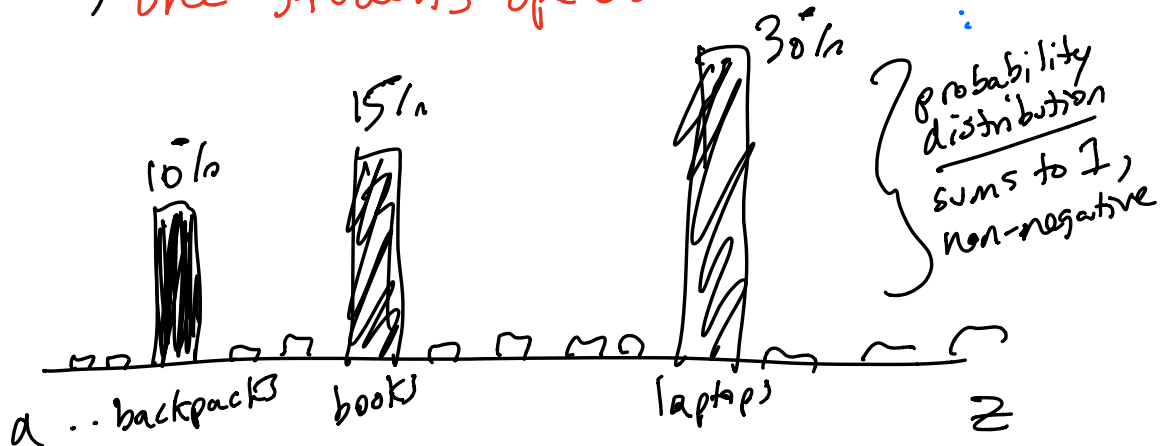


# Language models

↳ given a **prefix**, a LM predicts the **next word**

↳ the students opened their **backpacks**



given  $w_1, w_2, w_3, \dots, w_{n-1}$  we want to compute

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

↳ any model that computes this is called a **language model**

---

n-gram models :

$$p(\text{backpacks} | \text{"the students opened their"})$$

we have a training dataset that's fairly large  
let's extract all occurrences of the prefix

the students opened their books  
the students opened their backpacks  
the students opened their laptops  
the students opened their laptops  
the students opened their laptops  
the students opened their backpacks

6 occurrences of this prefix

$$\left. \begin{aligned} P(\text{books} | \text{prefix}) &= \frac{1}{6} \\ P(\text{laptops} | \dots) &= \frac{1}{2} \\ P(\text{backpacks} | \dots) &= \frac{1}{3} \end{aligned} \right\} \begin{array}{l} \text{Maximum} \\ \text{likelihood} \\ \text{estimate,} \\ \text{maximize the prob.} \\ \text{of training dataset} \end{array}$$

problems?

↳ dataset size and diversity

↳ generalization

the students lazily opened their \_\_\_\_\_

↳ storage  $w_1$   $w_2$   $w_3$  ...  $w_{100k}$   
prefix  
 $\theta_1$  100  
 $\theta_2$  0  
 $\theta_3$  0

∴ ∅ ∅ 5000

↳ Sparse

↳ table size is infinite

Storage issue can be addressed by truncating the prefix

$p(\text{backpacks} \mid \text{the students opened their})$

$\approx p(\text{backpacks} \mid \text{students opened their})$

$\approx p(\text{backpacks} \mid \text{opened their}) \Rightarrow \text{trigram model}$

$\approx p(\text{backpacks} \mid \text{their}) \Rightarrow \text{bigram}$

$\approx p(\text{backpacks}) \Rightarrow \text{unigram}$

vocabulary contains  $V$  words

Storage cost for  $n$ -gram model

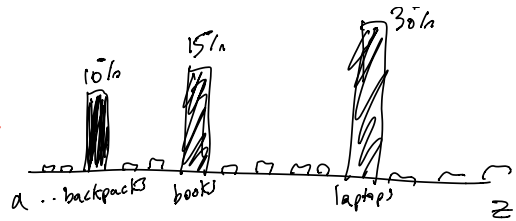
$\propto V^n$

---

decoding: how do we generate a multi-word response from a LM given a prefix?

evaluation: how do we know how good a given LM is?

the students opened their



↳ how do I generate the next word?

↳ just choose the word with the highest prob: **greedy decoding**

↳ randomly choose a word from vocab where prob of choosing  $w_i$  is proportional to  $P_{LM}(w_i | w_{1...i-1})$ : **sampling**

↳ **temperature sampling**

↳ changes the "peakiness" of distribution

↳ **truncation sampling**

↳ do not consider any words  $w_i$  where  $p(w_i | \text{prefix}) < \chi$

the students opened their laptops and

↳ 2% navigated to Google . <EOS>

evaluation: perplexity

let's say we estimated our LM  
on a training dataset

now, we want to see how well  
it does at estimating prob. of a test set

test set:  $t_1 t_2 t_3 \dots t_n$

We want the LM to guess each word  
in the test set without seeing it

We want:

$$P_{LM}(t_3 | t_1, t_2) \Rightarrow \text{to be high}$$

$$P_{LM}(t_4 | t_1, t_2, t_3) \Rightarrow \text{to be high}$$

⋮

$$P_{LM}(t_i | t_1 \dots t_{i-1}) \Rightarrow \text{to be high}$$

extremely important!

average negative log-likelihood

$$\text{test perplexity} = \exp\left(-\frac{1}{n} \sum_i^n \log p(t_i | t_{1, \dots, i-1})\right)$$

4. allows to interpret as "branching factor"

averaging over all words in test set

1. if LM is confident + correct, this prob is very high  
if confident + wrong, prob is very low

2. log will be a tiny negative number  
log will be a very large neg. number

perplexity (PPL):

- on average, how many equally likely words is the LM choosing between?

↳ if low PPL: model is more certain

↳ if high PPL: model is less certain

PPL of a model that gets 100% prob. on the test set is 1

↳ is this achievable?

let's say you estimate a bigram model on some training dataset. you want to eval test ppl.

Your test set contains a two-word phrase that never occurred in training dataset.

↳ what is your test PPL?