

LSTM:

9 weight matrices and 5 bias vectors to learn

Act as - encoding
on input token

Act on previous
timestep

Remember: W_r, B_r

V_r

New: W_n, B_n

V_n

output: W_o, B_o

V_o

Context: $W_{c_{temp}}, B_{c_{temp}}$

$V_{c_{temp}}$

classify: W_y, B_y

$$x^T = (d_{in}, N)$$

$$h^T = (d_h, N)$$

Dimensions:

d_{in} : input dim
of token

N : L or size

d_h : hidden dim

d_y : number
output lab

d_s : seq length

Variables to
learn:

$$W_{r,n,d} \text{ dim} = (d_h, d_{in})$$

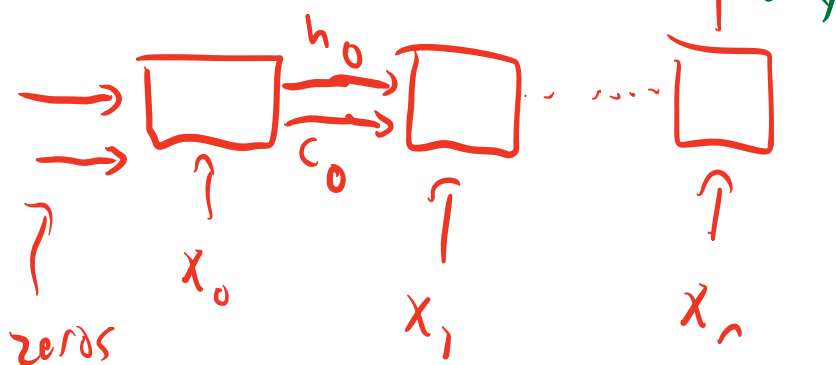
$$B_{r,n,d} \text{ dim} = (d_h, 1)$$

$$V \text{ dim} = (d_h, d_h)$$

$$W_y \text{ dim} = (d_y, d_h)$$

$$B_y \text{ dim} = (d_y, 1)$$

model
topology



Forward Pass: computing units

compute for $t=0$ to $t=d_s$ (whole sequence)
(d_h, N)

current Content: $C_{temp}^+ = \tanh(W_{C_{temp}} x^+ + U_{C_{temp}} h^{+-} + B_{C_{temp}})$

Remember: $r^+ = \sigma(W_r x^+ + U_r h^{+-} + B_r)$

New: $n^+ = \sigma(W_n x^+ + U_n h^{+-} + B_n)$

Content: $C^+ = r^+ \circ C^{+-1} + n^+ \circ C_{temp}^+$

output: $O^+ = \sigma(W_o x^+ + U_o h^{+-} + B_o)$

hidden: $h^+ = O^+ \circ \tanh(C^+)$

Get label after sequence run-through

linear output: $X_{out} = W_y h^{d_s-1} + B_y$

label distribution: $\hat{y} = \text{Softmax}(X_{out})$

Backwards Pass:

dim (dy, N)

Let output be Y where Y_i is one-hot encoding of class label of sample i . (samples are columns)

Let's use categorical cross-entropy loss.
The avg loss for given batch:

$$L = -\frac{1}{N} \sum_{i=0}^N \sum_{c=0}^{d_y} Y_{c,i} \cdot \log(\hat{Y}_{c,i})$$

Now with softmax and cross-entropy loss we have that

Output layer:

$$\frac{dL}{dx_{out}} = \underbrace{\hat{Y} - Y}_{(d_y, N)} \cdot \underbrace{\frac{1}{N}}_{\text{think should average over batch}}$$

Label layers

$$\bullet \frac{dL}{dW_y} = \frac{dL}{dx_{out}} \cdot (h^{d_s-1})^T$$

$$\bullet \frac{dL}{dB_y} = \frac{dL}{dx_{out}} \cdot \mathbf{1}_{(d_h, 1)}$$

$$\bullet \frac{dL}{dh^{d_s-1}} = (W_y)^T \cdot \frac{dL}{dx_{out}} \quad \left. \begin{array}{l} \text{link to recursive} \\ \text{structure} \\ \text{(first } \frac{dL}{dh^+} \text{)} \end{array} \right\}$$

$(d_h, N) \quad (d_h, d_y) \quad (d_x, N)$

Recursive layers

output:

$$\bullet \frac{dL}{d\sigma^+} = \frac{dL}{dh^+} \circ \tanh(c^+)$$

★ All W, B, V needed to be added each iteration on 3 seconds

Content:

$$\bullet \frac{dL}{dc^+} = \frac{dL}{dh^+} \circ o^+ \circ (1 - \tanh(c_+)^2)$$

$$\bullet \frac{dL}{dr^+} = \frac{dL}{dc^+} \circ c^{+-1}$$

$$\bullet \frac{dL}{dn_+^+} = \frac{dL}{dc^+} \circ c_{temp}^+$$

$$\bullet \frac{dL}{dc_{temp}^+} = \frac{dL}{dc^+} \circ n^+$$

General Form for W, V, B

For component of type $k \in K$

from $k = \{c_{\text{temp}}, r, n, o\}$

$$\bullet \frac{dL}{dW_k^+} = \left[\underbrace{\frac{dL}{dk^+}}_{\text{upstream}} \circ \underbrace{(k^+ \circ (1 - k^+))}_{\text{sigmoid activ deriv}} \right] \cdot \underbrace{(x^+)^T}_{(N, d_{in})}$$

(d_n, N)

$$\bullet \frac{dL}{dB_k^+} = \left[\frac{dL}{dk^+} \circ (k^+ \circ (1 - k^+)) \right] \cdot \underbrace{1}_{(N, 1)}$$

$$\bullet \frac{dL}{dV_k} = \left[\frac{dL}{dk^+} \circ (k^+ \circ (1 - k^+)) \right] \cdot \underbrace{(h^{+1})^T}_{(N, d_n)}$$

$$\bullet \frac{dk^+}{dh^{+1}} = \underbrace{v_k}_{(d_n, d_n)} \cdot \underbrace{(k^+ \circ (1 - k^+))}_{(d_n, N)} \quad \left. \begin{array}{l} \text{will be} \\ \text{used to} \\ \text{compute} \\ \text{overall} \end{array} \right\} \frac{dL}{dh^{+1}} \text{ (below)}$$

Gradient to prior cell

$$\frac{dL}{dc^{t-1}} = \frac{dL}{dc^t} \cdot r^+ \quad \left. \begin{array}{l} \text{will need} \\ \text{to add} \\ \text{to the} \\ \text{next passes} \\ \frac{dL}{dc^t} \end{array} \right\}$$

$$\begin{aligned} \frac{dL}{dh^{t-1}} = & \frac{dL}{d\theta^+} \cdot \frac{d\theta^+}{dh^{t-1}} \\ & + \frac{dL}{dc_{temp}^+} \cdot \frac{dc_{temp}^+}{dh_{t-1}} \\ & + \frac{dL}{dr^+} \cdot \frac{dr^+}{dh_{t-1}} \\ & + \frac{dL}{dn^+} \cdot \frac{dn^+}{dh_{t-1}} \end{aligned} \quad \left. \begin{array}{l} \text{will save} \\ \text{as starting} \\ \text{point} \\ \text{for next} \\ \text{pass} \\ \text{back's} \\ \frac{dL}{dh^t} \end{array} \right\}$$

What data structs are
updated through each
pass back through
cell?

• Each cell backwards add

$$\frac{dL}{dw_k}, \frac{dL}{db_k}, \frac{dL}{du_k} \quad \Delta k$$

(add to 12 "global" variables)

• Pass $\frac{dL}{dh^{t-1}}, \frac{dL}{dc^{t-1}}$, but these

are over-written each iteration
backwards

How to efficiently
store all temporary
computations which are
local to a single
iteration through time?

• If we care about memory
can just globally store

x^t , h^t and recompute k^t
on the way back,

LSTM Intuition

Idea is that we start with hidden state and output of previous cell.

Then based on previous hidden state and current input token we compute vector of percentages

(0,1) of how much to remember (analogously, forget) of previous output. [remember]

We also create vector of percentages for how much

to incorporate value from
current iteration, **[new]**

We set current iteration using
[content]

We apply these percentages and
sum the two set the

output of the coll.

Moreover we have another

vector of percentages

of how much we want

the current output to

contribute to the next

hidden state **[output]**