

# RNN Derivatives and Code sample

## By Mohit Kumar

# Recurrent Neural Networks:ex-1:BPTT:Derivatives

## Softmax Derivative

- It is fundamentally a vector function
  - So there is no derivative of softmax
  - Instead
    - Which component (output element) of softmax needs to be specified for derivative
    - Since softmax has multiple input, the derivative need to be calculated w.r.t which input element.
  - $\frac{\partial s_i}{\partial a_j}$  (partial derivative of  $i$ th output w.r.t  $j$ th input)
  - We say we compute its Jacobian Matrix.

# Recurrent Neural Networks: ex-1: BPTT: Derivatives

Softmax Derivative

$$\frac{\partial S_i}{\partial a_j} = \frac{\partial \frac{e^{a_j}}{\sum_{k=1}^N e^{a_k}}}{\partial a_j}$$

- For  $h_i$ , no matter for which  $a_j$  the derivative is computed, the answer will always be  $e^{a_j}$ .
- For  $g(i)$ , the derivative w.r.t  $a_j$  is  $\boxed{e^{a_j}}$  only if  $i=j$ , because only then  $g(i)$  has  $a_j$  anywhere in it. Otherwise the derivative is  $\boxed{0}$ .

Quotient Rule

$$f(x) = \frac{g(x)}{h(x)}$$

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

- in our case

$$g_i = e^{a_i}$$

$$h_i = \sum_{k=1}^N e^{a_k}$$

# Recurrent Neural Networks:ex-1:BPTT:Derivatives

• When  $i=j$  Softmax derivative

$$\begin{aligned}\frac{\partial \frac{e^{a_j}}{\sum_{k=1}^N e^{a_k}}}{\partial a_j} &= \frac{e^{a_i} \cdot \Sigma - e^{a_j} \cdot e^{a_i}}{\Sigma^2} \\ &= \frac{e^{a_i}}{\Sigma} \cdot \frac{\Sigma - e^{a_j}}{\Sigma} \\ &= s_i (1 - s_j)\end{aligned}$$

• When  $i \neq j$

$$\begin{aligned}\frac{\partial \frac{e^{a_j}}{\sum_{k=1}^N e^{a_k}}}{\partial a_j} &= \frac{0 - e^{a_j} e^{a_i}}{\Sigma^2} = -\frac{e^{a_j}}{\Sigma} \cdot \frac{e^{a_i}}{\Sigma} \\ &= -s_j s_i\end{aligned}$$

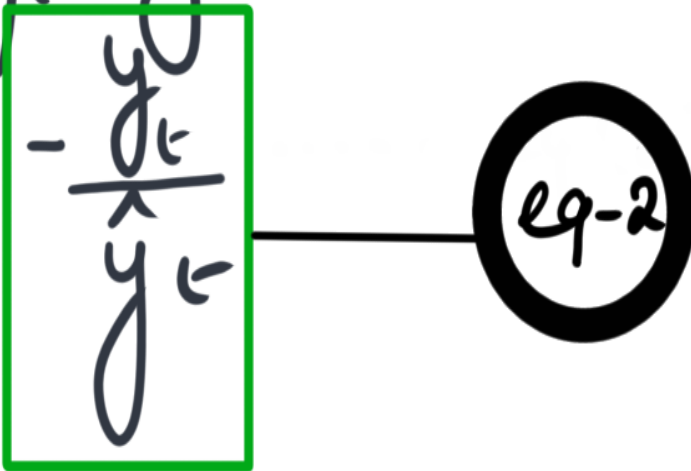
$$\begin{array}{ll} s_i (1 - s_j) & i=j \\ -s_j \cdot s_i & i \neq j \end{array}$$

eq-1

# Recurrent Neural Networks:ex-1:BPTT:Derivatives

$$\boxed{\frac{\partial E}{\partial v}} = \frac{\partial E_t}{\partial \hat{y}_t} \cdot \frac{\partial \hat{y}_t}{\partial q_t} \cdot \frac{\partial q_t}{\partial v} \quad q_t = v \cdot h_t$$

- $E = -y_t \cdot \log \hat{y}_t$  (cross entropy loss)

- $\frac{\partial E}{\partial \hat{y}_t} = -\frac{y_t}{\hat{y}_t}$  

- $\hat{y}$  is the softmax function

# Recurrent Neural Networks:ex-1:BPTT:Derivatives

$$\boxed{\frac{\partial E}{\partial v}} = \frac{\partial E_t}{\partial \hat{y}_t} \cdot \frac{\partial \hat{y}_t}{\partial q_t} \cdot \frac{\partial q_t}{\partial v}$$

*Combining 2 + 1*  $q_t = v \cdot \tanh_t$

$$\frac{\partial E}{\partial q_t} = \frac{\partial E_t}{\partial \hat{y}_t} \cdot \frac{\partial \hat{y}_t}{\partial q_t} = -\frac{y_t}{\hat{y}_t} \cdot \hat{y}_t (1 - \hat{y}_t) + \sum_{i=j} \left( -\frac{y_t}{\hat{y}_t} \right) (-\hat{y}_t \hat{y}_t)$$

$$= -y_t + y_t \hat{y}_t + \sum_{i=j} y_t \hat{y}_t$$

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

$$= \boxed{\hat{y}_t - y_t} \quad (\sum y_t = 1 \text{ because } y \text{ is one hot vector})$$

eq-3

# Recurrent Neural Networks:ex-1:BPTT:Derivatives

$$\boxed{\frac{\partial E}{\partial v}} = \frac{\partial E_t}{\partial \hat{y}_t} \cdot \frac{\partial \hat{y}_t}{\partial q_t} \cdot \frac{\partial q_t}{\partial v} \quad q_t = v \cdot h s_t$$

$$\boxed{\frac{\partial q_t}{\partial v} = \frac{\partial v \cdot h s_t}{\partial v} = h s_t} \quad \text{eq-4}$$

$$\frac{\partial E_t}{\partial \hat{y}_t} \cdot \frac{\partial \hat{y}_t}{\partial q_t} \cdot \frac{\partial q_t}{\partial v} = \hat{y} - y \cdot h s_t$$

$$\boxed{\frac{\partial E}{\partial v} = \hat{y}_t - y_t \cdot h s_t} \quad \text{eq-5}$$

# Recurrent Neural Networks:ex-1:BPTT:Derivatives

$$\frac{\partial E}{\partial u} = \frac{\partial E}{\partial \hat{y}_t} \cdot \frac{\partial \hat{y}_t}{\partial q_t} \cdot \frac{\partial q_t}{\partial h_{s_t}} \cdot \frac{\partial h_{s_t}}{\partial z_t} \cdot \frac{\partial z_t}{\partial u}$$

where  $q_t = v \cdot h_{s_t}$   
 $z_t = u \cdot x_t + w \cdot h_{s_{t-1}}$   
 $h_{s_t} = \tanh(z_t)$

$$\frac{\partial q_t}{\partial h_{s_t}} = \frac{\partial v \cdot h_{s_t}}{\partial h_{s_t}} = v$$

eq-6

$$\frac{\partial h_{s_t}}{\partial z_t} = (1 - h_{s_t}^2)$$

eq-7

$$\frac{\partial z_t}{\partial u} = \frac{\partial (u \cdot x_t + w \cdot h_{s_{t-1}})}{\partial u} = x_t$$

eq-8

$$\frac{\partial E}{\partial u} = \underbrace{(\hat{y}_t - y_t)}_{\text{eq-3}} \cdot \underbrace{v}_{\text{eq-6}} \cdot \underbrace{(1 - h_{s_t}^2)}_{\text{eq-7}} \cdot \underbrace{x_t}_{\text{eq-8}}$$

eq-9



# Recurrent Neural Networks:ex-1:BPTT:Derivatives

$$\frac{\partial E}{\partial \omega} = \frac{\partial E}{\partial \hat{y}_t} \cdot \frac{\partial \hat{y}_t}{\partial q_t} \cdot \frac{\partial q_t}{\partial h_{s_t}} \cdot \frac{\partial h_{s_t}}{\partial z_t} \cdot \frac{\partial z_t}{\partial \omega}$$

where  $q_t = v \cdot h_{s_t}$   
 $z_t = u \cdot x_t + \omega \cdot h_{s_{t-1}}$   
 $h_{s_t} = \tanh(z_t)$

$$\frac{\partial z}{\partial \omega} = \frac{\partial (u \cdot x + \omega \cdot h_{s_{t-1}})}{\partial \omega} = h_{s_{t-1}}$$

eq-10

$$\frac{\partial E}{\partial \omega} = \underbrace{(\hat{y} - y)}_{eq-3} \cdot \underbrace{v}_{eq-6} \cdot \underbrace{(1 - h_{s_t}^2)}_{eq-7} \cdot \underbrace{h_{s_{t-1}}}_{eq-10}$$

eq-11

# Recurrent Neural Networks:ex-1:BPTT:Derivatives

$$\frac{\partial E}{\partial h_{s_{t-1}}} = \frac{\partial E}{\partial \hat{y}_t} \cdot \frac{\partial \hat{y}_t}{\partial q_t} \cdot \frac{\partial q_t}{\partial h_{s_t}} \cdot \frac{\partial h_{s_t}}{\partial h_{s_{t-1}}}$$

where  $q_t = v \cdot h_{s_t}$   
 $g_t = u \cdot x_t + w \cdot h_{s_{t-1}}$   
 $h_{s_t} = \tanh(g_t)$

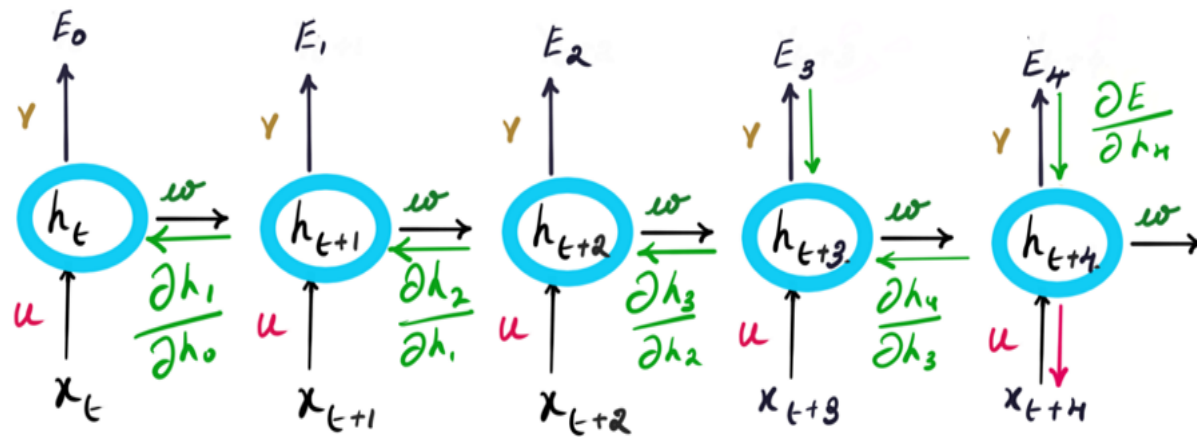
$$\frac{\partial g}{\partial h_{s_{t-1}}} = \frac{\partial (u \cdot x + w \cdot h_{s_{t-1}})}{\partial u} = w$$

eq-12

$$\frac{\partial E}{\partial h_{s_{t-1}}} = \underbrace{(\hat{y} - y)}_{eq-3} \cdot \underbrace{v}_{eq-6} \cdot \underbrace{(1 - h_{s_t}^2)}_{eq-7} \cdot \underbrace{w}_{eq-12}$$

eq-13

# Recurrent Neural Networks:ex-1:BPTT:Derivatives



$$h_t = \tanh(u \cdot x_t + w \cdot h_{t-1} + b_h)$$

$$y_t = \text{softmax}(v \cdot h_t)$$

$$\frac{\partial E}{\partial v} = \hat{y}_t - y_t \cdot h_{t+4} \quad \text{eq-5}$$

$$\frac{\partial E}{\partial u} = \underbrace{(\hat{y}_t - y_t)}_{\text{eq-3}} \cdot \underbrace{v}_{\text{eq-6}} \cdot \underbrace{(1 - h_{t+4}^2)}_{\text{eq-7}} \cdot \underbrace{x_t}_{\text{eq-8}} \quad \text{eq-9}$$

$$\frac{\partial E}{\partial w} = \underbrace{(\hat{y}_t - y_t)}_{\text{eq-3}} \cdot \underbrace{v}_{\text{eq-6}} \cdot \underbrace{(1 - h_{t+4}^2)}_{\text{eq-7}} \cdot \underbrace{h_{t-1}}_{\text{eq-10}} \quad \text{eq-11}$$

$$\frac{\partial E}{\partial h_{t-1}} = \underbrace{(\hat{y}_t - y_t)}_{\text{eq-3}} \cdot \underbrace{v}_{\text{eq-6}} \cdot \underbrace{(1 - h_{t+4}^2)}_{\text{eq-7}} \cdot \underbrace{w}_{\text{eq-12}} \quad \text{eq-13}$$

# Recurrent Neural Networks:ex-1:BPTT:Derivatives

```
du, dw, dv = np.zeros_like(self.u), np.zeros_like(self.w), np.zeros_like(self.v)
dbh, dby = np.zeros_like(self.bh), np.zeros_like(self.by)
dhnext = np.zeros_like(hs[0])
```

• calculate in reverse order.

```
for t in reversed(range(c)):
    label = labels[:, t * 1:(t + 1) * 1]
    dy = np.copy(ps[t])
    labelnhot = self.onehot(label)
    dy = dy - (labelnhot)
    dv += np.dot(dy, hs[t].T)
    dby += np.sum((dy), 1).reshape(self.numclasses, 1)
    dh = np.dot(self.v.T, dy) + dhnext # backprop into h
    dhraw = (1 - hs[t] * hs[t]) * dh # backprop through tanh nonlinearity
    dbh += np.sum((dhraw), 1).reshape(self.statesize, 1)
    du += np.dot(dhraw, xs[t].T)
    dw += np.dot(dhraw, hs[t - 1].T)
    dhnext = np.dot(self.w.T, dhraw)

for dparam in [dw, du, dv, dbh, dby]:
    np.clip(dparam, -5, 5, out=dparam) # clip to mitigate exploding gradients
self.state = hs[c - 1]
```

```
for param, dparam, mem in zip([self.w, self.u, self.v, self.bh, self.by], [dw, du, dv, dbh, dby],
                               [self.mw, self.mu, self.mv, self.mbh, self.mby]):
    mem += dparam * dparam
    param += -self.learningrate * dparam / np.sqrt(mem + 1e-8) # adagrad update
```

• Apply gradients.

$$\frac{\partial E}{\partial v} = (\hat{y}_t - y_t) \cdot h_{s_t} \quad \text{eq-5}$$

$$\frac{\partial E}{\partial u} = (\hat{y}_t - y_t) \cdot \underbrace{v}_{\text{eq-6}} \cdot \underbrace{(1 - h_{s_t}^2)}_{\text{eq-7}} \cdot \underbrace{x_t}_{\text{eq-8}} \quad \text{eq-9}$$

$$\frac{\partial E}{\partial w} = (\hat{y}_t - y_t) \cdot \underbrace{v}_{\text{eq-6}} \cdot \underbrace{(1 - h_{s_t}^2)}_{\text{eq-7}} \cdot \underbrace{h_{s_{t-1}}}_{\text{eq-10}} \quad \text{eq-11}$$

$$\frac{\partial E}{\partial h_{s_{t-1}}} = (\hat{y}_t - y_t) \cdot \underbrace{v}_{\text{eq-6}} \cdot \underbrace{(1 - h_{s_t}^2)}_{\text{eq-7}} \cdot \underbrace{w}_{\text{eq-12}} \quad \text{eq-13}$$

• softmax output or  $\hat{y}$   
 • dhnext is the portion that propagates back to step-0. 0 to begin with.