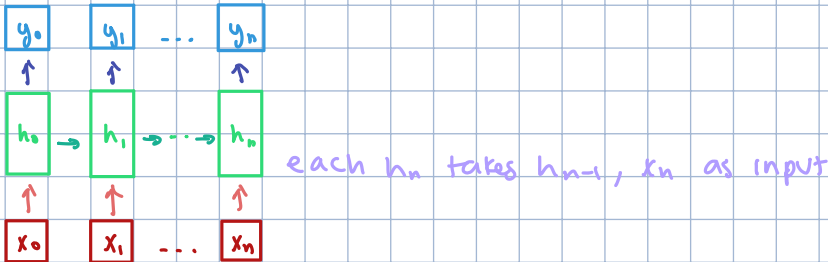


RNN: why

- vanilla neural nets (& CNNs) take fixed len in, produce fixed len out
- RNN can do variable len

How work? (many to many example)

- input: x_0, x_1, \dots, x_n
- output: y_0, y_1, \dots, y_n
- steps: iteratively update hidden state h (matrix w/ arbitrary dim). At step t :
 - 1) next hidden state calculated using h_{t-1} & next input x_t



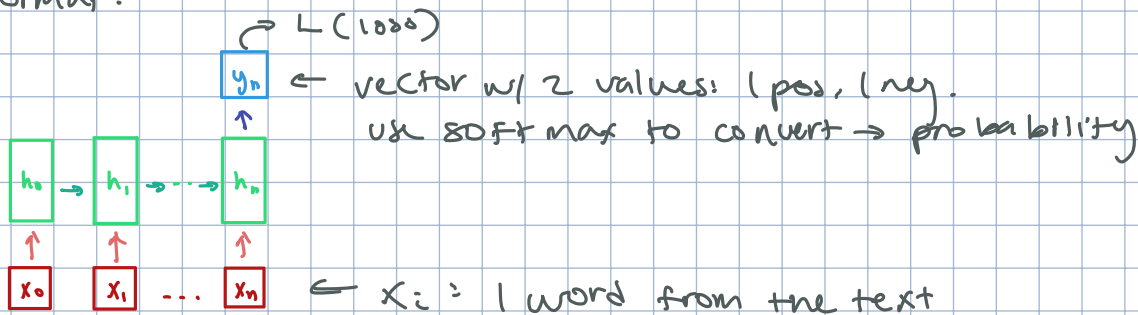
- vanilla RNN: 3 sets of weights
 - w_{xh} : $x_t \rightarrow h_t$ links
 - w_{hh} : $h_{t-1} \rightarrow h_t$
 - w_{hy} : $h_t \rightarrow y_t$
 - biases: b_h (for calc h_t)
 b_y (for calc y_t)
- representation: weights as MATRICES, biases as VECTORS
- calculation:
 - $h_t = \tanh(w_{xh}x_t + w_{hh}h_{t-1} + b_h)$ \Leftarrow can use another activation function
 - $y_t = w_{hy}h_t + b_y$ \Leftarrow like sigmoid

The Problem

- sentiment analysis: determine if text is pos. or neg.
- ex: I am good \Rightarrow good
I am not happy \Rightarrow bad

Plan

- use "many to one" RNN
- final output will just be the last y_n
- format:



sizes: x_i : input_size $\times 1$ \Leftarrow input_size = vocab_size
 w_{xh} : hidden_size \times input_size b_h : hidden_size $\times 1$
 w_{hh} : hidden_size \times hidden_size w_{hy} : output_size \times hidden_size
 b_y : output_size $\times 1$
 h_t, h_{t-1} : hidden_size $\times 1$

Backprop

- can use cross-entropy loss, which is often paired w/ softmax

$$L = -\ln(p_c)$$

p_c = RNN's predicted probability for the correct class

ex: if pos. text is predicted as 90% positive by RNN, loss is:

$$L = -\ln(.90) = .105$$

using cross-entropy guides the RNN to 1) get the correct ans & 2) be more confident in its answer

ex: if pos text is predicted as 10% positive by RNN:

1) it's wrong

2) it's not confident.

∴ incurs large loss of 2.3026

- defn:

y = RAW outputs from RNN

p = probabilities ($p = \text{softmax}(y)$)

c = correct label

L = cross entropy loss ($L = -\ln(p_c)$)

W_{xh}, W_{hh}, W_{hy} = 3 weight matrices

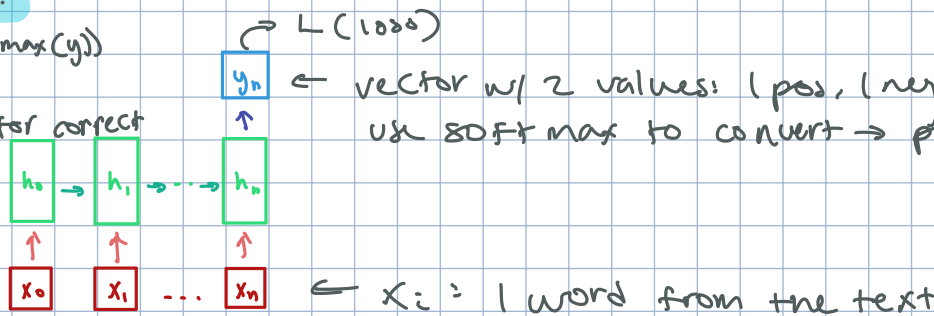
b_h, b_y = 2 bias vectors

- gradients:

start by calc dL/dy :

$$L = -\ln(p_c) = -\ln(\text{softmax}(y))$$

p_c = RNN's predicted prob for correct class.



y = [probability of positive, prob. of neg]

ex: [.90, .10] means 90% chance positive

p = probability of positive

$$L = c(-\ln(p)) + (1-c)(-\ln(1-p))$$

↳ if c (correct label) is 1 (aka positive) return $(-\ln(p))$. the $(1-c)(-\ln(1-p))$ part eval to 0
if c is 0, return $-\ln(1-p)$, aka $-\ln(\text{prob-neg})$
the $c(-\ln(p))$ term evaluates to 0

y_p = RNN's predicted chance of phrase being positive

$$\frac{dL}{dy_p} = \frac{d}{dy_p} [c(-\ln p) + (1-c)(-\ln(1-p))]$$

$$= -c \frac{d}{dy_p} \ln(\text{softmax}(y_p)) - (1-c) \frac{d}{dy_p} (\ln(1-p))$$

$\frac{\partial L}{\partial y_p}$

$$\begin{aligned}\frac{d}{dy_p} \ln(\text{sm}(y_p)) &= \frac{1}{\text{sm}(y_p)} \cdot \frac{d}{dy_p} \text{sm}(y_p) \\ &= \frac{\sum_{i=1}^n e^{y_i}}{e^{y_p}} \cdot \frac{d}{dy_p} \left(\frac{e^{y_p}}{\sum_{i=1}^n e^{y_i}} \right) \quad \text{quotient rule} \\ &= \frac{\sum_{i=1}^n e^{y_i}}{e^{y_p}} \cdot \frac{e^{y_p} (\sum_{i=1}^n e^{y_i}) - (e^{y_p})^2}{(\sum_{i=1}^n e^{y_i})^2} \\ &= \frac{1 - \text{sm}(y_p)}{1 - \text{sm}(y_p)}\end{aligned}$$

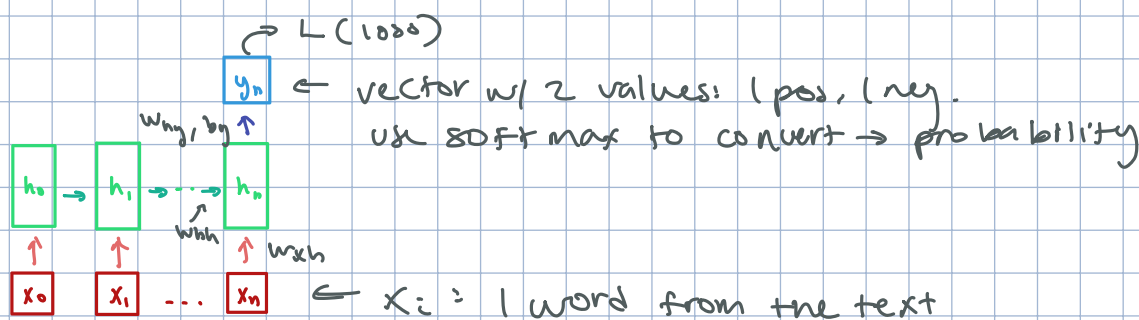
$$\frac{dL}{dy_p} = c(\text{sm}(y_p) - 1) - (1-c) \frac{1}{1 - \text{sm}(y_p)} \cdot \frac{d}{dy_p} (1 - \text{sm}(y_p))$$

$$\begin{aligned}\frac{dL}{dy_p} &= c(\text{sm}(y_p) - 1) - \frac{(1-c)}{1 - \text{sm}(y_p)} \left[\frac{-e^{y_p} (\sum_{i=1}^n e^{y_i} - e^{y_p})}{(\sum_{i=1}^n e^{y_i})^2} \right] \\ &\quad + \frac{(1-c)}{\sum_{i=1}^n e^{y_i} - e^{y_p}} \left[\frac{+e^{y_p} (\sum_{i=1}^n e^{y_i} - e^{y_p})}{(\sum_{i=1}^n e^{y_i})^2} \right]\end{aligned}$$

$$= c(\text{sm}(y_p) - 1) + (1-c)(\text{sm}(y_p))$$

or in other words:

$$\frac{\partial L}{\partial y_i} = \begin{cases} p_i - 1 & \text{if } i = c \\ p_i & \text{if } i \neq c \end{cases}$$



Compute gradient for w_{ny}, b_y :

$$\frac{\partial L}{\partial w_{ny}} = \frac{\partial L}{\partial y} \cdot \frac{\partial y}{\partial w_{ny}}$$

$$y = w_{ny} h_n + b_y$$

$$\frac{\partial y}{\partial w_{ny}} = h_n$$

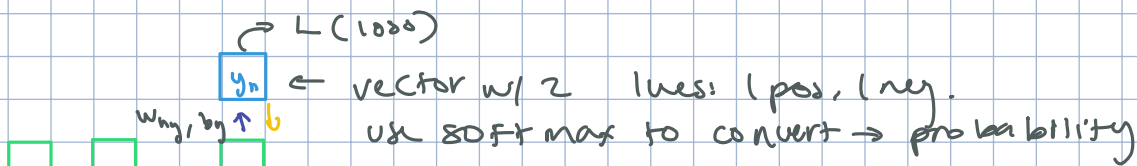
$$\therefore \frac{\partial L}{\partial w_{ny}} = \frac{\partial L}{\partial y} \cdot h_n$$

$$\frac{\partial L}{\partial b_y} = \frac{\partial L}{\partial y} \cdot \frac{\partial y}{\partial b_y}$$

$$y = w_{ny} h_n + b_y$$

$$\frac{\partial y}{\partial b_y} = 1$$

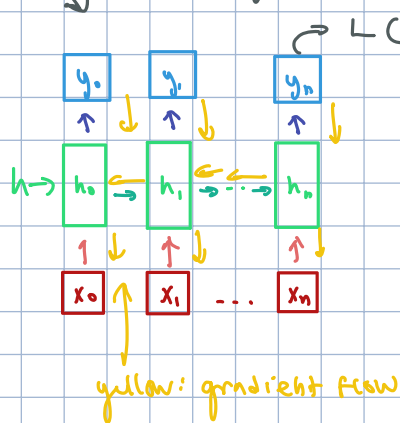
$$\therefore \frac{\partial L}{\partial b_y} = \frac{\partial L}{\partial y}$$



x_i : 1 word from the text

how do gradient for the weights?

full diagram (if we calculated each y_n) would be:



so aggregate loss (FL) becomes

$$FL = \sum L(y_t)$$

$$\frac{\partial FL}{\partial w_{xh}} = \sum \frac{\partial L(y_t)}{\partial y_t} \cdot \frac{\partial y_t}{\partial h_t} \cdot \frac{\partial h_t}{\partial w_{xh}}$$

as approximation, just have all the y 's equal the final y

$$\Rightarrow \frac{\partial FL}{\partial w_{xh}} \approx \frac{\partial L}{\partial w_{xh}} = \frac{\partial L}{\partial y} \approx \frac{\partial y}{\partial h_t} \cdot \frac{\partial h_t}{\partial w_{xh}}$$

(similarly,)

calculate $\frac{\partial h_t}{\partial w_{xh}}$:

$$h_t = \tanh h (w_{xh} x_t + w_{hh} h_{t-1} + b_h)$$

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

\therefore

$$\frac{\partial h_t}{\partial w_{xh}} = (1 - h_t^2) x_t$$

calculate $\frac{\partial h_t}{\partial w_{hh}}$:

$$\frac{\partial h_t}{\partial w_{hh}} = (1 - h_t^2) h_{t-1}$$

calculate $\frac{\partial b_h}{\partial h_t}$:

$$\frac{\partial h_t}{\partial b_h} = (1 - h_t^2)$$

wish: $\frac{\partial L}{\partial y} \cdot \frac{\partial y}{\partial h} = \frac{\partial L}{\partial h}$

$$\frac{\partial L}{\partial w_{hh}} = \frac{\partial L}{\partial y} \approx \frac{\partial y}{\partial h_t} \cdot \frac{\partial h_t}{\partial w_{hh}}$$

$$\frac{\partial L}{\partial b_h} = \frac{\partial L}{\partial y} \approx \frac{\partial y}{\partial h_t} \cdot \frac{\partial h_t}{\partial b_h}$$

calculate $\frac{\partial y}{\partial h_t}$

$$\frac{\partial y}{\partial h_t} = \frac{\partial y}{\partial h_{t+1}} \cdot \frac{\partial h_{t+1}}{\partial h_t}$$

$$\text{ex: } \frac{\partial y}{\partial h_0} = \frac{\partial y}{\partial h_1} \cdot \frac{\partial h_1}{\partial h_0}$$

$$\frac{\partial y}{\partial h_t} = \frac{\partial y}{\partial h_{t+1}} \cdot \frac{\partial h_{t+1}}{\partial h_t}$$

$$= \frac{\partial y}{\partial h_{t+1}} \cdot (1 - h_{t+1}^2) w_{hh}$$

but $\frac{\partial y}{\partial h_n} = w_{hy}$ ← only true for final hidden state

$$\frac{\partial L}{\partial w_{xh}} = \sum \frac{\partial L}{\partial h} (1-h_t^2) x_t$$

$$\frac{\partial L}{\partial w_{nh}} = \sum \frac{\partial L}{\partial h} (1-h_t^2) h_{t-1}$$

$$\frac{\partial L}{\partial b_h} = \sum \frac{\partial L}{\partial h} (1-h_t^2)$$

$$\frac{\partial L}{\partial h} = \frac{\partial L}{\partial y} \cdot \frac{\partial y}{\partial h}$$

$$\text{temp} = (1-h_{t+1}^2) \frac{\partial L}{\partial h}$$

$$\text{next } \frac{\partial y}{\partial h} = \frac{\partial y}{\partial h} (1-h_t^2) w_{nh}$$

$$\begin{aligned} \text{next } \frac{\partial L}{\partial h} &= \frac{\partial L}{\partial y} \frac{\partial y}{\partial h} (1-h_t^2) w_{nh} \\ &= \frac{\partial L}{\partial h} (1-h_t^2) w_{nh} \end{aligned}$$

$$= (\text{temp}) w_{nh}$$