

# The Language Modeling problem

- Create a probability distribution over all sequences of words
  - finite vocabulary:  $\Sigma$
  - infinite set of sequences:  $\Sigma^*$ 
    - Any sentence/sequence of words  $e = w_1 w_2 \dots w_n$  is an element of  $\Sigma^*$

$$\sum_{e \in \Sigma^*} P_{LM}(e) = 1$$

$$P_{LM}(e) \geq 0 \quad \forall e \in \Sigma^*$$

## Our First Attempt

- Assume we have  $N$  training sentences
- Let  $w_1 w_2 \dots w_n$  be a sentence, and  $\text{count}(w_1, w_2, \dots, w_n)$  be the number of times it appeared in the training data.
- Define a language model:

$$P(w_1, \dots, w_n) = \frac{\text{count}(w_1, \dots, w_n)}{N}$$

## Unigram Language Model

*“I have a dog whose name is Lucy. I have two cats, they like playing with Lucy.”*

- $\sum_{w \in \Sigma} \text{count}(w) =$   $P(w_k | w_{1:k-1}) \approx P(w_k)$
- $P(\text{Lucy}) =$
- $P(\text{cats}) =$   $\hat{P}(w) = \frac{\text{count}(w)}{\sum_{v \in \Sigma} \text{count}(v)}$

## Bigram Language Model

*“I have a dog whose name is Lucy. I have two cats, they like playing with Lucy.”*

- $P(\text{have} | \text{I}) =$   $P(w_k | w_{1:k-1}) \approx P(w_k | w_{k-1})$
- $P(\text{two} | \text{have}) =$
- $P(\text{eating} | \text{have}) =$   $\hat{P}(w_2 | w_1) = \frac{\text{count}(w_1 w_2)}{\text{count}(w_1)}$