

## Logistics

↳ HW 0 due 9/17

↳ Form project group by 9/17

↳ 3-5 students

↳ project proposal due 10/6

↳ has to deal w/ language

---

## Language models

↳ a language model assigns a probability to a piece of text

$P(\text{"students opened their books"})$   
 $> P(\text{"students opened their doorknobs"})$

given a seq of words  $S = w_1, w_2, \dots, w_n$

a LM computes

$$P(S) = P(w_1, w_2, \dots, w_n) = P(w_{1:n})$$

What about "predict the next word"?

$$P(w_n | w_1, w_2, \dots, w_{n-1})$$

prefix

$$= P(w_n | w_{<n}) = P(w_n | w_1 \dots n-1)$$

remember the chain rule of probability

$$P(B|A) = \frac{P(A, B) \rightarrow \text{joint prob}}{P(A)}$$

cond. prob

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

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

$$P(S) = P(w_1, \dots, w_n)$$
$$= P(w_1) P(w_2 | w_1) P(w_3 | w_1, w_2), \dots$$

$$= \prod_i P(w_i | w_{<i})$$

let's say we are given a **training dataset** of documents to estimate these probabilities.  
let's just count and divide!

My favorite LLM is \_\_\_\_\_

$P(\text{Gemini} | \text{My fav. LLM is})$

$$\approx \frac{\text{Count}(\text{My favorite LLM is Gemini})}{\text{Count}(\text{My fav. LLM is})}$$

issues:

↳ Sparsity

↳ no sharing counts between seqs w/ similar meaning

↳ Storage

n-gram models:

↳ approximate these probabilities  
by dropping context from the prefix  
⇒ Markov assumption

$P(\text{Gemini} | \text{My fav LLM is})$

$$\approx P(\text{Gemini} | \text{is}) \Rightarrow \frac{\text{count}(\text{is Gemini})}{\text{count}(\text{is})} \Rightarrow \text{bigram 2-gram}$$

$$\approx P(\text{Gemini} | \text{LLM is}) \Rightarrow \frac{\text{count}(\text{LLM is Gemini})}{\text{count}(\text{LLM is})} \Rightarrow \text{trigram 3-gram}$$

how many counts do we have to store for a  $n$ -gram model?  $\Rightarrow V^n$  where  $V$  is the size of the vocab

$\langle s \rangle$  I am Sam  $\langle /s \rangle$

$\langle s \rangle$  Sam I am  $\langle /s \rangle$

$\langle s \rangle$  I do not like green eggs and ham  $\langle /s \rangle$

vocab:  $\{ \langle s \rangle, \langle /s \rangle, \text{I}, \text{am}, \text{Sam}, \text{do}, \text{not}, \text{like}, \text{green}, \text{eggs}, \text{and}, \text{ham} \}$

↳ each unique element of vocab is called a **type**

↳ instances of types in a dataset are called **tokens**

$$P(\text{I} | \langle s \rangle) = \frac{2}{3} \quad P(\text{ham} | \langle s \rangle) = 0$$

$$P(\text{Sam} | \langle s \rangle) = \frac{1}{3} \quad P(\text{Sam} | \text{am}) = \frac{1}{2}$$

$$P(w_n | w_{<n}) \approx \prod_i P(w_i) \Rightarrow \text{unigram}$$

## Evaluating language models

<S> Ham and eggs </S>

$$\vdots P(\text{Ham} | \text{<S>}) \cdot P(\text{and} | \text{ham}) \cdot \dots$$

as seqs get longer, the product of these cond. probs gets smaller and smaller

$$\log \prod_i P(w_i | w_{<i}) = \sum \log P(w_i | w_{<i})$$

in general, given a test doc  $x_1, x_2, \dots, x_t$

we want  $P(x_1, \dots, x_t)$  to be high =

$$\frac{1}{t} \cdot \sum_i \log P(x_i | x_{<i}) \text{ to be high}$$

PPL;  
perplexity:  $\exp\left(-\frac{1}{t} \sum_i \log P(x_i | x_{<i})\right)$

↳ want this to be low

intuition: perplexity measures how many equally likely next words is the model choosing from given a prefix

$P(\text{Gemini} | \text{is}) \Rightarrow \text{PPL is high}$

$P(\text{Gemini} | \text{LLM is}) \Rightarrow \text{PPL is lower}$

$P(\text{Gemini} | \text{My favorite LLM is}) \Rightarrow \text{PPL is very low}$