

From last time:

LMs measure:

↳  $P(w_1, w_2, \dots, w_n) \Rightarrow$  prob of a text

↳  $P(w_n | w_1, w_2, \dots, w_{n-1}) \Rightarrow$  cond. prob of next word given prefix

n-gram models:

↳ cheap way to get these probabilities by counting/dividing over a training dataset

Perplexity

↳ a metric used to eval LMs

↳ PPL is exponentiated avg. neg log likelihood

↳ ideal case:

PPL (train) is low

PPL (test) is also low

} generalization

↳ bad :

$\left. \begin{array}{l} \text{PPL (train) is low} \\ \text{PPL (test) is high} \end{array} \right\} \text{overfitting}$

↳ to compute PPL, you need an existing dataset of text

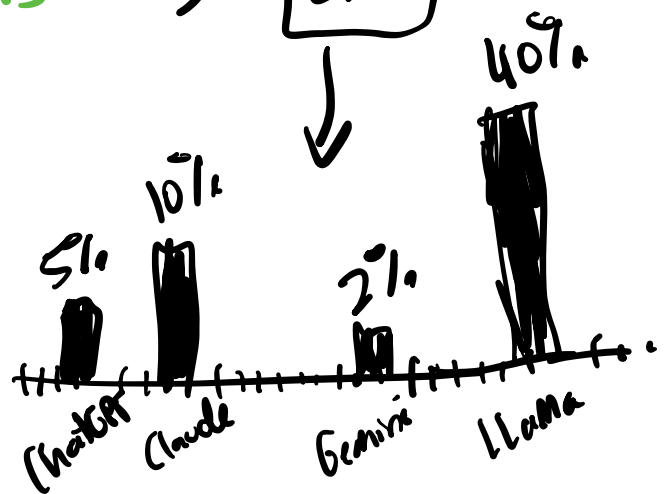
how do we use a LM to generate new text?

↳ this process is called **decoding** from the LM

↳ part of the higher-level stage known as **inference**

My favorite LLM is  $\rightarrow$  **LM**

$P(w_i | \text{prefix})$   
distribution over vocabulary



deterministic

↳ **greedy decoding** : choose  $\underset{w_i}{\operatorname{argmax}} P(w_i | \text{prefix})$

↳ **sampling** word  $w_i$  from this cond. distribution

↳ 40% of the time choose LLaMa

↳ 50% we choose ChatGPT

⋮

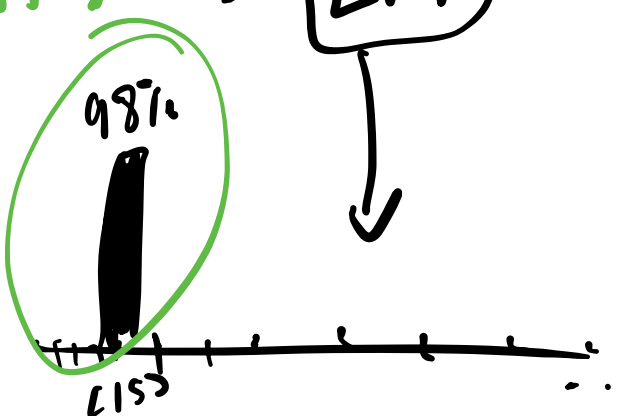
after generating a word/token, we add it  
to our prefix and then generate the next word

My favorite LLM is ChatGPT → **LM**



A bar chart representing the probability distribution of tokens for the prefix "My favorite LLM is ChatGPT". The x-axis lists tokens: "the", ";", "but", "if", "in", "a", "when". The y-axis represents probability. The bars for "the" and ";" are circled in green, indicating they are the most likely next tokens.

My fav. LLM is ChatGPT. → **LM**



A bar chart representing the probability distribution of tokens for the prefix "My fav. LLM is ChatGPT.". The x-axis lists tokens: "the", ";", "but", "if", "in", "a", "when". The y-axis represents probability. The bar for "the" is circled in green and labeled "98%", indicating it is the most likely next token.

My fav. LLM is ChatGPT,  $\langle /s \rangle$   
 $\langle eos \rangle$

↳ autoregressive decoding

## neural language models:

↳ move away from "count + divide"

↳ core concepts:

↳ forward propagation

↳ how we go from input seq to output prob. dist

↳ backpropagation

↳ how we "train" the model to make better predictions

↳ gradient descent

↳ embedding } today

↳ composition }

## embeddings:

↳ in an n-gram model, "movie" and "film" are treated as completely different

↳ "one-hot" representation of a word

$\nearrow V\text{-dimensional}$

$$\text{movie} = \langle 0, 0, 0, \dots, 1, 0, 0, \dots \rangle$$

$\uparrow$   
movie

$$\text{film} = \langle 0, 0, 1, \dots, 0, 0, 0, \dots \rangle$$

$\uparrow \qquad \qquad \uparrow$   
film                      movie

$\swarrow$  dot product

$$\vec{\text{movie}} \cdot \vec{\text{film}} = 0 \quad \left. \vphantom{\vec{\text{movie}} \cdot \vec{\text{film}}} \right\} \text{orthogonal}$$

ideally, we want a vector space  
where words/phrases/docs w/ similar  
meanings have similar representations

↳ what if we switch from sparse  $V\text{-dim}$   
vector to a dense, low-dim vector

$\hookrightarrow d = \{50, 100, 768, 1024\}$  common values

$$\text{Movie} = \langle 0.3, -1.4, 5.8 \rangle$$

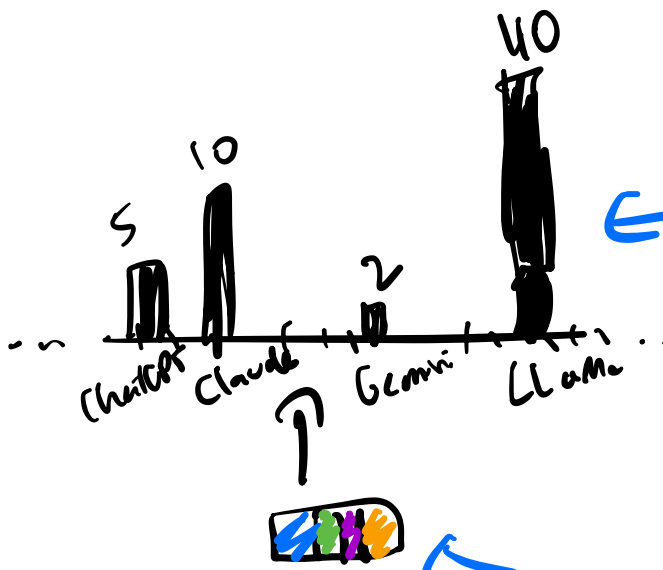
$$\text{film} = \langle 0.2, -1.8, 6.3 \rangle$$

$$\text{zebra} = \langle 7.9, 3.2, -17.5 \rangle$$

→ these are examples of word embeddings

Now, given a sequence of word embs associated w/ some prefix, we need a way to compose them together

composition



← output:

V-dim prob dist over next word

→ sum to 1

→ non-negative

← d-dim prefix vector

Softmax layer:

→ convert a <sup>d-dim</sup> dense vector to a V-dim prob. dist

example : output vocab = { ChatGPT,  
Claude,  
Gemini,  
Llama }

given prefix  
vector   $\leftarrow X$

$$X = \langle -2.3, 0.9, 5.4 \rangle$$

$$d=3, V=4$$

1. project  $X$  to 4-d  
using a weight matrix  $W$

↳ weight matrices contain  
parameters that we update in backprop  
to change the next word dist.

$$W = \begin{pmatrix} 1.2 & -0.3 & 0.9 \\ 0.2 & 0.4 & -2.2 \\ 8.9 & -1.9 & 6.5 \\ 4.5 & 2.2 & -0.1 \end{pmatrix}$$

this is a  $4 \times 3$  matrix  
and  $X$  is 3 dimensional

↳ matrix-vector product

$Wx$  is a 4-d vector

$W = \begin{pmatrix} 1.2 & -0.3 & 0.9 \\ 0.2 & 0.4 & -2.2 \\ 8.9 & -1.9 & 6.5 \\ 4.5 & 2.2 & -0.1 \end{pmatrix}$

feature weights for  
→ ChatGPT  
→ Claude  
→ Gemini  
→ LLaMa

$x = \langle -2.3, 0.9, 5.4 \rangle$

↑  
each dim of  $x$  intuitively  
corresponds to a "feature"  
of the prefix

$$Wx = \langle 1.8, -11.9, 12.9, -8.9 \rangle$$

enter the softmax fn

$$\text{softmax}(x) = \frac{e^x}{\sum_i e^{x_i}}$$



$\uparrow$   
 $x$  is  
a vector

$\uparrow$   
 $x_j$  is  
dimension  $j$  of  $x$

$\text{Softmax}(Wx) =$

$\langle 1.5 \cdot 10^{-5}, 1.6 \cdot 10^{-11}, 0.999984, 3.4 \cdot 10^{-10} \rangle$

$\downarrow$   
 $P(\text{ChatGPT} | \dots)$

$\downarrow$   
 $P(\text{Claude} | \text{prefix})$

$\downarrow$   
 $P(\text{Gemini} | \text{prefix})$

$\downarrow$   
 $P(\text{Llama} | \text{prefix})$