



Why?

$$P(W) = P(w_1, w_2, w_3, w_4, w_5 \dots w_n)$$

lated task: probability of an upcoming word:
 $P(w_5 | w_1, w_2, w_3, w_4)$

model that computes either of these:

$P(W)$ or $P(w_n | w_1, w_2 \dots w_{n-1})$ is called a **language model**.

- Better: **the grammar** But **language model** or **LM** is standard

- How to compute this joint probability:
 - $P(\text{its, water, is, so, transparent, that})$
- Intuition: let's rely on the Chain Rule of Probability

- Recall the definition of conditional probabilities

$$P(\mathbf{B} | \mathbf{A}) = P(\mathbf{A}, \mathbf{B}) / P(\mathbf{A}) \quad \text{Rewriting: } P(\mathbf{A}, \mathbf{B}) = P(\mathbf{A})P(\mathbf{B} | \mathbf{A})$$
- More variables:

$$P(\mathbf{A}, \mathbf{B}, \mathbf{C}, \mathbf{D}) = P(\mathbf{A})P(\mathbf{B} | \mathbf{A})P(\mathbf{C} | \mathbf{A}, \mathbf{B})P(\mathbf{D} | \mathbf{A}, \mathbf{B}, \mathbf{C})$$
- The Chain Rule in General

$$P(x_1, x_2, x_3, \dots, x_n) = P(x_1)P(x_2 | x_1)P(x_3 | x_1, x_2) \dots P(x_n | x_1, \dots, x_{n-1})$$

$$p(B | A) = P(A,B)/P(A) \quad \text{Rewriting: } P(A,B) = P(A)P(B | A)$$

- More variables:

$$P(A,B,C,D) = P(A)P(B|A)P(C|A,B)P(D|A,B,C)$$

- The Chain Rule in General

$$P(x_1, x_2, x_3, \dots, x_n) = P(x_1)P(x_2|x_1)P(x_3|x_1, x_2) \dots P(x_n|x_1, \dots, x_{n-1})$$
$$P(\text{its}) \times P(\text{water}|\text{its}) \times P(\text{is}|\text{its water}) \\ \times P(\text{so}|\text{its water is}) \times P(\text{transparent}|\text{its water is so})$$



How to estimate these probabilities

- Could we just count and divide?

$$P(\text{the l its water is so transparent that}) = \frac{\text{Count}(\text{its water is so transparent that the})}{\text{Count}(\text{its water is so transparent that})}$$

- No! Too many possible sentences!
- We'll never see enough data for estimating these



Markov Assumption



- Simplifying assumption:

$$P(\text{the l its water is so transparent that}) \approx P(\text{the l that})$$

- Or maybe

$$P(\text{the l its water is so transparent that}) \approx P(\text{the l transparent that})$$



Markov Assumption

$$P(w_1 w_2 \dots w_n) \approx \prod_i P(w_i | w_{i-k} \dots w_{i-1})$$

- In other words, we approximate each component in the product

$$P(w_i | w_1 w_2 \dots w_{i-1}) \approx P(w_i | w_{i-k} \dots w_{i-1})$$



Simplest case: Unigram model

$$P(w_1 w_2 \dots w_n) \approx \prod_i P(w_i)$$

Some automatically generated sentences from a unigram model

fifth, an, of, futures, the, an, incorporated, a,
a, the, inflation, most, dollars, quarter, in, is,
mass

thrift, did, eighty, said, hard, 'm, july, bullish

that, or, limited, the



Bigram model

- Condition on the previous word:

$$P(w_i | w_1 w_2 \dots w_{i-1}) \approx P(w_i | w_{i-1})$$

texaco, rose, one, in, this, issue, is, pursuing, growth, in,
a, boiler, house, said, mr., gurria, mexico, 's, motion,
control, proposal, without, permission, from, five, hundred,
fifty, five, yen

outside, new, car, parking, lot, of, the, agreement, reached

this, would, be, a, record, november



N-gram models

- We can extend to trigrams, 4-grams, 5-grams
- In general this is an insufficient model of language
 - because language has **long-distance dependencies**:

"The computer(s) which I had just put into the machine room
on the fifth floor is (are) crashing."

- But we can often get away with N-gram models