

## How to compute the probabilities?

Language Models are defined by the following objective :

1. Objective - compute the probability of a sentence or sequence of words.  $P(W) = P(w_1, w_2, w_3, \dots, w_n)$
2. Related task - computing the probability of the upcoming word.  $P(w_4 | w_1, w_2, w_3)$

**Probability of the entire sentence :  $P(W) = P(w_1, w_2, w_3 \dots w_n)$**  = Probability when all the words are happening together in sequence.

$P(\text{The, water, of, Walden, Pond, is, so, beautifully, blue})$  : The probability of when the sentence "The water of Walden Pond is so beautifully blue" will occur in a corpus.

There is a difference between the following -

1. **Conditional Probability  $P(B|A)$**  : Probability of B given A : Event B has happened in past. A is happening now. This is known as Conditional Probability.
2. **Joint Probability  $P(A,B)$**  : Joint probability of A and B. Or probability when both the events A and B are happening simultaneously.

$$P(B|A) = P(A \cap B) / P(A) \text{ or, } P(B|A) = P(A,B) / P(A) \text{ or, } P(A,B) = P(B|A) \times P(A)$$

$$\text{or, } P(A,B) = P(A) \times P(B|A)$$

ie. probability of two events A and B happening together (joint probability) is probability of A multiplied by probability of B when A has already happened.

### 3. Extending it to multiple events we can write

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

To get the intuition, following is the chain of thought -

- first the event A happened. So, the probability is  $P(A)$  as nothing else has happened now.
- second the event B happened. Event A has already happened in the last step. so the probability of B, we need to compute  $P(B|A)$ , i.e., probability of B when A has already happened.
- now, the third event C happened. A and B has already happened. So, probability of C would be,  $P(C|A,B)$ , because A and B has already happened. So,
- now, the fourth event D happened. A, B and C has already happened by now. So, probability of D when A, B and C has already happened is  $P(D|A,B,C)$

$$P(\text{blue} | \text{the water of walden pond is so beautifully}) = P(\text{the water of walden pond is so beautifully blue}) / P(\text{the water of walden pond is so beautifully})$$

### 5. Generalizing the above

$$P(w_{1:n}) = P(w_1) \times P(w_2 | w_1) \times P(w_3 | w_{1:2}) \dots P(w_n | w_{1:n-1}) = \prod_{k=1}^n P(w_k | w_{1:k-1})$$

Eg :  $P(\text{"The \ water \ of \ walden \ pond \ is \ so \ beautifully \ blue"}) = P(\text{The}) \times P(\text{water} | \text{The}) \times P(\text{of} | \text{The \ water}) \times P(\text{walden} | \text{The \ water \ of}) \times P(\text{pond} | \text{The \ water \ of \ walden}) \times P(\text{is} | \text{The \ water \ of \ walden \ pond}) \times P(\text{so} | \text{The \ water \ of \ walden \ pond \ is}) \times P(\text{beautifully} | \text{The \ water \ of \ walden \ pond \ is \ so}) \times P(\text{blue} | \text{The \ water \ of \ walden \ pond \ is \ so \ beautifully})$

water \ of \ walden \ pond) \times P(so|The \ water \ of \ walden \ pond \ is) \times P(beautifully|The \ water \ of \ walden \ pond \ is \ so) \times P(blue|The \ water \ of \ walden \ pond \ is \ so \ beautifully)\$

1. We will never see enough data for estimating all these probabilities. Hence **Markov assumption** is used to simplify the matter. It states the following -

#### Markov Assumption

$P(\text{blue} \mid \text{The water of Walden Pond is so beautifully}) \approx P(\text{blue} \mid \text{beautifully})$   $P(w_n|w_{1:n-1}) \approx P(w_n|w_{n-1})$   
 The approximation is known as *bi-gram assumption* or *1st order markov assumption*. It is considering the probability of  $(n-1)^{\text{th}}$  word, instead of probability of  $n-1$  words, to determine the Probability of  $n^{\text{th}}$  word.

$k = 2$  : bi-gram model : probability of the current word depends on the previous 2-1 words.  $k = n$  : n-gram model : probability of the current word depends on the previous  $n-1$  words.  $k = k$  : k-gram model : probability of the current word depends on the previous  $k-1$  words.

$$P(w_{1:n}) \approx \prod_{k=1}^n P(w_k|w_{k-1})$$

Going more generally, **n-gram** would consider  $n-1$  words prior to the  $n^{\text{th}}$  word to determine the probability of the  $n^{\text{th}}$  word. bi-gram was considering 2-1 words prior.

Therefore bi-gram model, the probability of the  $n^{\text{th}}$  word would be determined by the  $n-k^{\text{th}}$  word.

Therefore for an n-gram model,

$$P(w_{1:n}) \approx \prod_{k=1}^n P(w_k|w_{k-(n-1)})$$