

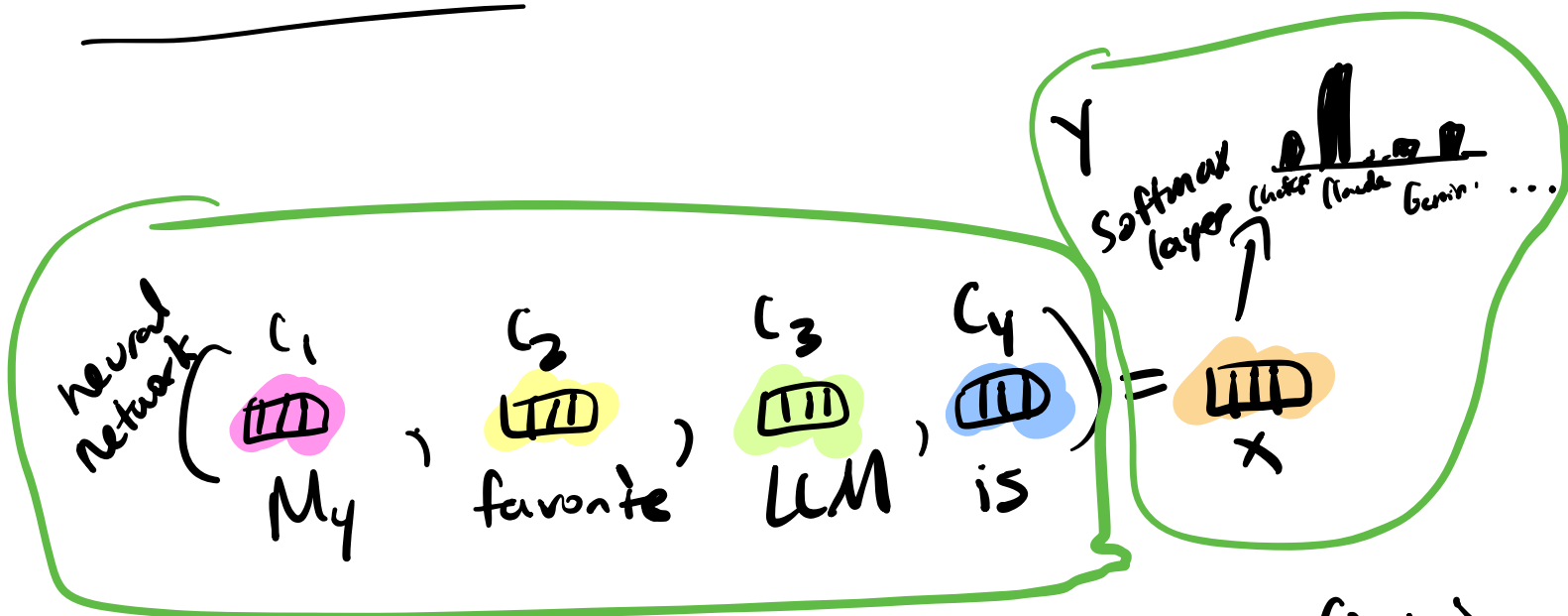
Logistics

↳ Hwo, ^{final} project group assignments due Wed

↳ Exam 1: 10/20

Exam 2: 12/10

Final project report: 12/20



Composition



1. element-wise functions
2. fixed-window neural LLMs
3. recurrent neural LLMs
4. Transformers
 - ↳ self-attention

$$Y = \text{Softmax}(Wx)$$

Diagram illustrating the Softmax function: $Y = \text{Softmax}(Wx)$. The input vector x is $V \times d$ (V-dim, d-dim). The weight matrix W is $V \times d$ (V-dim, d-dim). The output vector Y is V -dim.

takes an input vector and outputs a vector that is positive and sums to 1

Composition function:

↳ input: sequence of d -dim embeddings corresponding to tokens of the prefix

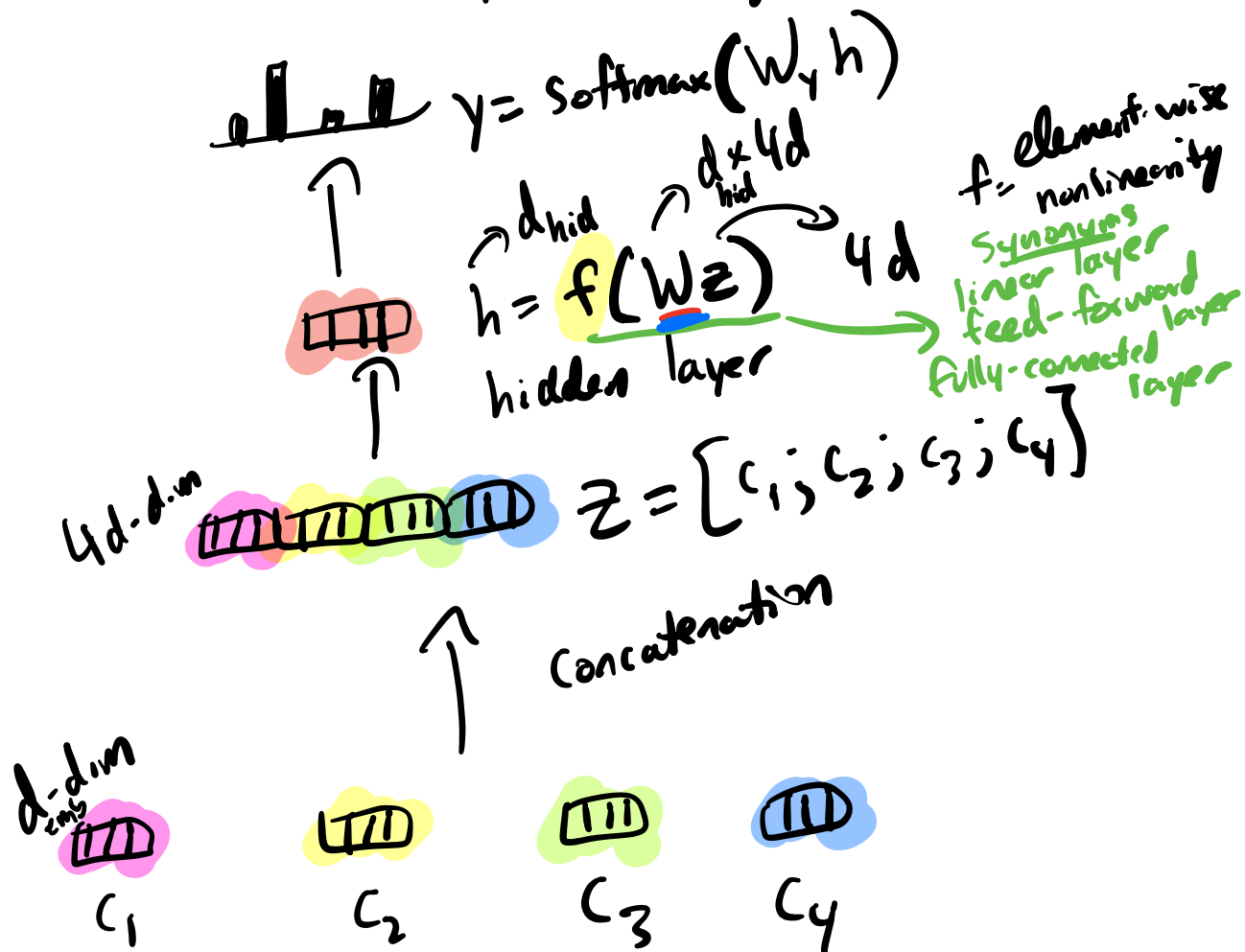
↳ output: single vector representing the entire prefix

Simplest fn:

↳ element-wise sum:

$$x = \sum_i^n c_i \quad \text{problems?}$$

↳ concatenation: fixed-window neural LM neural n-gram model



Why nonlinearity?

↳ in LM (and many other tasks), the relationship between the input and the output is highly complex and nonlinear

↳ without nonlinearity, you can only model linear relationships

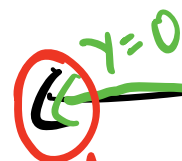
$$\begin{aligned} b &= W_1 a \\ c &= W_2 b \\ d &= W_3 c \\ &\vdots \end{aligned} \quad \left. \vphantom{\begin{aligned} b &= W_1 a \\ c &= W_2 b \\ d &= W_3 c \\ &\vdots \end{aligned}} \right\} c = W_2 W_1 a \quad \begin{aligned} &\underbrace{\hspace{1cm}} \\ &\rightarrow W_{\text{new}} a \end{aligned}$$

↑
Simplification doesn't work if

$$\begin{aligned} b &= f(W_1 a) \\ c &= f(W_2 b) \end{aligned}$$

ReLU (rectified linear unit)

$$\text{ReLU}(x) = \max(0, x)$$

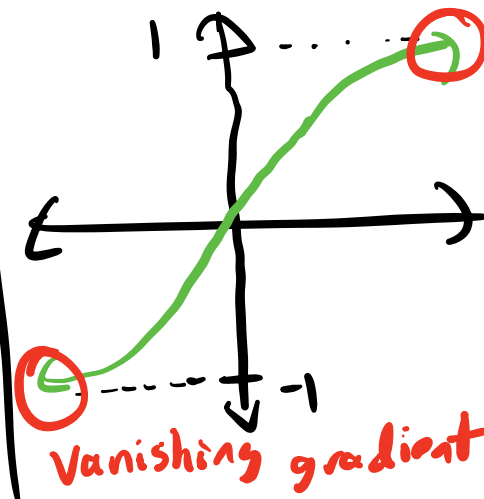


dead neurons

$$x = \langle -2.3, 5.7, -5.0 \rangle$$

$$\text{ReLU}(x) = \langle 0, 5.7, 0 \rangle$$

tanh(x)



Vanishing gradient

Comparing fixed-window NLM to n-gram model:

↳ storage: model size scales linearly with n rather than exponentially

↳ reduced sparsity problem

↳ hard to interpret

↳ fixed prefix window

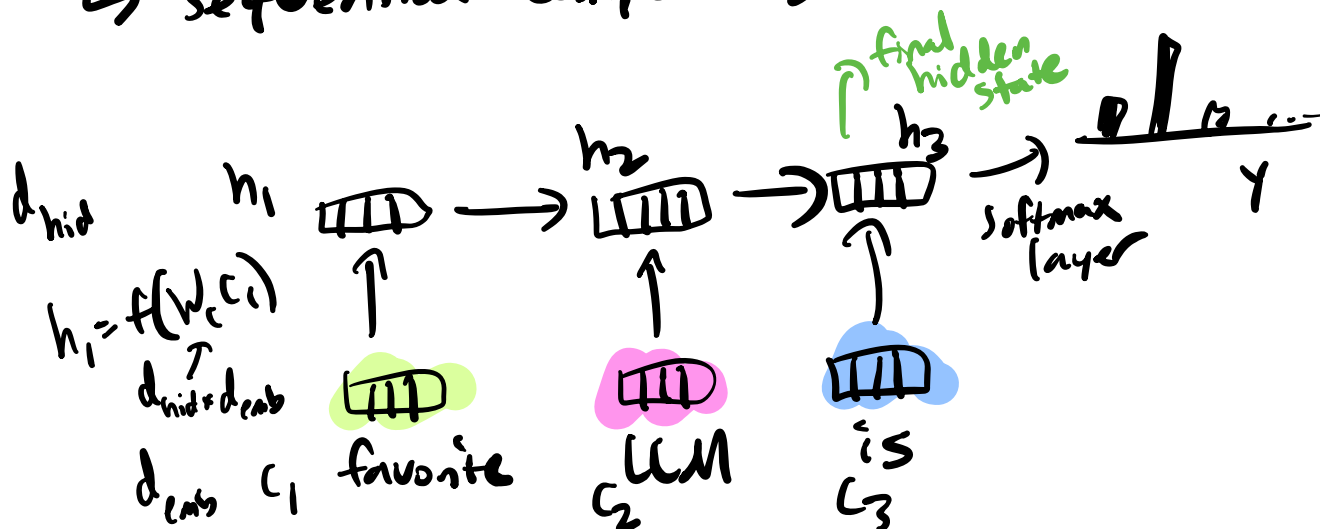
↳ impossible to handle long-range deps

↳ doesn't share weights between different positions in the prefix

"My favorite LLM is"
vs. "favorite LLM really is"
2

Recurrent neural networks:

↳ sequential composition, one word at a time



↳ $h_t \Rightarrow$ hidden state at **timestamp t**
position

↳ h_t is a fn of h_{t-1} and c_t
 \uparrow \uparrow
prev. hidden state current token embedding

↳ h_t and c_t don't have to be the same dimensionality

$$h_t = \underset{\text{tanh}}{f}(\underset{\substack{\uparrow \\ d_{\text{hid}} \times d_{\text{hid}}}}{W_h} h_{t-1} + \underset{\substack{\hookrightarrow d_{\text{hid}} \times d_{\text{emb}}}}{W_c} c_t)$$

high-level training overview:

↳ NLMs contain **trainable parameters**

$$\Theta_{\text{nn}} = \{ W_h, W_c, c_{\dots v}, W_y \}$$



1. **randomly** initialized
2. **trained** to better predict next word on a training dataset

↳ define a loss fn

↳ tells us how bad we are doing at predicting the next word

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

↳ neg log prob of correct next word in training dataset

↳ given $L(\theta)$, we compute the gradient of L wRT θ

↳ gradient provides direction of steepest ascent

↳ take a step in direction of negative gradient

$$w_{h_{\text{new}}} = w_{h_{\text{old}}} - \eta \frac{dL}{dw_h}$$

↳ learning rate