

in an n-gram model:

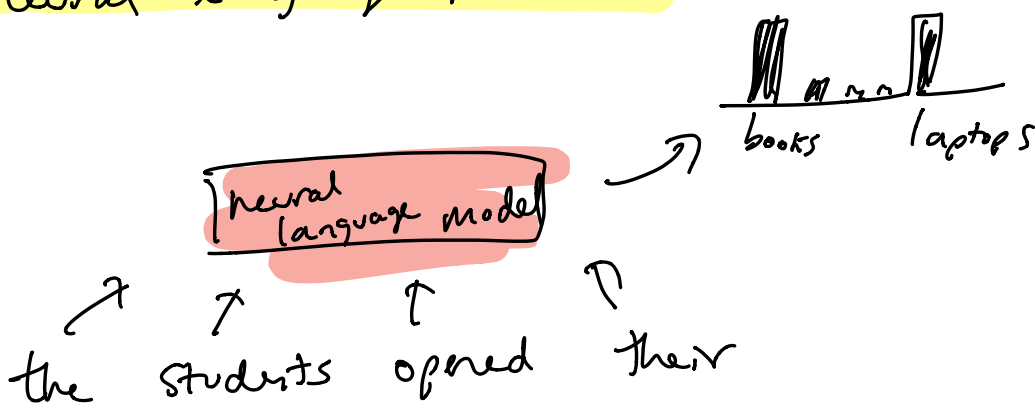
the students opened their _____

the student opened her _____

the iPad exploded because _____

↳ how do we share info between semantically-similar prefixes?

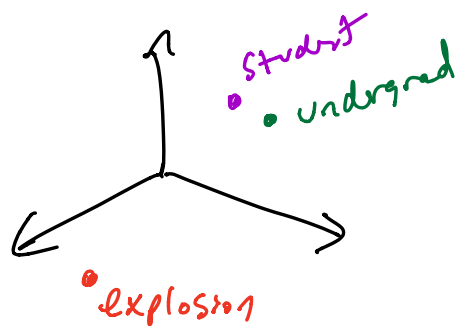
neural language models



↳ NLMs need to be trained

↳ NLMs represent words, phrases, sentences, documents w/ dense real-valued vectors embeddings

student: $[-0.43, 0.72, 1.1]$



word embeddings

↳ **type** vs. **tokens**

↓
unique
entry in
vocabulary

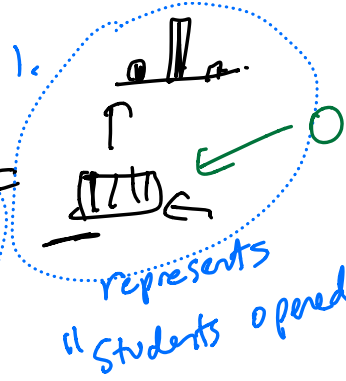
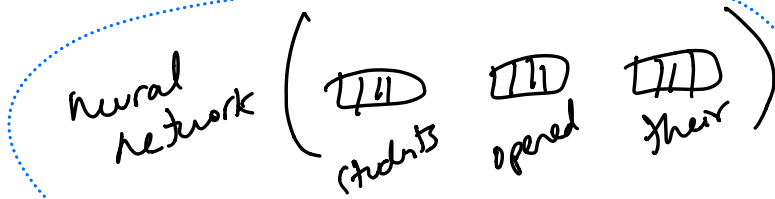
↓
occurrence
of a type
in a text

the blue ants hated the red ants.
 1 2 3 4 5 6 7 8

↳ embs are usually 50-1000d

↳ embs are stored in a $V \times d$ matrix

Composition

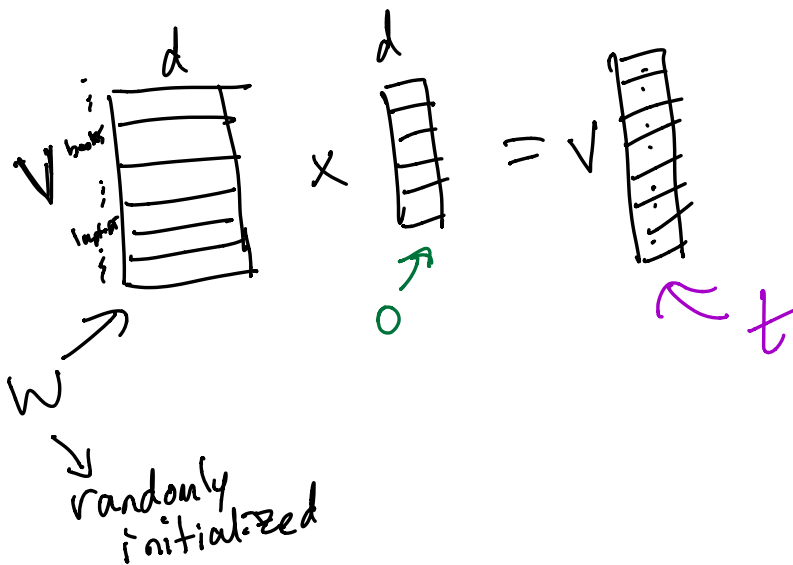


2.

↳ vocab has V word types
embs are d dimensional

↳ if we do a matrix vector product, we can project our emb. into a V -dim vector

→ W is a $V \times d$ matrix



$$t = [1.8, -1.9, 2.9, \dots]$$

t_0 t_1

↳ we use the softmax function to squash any input vector into a probability distribution

$$\text{Softmax}(t) = \frac{e^t}{\sum_j e^{t_j}}$$

→ making everything positive
 → make everything sum to 1
 t_j is the value at the j^{th} dim of t

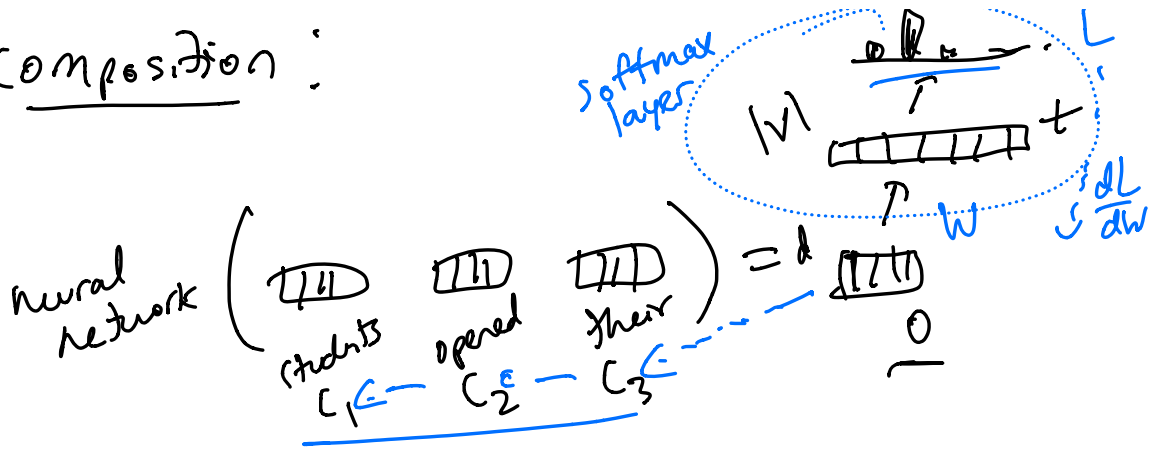
let's say $V = \{ \text{books, cars, laptops, zoo} \}$

if $t = [1.8, -1.9, 2.9, -0.9]$

$$\text{Softmax}(t) = [0.24, 0.006, 0.73, 0.02]$$

\downarrow \downarrow
 $p(\text{books} | \text{prefix})$ $p(\text{car} | \text{prefix})$

Composition:



architecture:

↳ Transformers

↳ recurrent neural networks

↳ element-wise / feed-forward NNs

$O = C_1 + C_2 + C_3 \Rightarrow$ element-wise addition

↳ issue #1: no position info

let's stick w/ simple addition
to define our NLM

↳ NLMs contain **parameters** θ

$$\theta = \{ W, c_1, c_2, c_3, \dots \}$$

↳ randomly initialized

↳ **trained** to minimize neg. log likelihood of training dataset

↳ how do we train NLM?

1. define a **loss fn**

↳ tell us how bad we are at predicting the next word

$$\hookrightarrow L(\theta) = -\log p(w_n | w_1, \dots, w_{n-1})$$

→ neg log prob of correct next word in training dataset

2. given $L(\theta)$, we compute the **gradient** of L w.r.t θ

→ gradient gives the direction of steepest ascent

$$\rightarrow \frac{dL}{d\theta} = \left\{ \frac{dL}{dW}, \frac{dL}{dc_1}, \frac{dL}{dc_2}, \frac{dL}{dc_3} \right\}$$

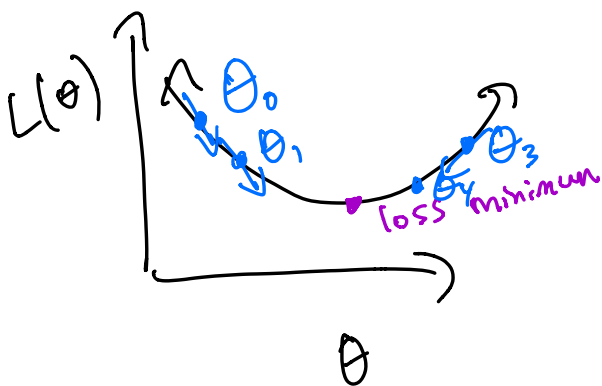
→ $\frac{dL}{dW}$ is same dim as W

→ intuition: for each param in Θ ,
 $\frac{dL}{d\theta}$ tells us how much L changes
if we increase that param by a
tiny amount

3. Given gradient $\frac{dL}{d\theta}$, we take a
step in the direction of negative gradient

$$\Theta_{\text{new}} = \Theta_{\text{old}} - \eta \frac{dL}{d\theta}$$

learning rate, "step size" → η
gradient → $\frac{dL}{d\theta}$



optimizers:

↳ Adam

↳ Sophia

↳ SGD

↑ this is called gradient descent

hyperparameters:

↳ learning rate

↳ batch size

→ how many training examples do you use to estimate $\frac{dL}{d\theta}$

→ algorithm to compute gradient efficiently in neural networks is called **backpropagation**

→ loss.backward()

↳ architectures for NLMs

→ recurrent neural network

→ Transformer

→ attention

→ state space models (Mamba)