

Can't we just use the chain rule?

- Can't we just use the chain rule of probability to factor the joint distribution into a distribution over its suffix conditioned on its prefix?

$$P(\text{sentence}) = P(\mathbf{x}^{(1)}, \dots, \mathbf{x}^{(T)}) = P(\mathbf{x}^{(1)}) \cdot P(\mathbf{x}^{(2)} | \mathbf{x}^{(1)}) \cdot P(\mathbf{x}^{(3)} | \mathbf{x}^{(1)}, \mathbf{x}^{(2)}) \dots P(\mathbf{x}^{(T)} | \mathbf{x}^{(1)}, \dots, \mathbf{x}^{(T-1)})$$

yes

yes

yes

No!

- No because we end up with (essentially) the same problem, now it's too many possible sequences over $T - 1$ tokens.
- We need to make a simplifying assumption regarding the independence of words in a sentence.

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

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

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

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