

Forward Pass

$$h_t = f_h(x_t, h_{t-1}) = \phi_h(W h_{t-1} + U x_t + b)$$

take 0 when $t-1$ is (-1)
take 0 when $t+1$ is (T)

$$\hat{y}_t = f_o(h_t) = \phi_o(V h_t + c), \quad o_t = V h_t + c$$

Back propagation through Time

I assumed that $f_h()$ is $\tanh()$

$f_o()$ is softmax() (Usually this has been being used).

$$L(\hat{y}, y) = \sum_{t=1}^T L_t(\hat{y}_t, y_t)$$

$$= -\sum y_t \log \hat{y}_t$$

$$= -\sum y_t \log (\text{softmax}(o_t))$$

$$\text{SoftMax}(x_t) = S_t = \frac{e^{x_t}}{\sum e^{x_k}} \quad \text{for } t=1, \dots, k$$

Since softMax is a $\mathbb{R}^k \rightarrow \mathbb{R}^k$ mapping function, most general Jacobian matrix for

$$\frac{\partial S}{\partial x} = \begin{bmatrix} \frac{\partial S_1}{\partial x_1} & \dots & \frac{\partial S_1}{\partial x_k} \\ \vdots & \ddots & \vdots \\ \frac{\partial S_k}{\partial x_1} & \dots & \frac{\partial S_k}{\partial x_k} \end{bmatrix}$$

Let's compute $\frac{\partial S_i}{\partial x_j}$ for some arbitrary i and j :

$$\frac{\partial S_i}{\partial x_j} = \frac{\partial}{\partial x_j} \frac{e^{x_i}}{\sum_k e^{x_k}}$$

$$f(x) = \frac{g(x)}{h(x)} ; \quad f'(x) = \frac{g'(x)h(x) - g(x)h'(x)}{h(x)^2}$$

In our case $g_i = e^{x_i}$ and $h_i = \sum e^{x_k}$. No matter which x_j , when we compute the derivative of h_i with respect to x_j , the answer will always be e^{x_j} .

$$\frac{\partial}{\partial x_j} h_i = \frac{\partial}{\partial x_j} \sum e^{x_k} = \sum \frac{\partial e^{x_k}}{\partial x_j} = e^{x_j}$$

because $\frac{\partial e^{x_k}}{\partial x_j} = 0$ for $k \neq j$. There are only meaningful derivatives for $i=j$

in $\frac{\partial S}{\partial x}$ matrices for our problem.

Let's look at $i=j$ case

$$\begin{aligned} \frac{\partial}{\partial x_j} \frac{e^{x_i}}{\sum e^{x_k}} &= \frac{e^{x_i} \sum e^{x_k} - e^{x_i} e^{x_i}}{\left(\sum e^{x_k}\right)^2} \\ &= \frac{e^{x_i} \left(\sum e^{x_k} - e^{x_i} \right)}{\left(\sum e^{x_k}\right)^2} \\ &= \frac{e^{x_i}}{\sum e^{x_k}} \cdot \left(\frac{\sum e^{x_k}}{\sum e^{x_k}} - \frac{e^{x_i}}{\sum e^{x_k}} \right) \\ &= S_i (1 - S_i) \end{aligned}$$

Now we found what derivative of softmax, let's go back to the loss function,

$$L(\hat{y}, y) = - \sum y_t \log (\text{softmax}(o_t))$$

and we know that

$$\frac{\partial S_i}{\partial x_j} = S_i(1-S_j) \quad \forall i \neq j$$

$$\frac{\partial L}{\partial o_t} = - \frac{\partial}{\partial o_t} \sum y_t \log (\text{softmax}(o_t))$$

$$= - \sum y_t \frac{\partial}{\partial o_t} \log (\text{softmax}(o_t))$$

!

$$f(x) = \log_y x$$

$$f'(x) = \frac{x}{\ln x} \cdot \log_y e$$

$$= - \sum y_t \frac{\partial \log (\text{softmax}(o_t))}{\partial o_t}$$

$$= - \sum \frac{y_t \frac{\partial \text{softmax}(o_t)}{\partial o_t}}{\text{softmax}(o_t)} \cdot \log_e e$$

$$= - \sum \frac{\partial \hat{y}_t}{\partial o_t} \cdot \frac{y_t}{\hat{y}_t}$$

$$= - \sum S_t(1-S_t) \cdot \frac{y_t}{\hat{y}_t}$$

$$= - \sum (1-S_t) y_t$$

$$= - \sum (y_t - S_t \hat{y}_t)$$

$$= \sum S_t \hat{y}_t - \sum y_t$$

$$= S_t \sum \hat{y}_t - \sum y_t$$

$$= \sum \hat{y}_t - y_t$$

j represents predict value
i represents real value at a time

We know that i always equals j

$$\begin{bmatrix} x_1 \\ x_2 \\ x_3 \\ x_4 \end{bmatrix} \rightarrow \begin{bmatrix} 0.25 \\ 0.30 \\ 0.40 \\ 0.65 \end{bmatrix}$$

We need one derivative for one index to that one index to S_t

Now we found $\frac{\partial L}{\partial o} = \sum (\hat{y}_t - y_t)$

Back propagation goes backwards by propagating over functions

$$\left. \begin{array}{l} a_t = Wh_{t-1} + Ux_t + b \\ h_t = \tanh(a_t) \\ o_t = Vh_t + c \\ \hat{y}_t = \text{softmax}(o_t) \end{array} \right\} \begin{array}{l} \text{Encoder} \\ \text{Decoder} \end{array}$$

$$\begin{aligned} \frac{\partial L}{\partial V} &= \sum \frac{\partial L_t}{\partial V} \\ &= \sum \frac{\partial L_t}{\partial \hat{y}_t} \underbrace{\frac{\partial \hat{y}}{\partial o_t}}_{\hat{y} - y_t} \underbrace{\frac{\partial o_t}{\partial V}}_{\text{simple coef of } V} \\ &= \sum (\hat{y} - y_t) \cdot h_t \quad \text{(T) this is just for integrating for multi-hidden layers} \end{aligned}$$

$$\begin{aligned} \frac{\partial L}{\partial c} &= \sum \frac{\partial L_t}{\partial c} \\ &= \sum \frac{\partial L_t}{\partial \hat{y}_t} \underbrace{\frac{\partial \hat{y}_t}{\partial o_t}}_{\hat{y} - y_t} \underbrace{\frac{\partial o_t}{\partial c}}_{\$ \text{ its constant number}} \\ &= \sum (\hat{y} - y_t) \end{aligned}$$

$$\frac{\partial L_{(t+1)}}{\partial W} = \frac{\partial L_t}{\partial \hat{y}_t} \frac{\partial \hat{y}_t}{\partial h_t} \frac{\partial h_t}{\partial W}$$

h_{t-1}

Derivation of W also depends on h_{t-1} according to recursive function

($h_t = \tanh(W_t h_{t-1} + U_t x_t + b)$)

↓

$$\frac{\partial L_t}{\partial W} = \sum_k \frac{\partial L_t}{\partial \hat{y}_t} \frac{\partial \hat{y}_t}{\partial h_t} \left(\prod_{j=k}^{T-1} \frac{\partial h_{j+1}}{\partial h_j} \right) \frac{\partial h_K}{\partial W}$$

where

$$\prod_{j=k}^t \frac{\partial h_{j+1}}{\partial h_j} = \frac{\partial h_{t+1}}{\partial h_k} = \frac{\partial h_{t+1}}{\partial h_t} \frac{\partial h_t}{\partial h_{t-1}} \dots \frac{\partial h_{k+1}}{\partial h_k}$$

Thus at the time-step $t+1$, we can compute the gradient and further use back propagation through time from $t+1$ to 1 to compute the overall gradient with respect to W . We will finally yield the following gradient with respect to W :

$$\frac{\partial L_{t+1}}{\partial W} = \sum_{k=1} \frac{\partial L_{t+1}}{\partial \hat{y}_t} \cdot \frac{\partial \hat{y}_t}{\partial h_{t+1}} \cdot \frac{\partial h_{t+1}}{\partial h_k} \cdot \frac{\partial h_k}{\partial W}$$

Here we need to look at derivative of $\tanh()$ because The Gradient related to $f(x) = \tanh(x)$. Take the derivative of $\tanh(x)$ w.r.t. x

$$\frac{\partial \tanh(x)}{\partial x} = \frac{\partial \frac{\sinh(x)}{\cosh(x)}}{\partial x} = \frac{\cosh(x) - \sinh(x)}{\cosh(x)^2} \frac{\partial \cosh(x)}{\partial x} = \frac{\cosh(x)^2 - \sinh(x)^2}{\cosh(x)^2} = 1 - \tanh(x)^2$$

$$\frac{\partial L_{t+1}}{\partial w} = \sum_{k=1}^t \frac{\partial L_{t+1}}{\partial \hat{y}_t} \cdot \frac{\partial \hat{y}_t}{\partial h_{t+1}} \cdot \frac{\partial h_{t+1}}{\partial h_k} \cdot \frac{\partial h_k}{\partial w}$$

$$\frac{\partial L}{\partial w} = \sum_t \sum_{k=1}^{t+1} \frac{\partial L_{t+1}}{\partial \hat{y}_t} \frac{\partial \hat{y}_t}{\partial h_{t+1}} \frac{\partial h_{t+1}}{\partial h_k} \frac{\partial h_k}{\partial w}$$

Let's assume that $\frac{\partial L_{t+1}}{\partial h_{raw_k}} = \frac{\partial L_{t+1}}{\partial \hat{y}_t} \frac{\partial \hat{y}_t}{\partial h_{t+1}} \frac{\partial h_{t+1}}{\partial h_k}$

Let's try to simplify it.

$$\frac{\partial L}{\partial h_{raw_k}} = \sum_t \sum_{k=1}^{t+1} \frac{\partial L_{t+1}}{\partial h_{raw_k}}$$

$$= \sum_t \frac{\partial L_{t+1}}{\partial h_{raw_1}} + \frac{\partial L_{t+1}}{\partial h_{raw_2}} + \dots + \frac{\partial L_{t+1}}{\partial h_{raw_{t+1}}}$$

$$\equiv \left(\frac{\partial L_{t+1}}{\partial h_{raw_1}} + \frac{\partial L_{t+1}}{\partial h_{raw_2}} + \dots + \frac{\partial L_{t+1}}{\partial h_{raw_t}} + \frac{\partial L_{t+1}}{\partial h_{raw_{t+1}}} \right) +$$

$$\left(\frac{\partial L_t}{\partial h_{raw_1}} + \frac{\partial L_t}{\partial h_{raw_2}} + \dots + \frac{\partial L_t}{\partial h_{raw_t}} \right)$$

$$+ \left(\frac{\partial L_1}{\partial h_{raw_1}} \right)$$

from t+1 to 0
from T to 0

We unfolded these $2 \sum$ s. Lets try to fold column by column

instead of rowwise

$$= \sum_{t=1}^T \frac{\partial L_t}{\partial h_{row_1}} + \sum_{t=2}^T \frac{\partial L_t}{\partial h_{row_2}} + \sum_{t=3}^T \frac{\partial L_t}{\partial h_{row_3}} \dots + \frac{L_t}{\partial h_{row_e}}$$

$$= \sum_{i=1}^T \sum_{t=i}^T \frac{\partial L_t}{\partial h_{row_i}}$$

We define it as $T-1$ since we are sure that it will be simplified by last element of products

It's time to combine them.

What we found is:

$$\frac{\partial L}{\partial W} = \sum_{t=0}^T \sum_{k=1}^{t+1} \frac{\partial L_t}{\partial \hat{y}_t} \frac{\partial \hat{y}_t}{\partial h_t} \left(\prod_{j=k}^{T-1} \frac{\partial h_{j+1}}{\partial h_j} \right) \frac{\partial h_T}{\partial W}$$

also we can write

$$= \sum_t^T \frac{\partial h_{T-1}^{(\text{current state})}}{\partial W} \sum_{k=1}^t \frac{\partial L_t}{\partial \hat{y}_t} \frac{\partial \hat{y}_t}{\partial h_t} \frac{\partial h_t}{\partial h_k}$$

$$\sum_{t=i}^T \frac{\partial L_t}{\partial h_{row_i}}$$

$$\frac{\partial L}{\partial W} = \sum_t K_{t+1}^T \cdot \partial h_{next,t} \leftarrow \left((1-h_t^2) \cdot V^T \cdot \frac{\partial L}{\partial \hat{y}_t} + \partial h_{next,t+1} \right)$$

finally we yield the following gradient with respect to w

Finally we yield the following gradient w.r.t. to w

$$\frac{\partial L}{\partial w} = \sum_{t=1}^{t=T} h_{t-1}^T \cdot \text{next}_i()$$

(1 - h_t^2) $\cdot V^T \frac{\partial L}{\partial g_t}$

next_{i+1}

Let's work on to derive the gradient w.r.t. U . Similarly
further, we can take derivative over the whole sequence as:

We can take derivate over the whole sequence as

$$\frac{\partial L}{\partial U} = \sum_t^T \left(\frac{\partial L_{t+1}}{\partial \hat{y}_{t+1}} \frac{\partial \hat{y}_t}{\partial h_{t+1}} \left(\prod_{j=k}^{T-1} \frac{\partial h_{j+1}}{\partial h_j} \right) \frac{\partial h_t}{\partial U} \right)$$

$$= \sum_t^T \frac{\partial h_{T-t}^{(\text{current})}}{\partial U} \sum_{k=1}^{t+1} \frac{\partial L_t}{\partial \hat{y}_t} \frac{\partial \hat{y}_t}{\partial h_t} \frac{\partial h_t}{\partial h_k}$$

t
 $h_t^0, h_t^1, \dots, h_t^T$
 w_0, w_1, \dots, w_T
 $t=1, \text{too}$
 $t=T$
 X_t
 $\text{next}_i()$
 $\text{next}_i() = (l - h_t^2) \cdot V^T \cdot \sum_u$

Let's work on to derive the gradient w.r.t. b . Similarly, further, we can take derivative over the whole sequence as:

take derivate over the whole sequence

$$\frac{\partial L}{\partial b} = \sum_t \left(\sum_{k=1}^T \frac{\partial L_{t+1}}{\partial \hat{y}_{t+1}} \frac{\partial \hat{y}_t}{\partial h_{t+1}} \left(\prod_{j=k}^T \frac{\partial h_j}{\partial h_j} \right) \frac{\partial h_T}{\partial u} \right)$$

1 (its coef is 1.)

$$= \sum_{t=1}^{t=1, t=0} \text{next}_i() \left\{ \begin{array}{l} \text{next}_i() = (1-h_t^2) V^T \cdot \frac{\partial L}{\hat{y}_t} + \text{next}_{i+1}() \end{array} \right.$$

Algorithm RNN BPTT($x, h, \hat{y}, U, W, V, b, c, \eta, \beta$)

$dU, dW, dV = \text{zeros_like}(U, W, V)$

$db, dc = \text{zeros_like}(b, c)$

$dh_{\text{next}} = \text{zeros_like}(h[0])$

for t from $\text{length}(x[0])$ to 0 :

$$dy \leftarrow \hat{y} - t$$

$$dV \leftarrow dV + dy \otimes h_t^T$$

$$dc \leftarrow dc + dy$$

$$dh \leftarrow V^T \otimes dy + dh_{\text{next}}$$

$$dh_{\text{raw}} \leftarrow (1 - h_t^2) \cdot dh$$

$$dW \leftarrow dW + dh_{\text{raw}} \otimes h_{t-1}^T$$

$$dU \leftarrow dU + dh_{\text{raw}} \otimes x_t$$

$$db \leftarrow db + dh_{\text{raw}}$$

$$U = U - \eta \cdot dU$$

$$W = W - \eta \cdot dW$$

$$b = b - \eta \cdot db$$

$$V = V - \eta \cdot dV$$

$$c = c - \eta \cdot dc$$

return U, W, b, V, c