

This class:

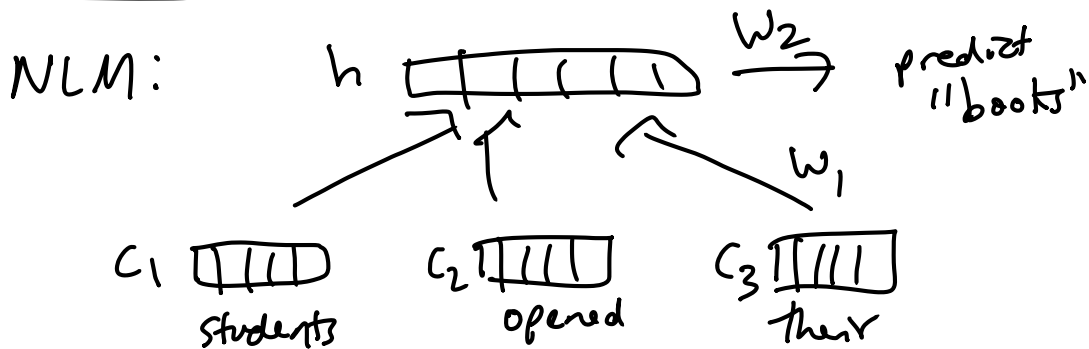
- gradient descent  $\Rightarrow$  cross-entropy loss

- backpropagation

$\hookrightarrow$  feed-forward NN } single neurons  
 $\hookrightarrow$  recurrent NN }

$\hookrightarrow$  see F2020 video for  
backprop thru linear layer

---



params:  $w_1, w_2, c_1, c_2, c_3$

$$h = f(w_1 [c_1; c_2; c_3])$$

$$o = \text{Softmax}(w_2 h)$$

$h, o$  are intermediate vars

---

how do we train this model

to make better predictions  $\Rightarrow$  gradient descent!

1. define loss fn  $L(\theta)$  <sup>all params</sup> that tells us how bad the model is doing at predicting the next word

$\hookrightarrow$  cross entropy loss

↳ say we have a training ex

students opened their  $\Rightarrow$  books  
input prefix target

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

↳ maximize this probability

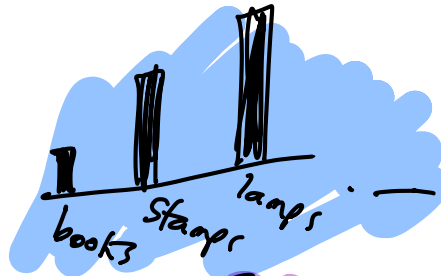
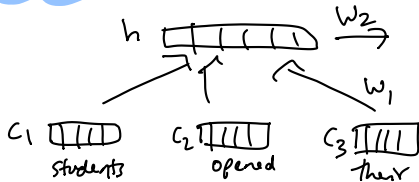
in practice,

we minimize the negative log prob of "books"

$$L = -\log (P(\text{books} | \text{"students opened their"}))$$

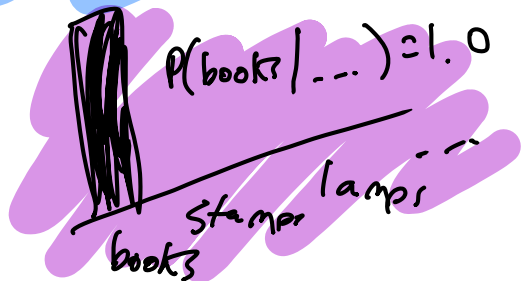
↳ cross entropy loss

Model:



data:

students opened their



Cross-entropy loss of two dists  $p$  and  $q$

$$-\sum_{w \in V} p(w) \log q(w)$$

↳ data

↳ model

1 for books, 0 for everything else

$$= -\log q(\text{books} \mid \text{students opened them})$$

(neg. log prob of correct word)

okay, so we have a loss fn  $L(\theta)$

need to compute **gradient** of  $L$  WRT  $\theta$

$$\frac{dL}{d\theta}$$

↳ model  
params

↳ gradient tells us the direction of steepest ascent of  $L$ , intuitively, tells us how  $L$  changes when we modify  $\theta$

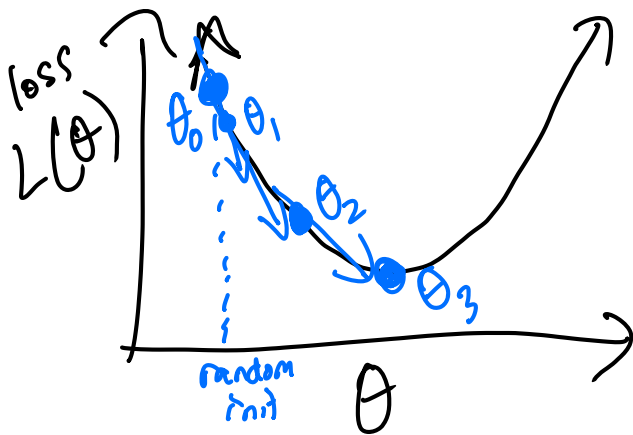
3. after computing  $\frac{dL}{d\theta}$ , take a step in direction of negative gradient, minimizing  $L$

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

↳ learning rate,  
controls step size

learning rate is a hyperparameter

batch size is important hyperparam



Simple example:



compute loss fn:  $L = \frac{1}{2} (y - o)^2$  square loss

$\hookrightarrow$  target  $\hookrightarrow$  model prediction

compute gradient  $\frac{dL}{d\theta} : \frac{dL}{dw_1}, \frac{dL}{dw_2}$  Chain rule of calculus

$$L = \frac{1}{2} (y - o)^2$$

$$o = \tanh(w_2 h)$$

intermediate var  
 $a = w_2 h$   
 $b = w_1 x$

$$o = \tanh(a)$$

$a = w_2 h, h = \tanh(b), b = w_1 x$

$$\frac{d}{dx} g(f(x)) = \frac{dg}{df} \cdot \frac{df}{dx}$$

$$\frac{d}{dx} \tanh(x) = 1 - \tanh^2(x)$$

$$\frac{dL}{dw_2} = \frac{dL}{do} \cdot \frac{do}{da} \cdot \frac{da}{dw_2}$$

$$\downarrow \quad \downarrow$$

$$-(y - o) \cdot (1 - o^2) \cdot h$$

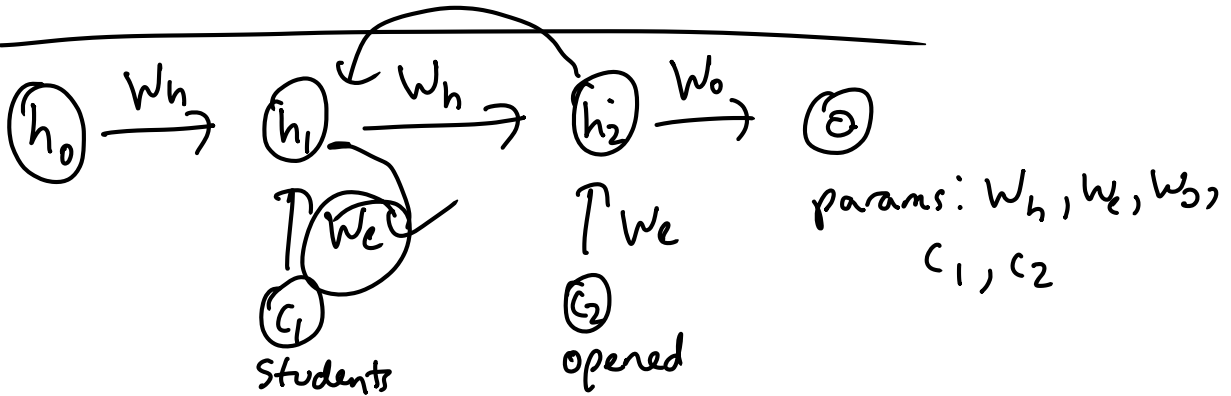
$$\frac{dL}{dw_1} = \frac{dL}{do} \cdot \frac{do}{da} \cdot \frac{da}{dh} \cdot \frac{dh}{db} \cdot \frac{db}{dw_1}$$

backprop = chain rule + caching prev. computed derivatives

Update params:

$$W_{2, \text{new}} = W_{2, \text{old}} - \eta \frac{dL}{dW_2}$$

$$W_{1, \text{new}} = W_{1, \text{old}} - \eta \frac{dL}{dW_1}$$



$$L = \frac{1}{2} (y - o)^2$$

$$o = W_o h_2$$

$$h_2 = \tanh(W_{c_2} c_2 + W_h h_1)$$

$$h_1 = \tanh(W_{c_1} c_1 + W_h h_0)$$

$$\frac{dL}{dW_o} = \frac{dL}{do} \cdot \frac{do}{dW_o} = -(y - o) \cdot h_2$$

$$\frac{dL}{dC_2} = \frac{dL}{do} \cdot \frac{do}{dh_2} \cdot \frac{dh_2}{da} \cdot \frac{da}{dC_2} = -(y - o) \cdot W_o \cdot (1 - h_2^2) \cdot W_{c_2}$$

$$\frac{dL}{dC_1} = \frac{dL}{do} \cdot \frac{do}{dh_2} \cdot \frac{dh_2}{da} \cdot \frac{da}{dh_1} \cdot \frac{dh_1}{db} \cdot \frac{db}{dC_1}$$

$\frac{dL}{dW_e}$  and  $\frac{dL}{dW_h}$  are trickier b/c of shared weights

"backpropagation thru time" allows us to compute these gradients by summing the contributions from diff. time steps

$$\frac{dL}{dW_e} = \frac{dL}{do} \cdot \frac{do}{dh_2} \cdot \frac{dh_2}{da} \cdot \frac{da}{dW_e} + \dots$$

we accumulate these  $\frac{dL}{dW_e}$  as we step back thru time.

the vanishing gradient problem occurs in RNNs when gradient contributions at faraway time steps go to zero