_{j=1}^N \text{count}(\mathbf{x}^{(t-n)}, \dots, \mathbf{x}_j^{(t)})}$$

Example from Jurafsky & Martin

<s> I am Sam </s>

<s> Sam I am </s>

<s> I do not like green eggs and ham </s>

Here are the calculations for some of the bigram probabilities from this corpus

$$P(\text{I} \mid \text{<s>}) = \frac{2}{3} = .67 \quad P(\text{Sam} \mid \text{<s>}) = \frac{1}{3} = .33 \quad P(\text{am} \mid \text{I}) = \frac{2}{3} = .67$$

$$P(\text{</s>} \mid \text{Sam}) = \frac{1}{2} = 0.5 \quad P(\text{Sam} \mid \text{am}) = \frac{1}{2} = .5 \quad P(\text{do} \mid \text{I}) = \frac{1}{3} = .33$$