

# N-grams and Probabilities

Before we start computing probabilities of certain sequences, we need to first define what is an N-gram language model:

An N-gram is a sequence of N words

Corpus: I am happy because I am learning

Unigrams: { I, am, happy, because, learning }

Bigrams: { I am, am happy, happy because ... }

✗ I happy

Trigrams: { I am happy, am happy because, ... }

Now given the those definitions, we can label a sentence as follows:

Corpus:  $w_1$   $w_2$   $w_3$  ...  $w_{498}$   $w_{499}$   $w_{500}$   $m = 500$

In other notation you can write:

- $w_1^m = w_1 w_2 w_3 \dots w_m$
- $w_1^3 = w_1 w_2 w_3$
- $w_{m-2}^m = w_{m-2} w_{m-1} w_m$

Given the following corpus: *I am happy because I am learning.*

- Size of corpus  $m = 7$ .
- $P(I) = \frac{2}{7}$
- $P(happy) = \frac{1}{7}$

To generalize, the probability of a unigram is  $P(w) = \frac{C(w)}{m}$

**Bigram Probability:**

Corpus: I am happy because I am learning

$$P(am|I) = \frac{C(I \text{ } am)}{C(I)} = \frac{2}{2} = 1$$

$$P(happy|I) = \frac{C(I \text{ } happy)}{C(I)} = \frac{0}{2} = 0 \quad \times \text{ I happy}$$

$$P(learning|am) = \frac{C(am \text{ } learning)}{C(am)} = \frac{1}{2}$$

Probability of a bigram: 
$$P(y|x) = \frac{C(x \text{ } y)}{\sum_w C(x \text{ } w)} = \frac{C(x \text{ } y)}{C(x)}$$

### Trigram Probability:

To compute the probability of a trigram:

- $P(w_3 | w_1^2) = \frac{C(w_1^2 w_3)}{C(w_1^2)}$
- $C(w_1^2 w_3) = C(w_1 w_2 w_3) = C(w_1^3)$

### N-gram Probability:

- $P(w_N | w_1^{N-1}) = \frac{C(w_1^{N-1} w_N)}{C(w_1^{N-1})}$
- $C(w_1^{N-1} w_N) = C(w_1^N)$