

The Markov approximation

- From previous slide, notice that can estimate **some** of the factors of the joint distribution:

$$P(\text{sentence}) = P(\mathbf{x}^{(1)}, \dots, \mathbf{x}^{(T)}) = \underbrace{P(\mathbf{x}^{(1)})}_{\text{yes}} \cdot \underbrace{P(\mathbf{x}^{(2)} | \mathbf{x}^{(1)})}_{\text{yes}} \cdot \underbrace{P(\mathbf{x}^{(3)} | \mathbf{x}^{(1)}, \mathbf{x}^{(2)})}_{\text{yes}} \dots \underbrace{P(\mathbf{x}^{(T)} | \mathbf{x}^{(1)}, \dots, \mathbf{x}^{(T-1)})}_{\text{No!}}$$

- The Markov assumption: $P(\mathbf{x}^{(t)} | \mathbf{x}^{(1)}, \dots, \mathbf{x}^{(t-1)}) = P(\mathbf{x}^{(t)} | \mathbf{x}^{(t-n)}, \dots, \mathbf{x}^{(t-1)})$

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

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

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

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

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.