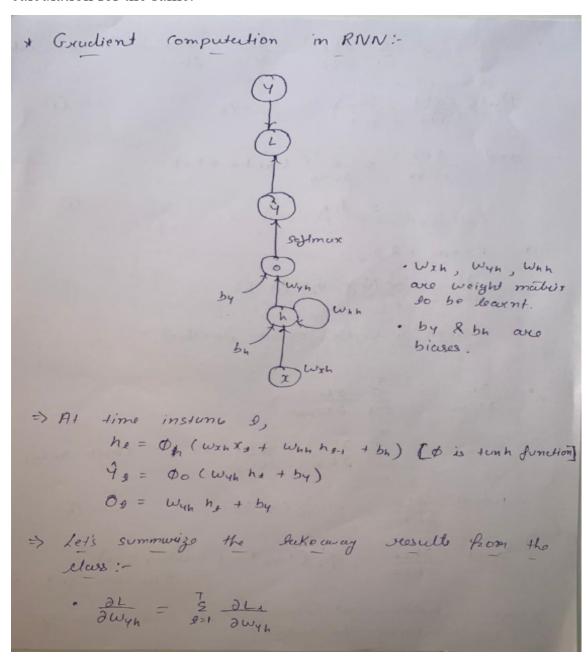
**Problem 1:** Derive expressions for  $\partial L/\partial bh$  and  $\partial L/\partial by$  for the RNN discussed in Lectures 2-3. Include the derived bias update equations in the RNN code shared and train the RNN for a sentence/word. Record relevant observations during training after adding bias terms.

**Solution:** To derive gradients with respect to the biases, I had used the result derived in the lectures. Below I have mentioned the complete calculation for the same.



But, 
$$\left(\frac{\partial L_{\pm}}{\partial \hat{Y}_{\pm}} \cdot \frac{\partial \hat{Y}_{\pm}}{\partial O_{\pm}} \cdot \frac{\partial O_{\pm}}{\partial W_{\pm}} \cdot \frac{\partial O_{\pm}}{\partial W_{\pm}}\right) = \hat{Y}_{1} - Y_{\pm} \qquad -- 0$$

and  $\frac{\partial O_{\pm}}{\partial W_{\pm}} = \frac{\partial}{\partial U_{\pm}} \left( \dot{W}_{\pm} \dot{W}_{\pm} + \dot{b}_{\pm} + \dot{b}_{\pm} \right)$ 

$$= \dot{h}_{\pm} \qquad -- 0$$

The entropy of the product of the p

(Derivations from class)

\* Now & let's use this results in olderwing 
$$\frac{\partial L}{\partial by}$$
 &  $\frac{\partial L}{\partial bh}$ .

•  $\frac{\partial L}{\partial by} = \frac{7}{9-1} \frac{\partial L}{\partial by}$  &  $\frac{\partial L}{\partial by}$  (: chain such)

• From eq 1,  $\left(\frac{\partial L}{\partial \dot{q}} + \frac{\partial \dot{q}}{\partial of}\right) = \dot{q}_1 - \dot{q}_1$ 

ond as  $O_4 = W_{4h} h_{+} + b_{4}$ 

•  $\frac{\partial O_7}{\partial by} = 1$ 

hence  $\frac{\partial L}{\partial by} = \frac{7}{9-1} \left(\dot{q}_4 - \dot{q}_4\right)$ 

(Derivation of  $\partial L/\partial by$ )

Now let's derive  $\partial L/\partial b$ h:

o 
$$\frac{\partial L}{\partial bn} = \frac{7}{9=1} \frac{\partial L_{\pm}}{\partial bn}$$

$$= \frac{5}{9=1} \frac{\partial L_{\pm}}{\partial h_{\pm}} \cdot \frac{\partial h_{\pm}}{\partial bn} \quad (:: Chain scale)$$

$$\Rightarrow Fxom eq 4, \frac{\partial L_{\pm}}{\partial h_{\pm}} = -w_{yh}^{T} (y-9)$$

And 
$$h_{4} = \Phi h \left( \chi \omega_{\chi_{h}} \chi_{4} + \omega_{h_{h}} h_{4-1} + b_{h} \right)$$

$$\frac{\partial h_{+}}{\partial b_{h}} = \Phi_{h}^{i} \left( \omega_{\chi_{h}} \chi_{4} + \omega_{h_{h}} h_{4-1} + b_{h} \right).$$

$$\frac{\partial}{\partial b_{h}} \left( \omega_{\chi_{h}} \chi_{4} + \omega_{h_{h}} h_{4-1} + b_{h} \right).$$

$$\frac{\partial}{\partial b_{h}} \left( \omega_{\chi_{h}} \chi_{4} + \omega_{h_{h}} h_{4-1} + b_{h} \right) = 1$$

$$\frac{\partial h_{+}}{\partial b_{h}} = \Phi_{h}^{i} \left( \omega_{\chi_{h}} \chi_{4} + \omega_{h_{h}} h_{4-1} + b_{h} \right).$$

$$\frac{\partial h_{+}}{\partial b_{h}} = \Phi_{h}^{i} \left( \omega_{\chi_{h}} \chi_{4} + \omega_{h_{h}} h_{4-1} + b_{h} \right).$$

$$\frac{\partial h_{+}}{\partial b_{h}} = \Phi_{h}^{i} \left( \omega_{\chi_{h}} \chi_{4} + \omega_{h_{h}} h_{4-1} + b_{h} \right).$$

$$\frac{\partial h_{+}}{\partial b_{h}} = \Phi_{h}^{i} \left( \omega_{\chi_{h}} \chi_{4} + \omega_{h_{h}} h_{4-1} + b_{h} \right).$$

$$\frac{\partial h_{+}}{\partial b_{h}} = \Phi_{h}^{i} \left( \omega_{\chi_{h}} \chi_{4} + \omega_{h_{h}} h_{4-1} + b_{h} \right).$$

$$\frac{\partial h_{+}}{\partial b_{h}} = \Phi_{h}^{i} \left( \omega_{\chi_{h}} \chi_{4} + \omega_{h_{h}} h_{4-1} + b_{h} \right).$$

$$\frac{\partial h_{+}}{\partial b_{h}} = \Phi_{h}^{i} \left( \omega_{\chi_{h}} \chi_{4} + \omega_{h_{h}} h_{4-1} + b_{h} \right).$$

$$\frac{\partial h_{+}}{\partial b_{h}} = \Phi_{h}^{i} \left( \omega_{\chi_{h}} \chi_{4} + \omega_{h_{h}} h_{4-1} + b_{h} \right).$$

$$\frac{\partial h_{+}}{\partial b_{h}} = \Phi_{h}^{i} \left( \omega_{\chi_{h}} \chi_{4} + \omega_{h_{h}} h_{4-1} + b_{h} \right).$$

$$\frac{\partial h_{+}}{\partial b_{h}} = \Phi_{h}^{i} \left( \omega_{\chi_{h}} \chi_{4} + \omega_{h_{h}} h_{4-1} + b_{h} \right).$$

(Derivation of  $\partial L/\partial bh$ )

Now we need to add this portion in to our code. So below I have mentioned all the changes we need to do to add biases.

## **Steps:**

• First initialize bias variables with respective size in init() method of RNN class.

```
#Added biases
self.bh = np.zeros((hidden_size, 1)) # hidden bias
self.by = np.zeros((vocab_size, 1)) # output bias
```

• Add bias terms in calculation in forward pass.

```
hs[t] = np.tanh(np.dot(self.W_hh,hs[t-1]) + np.dot(self.W_xh,xs[t]) + self.bh)
os[t] = np.dot(self.W yh,hs[t]) + self.by
```

• To calculate gradients, first declare gradient variables of appropriate size in backward() method of RNN class.

```
#Gradient for biases
dL_dbh = np.zeros_like(self.bh)
dL_dby = np.zeros_like(self.by)
```

• Using the expressions derived above, we will do gradient calculations. Add the below code in the loop to calculate  $\partial L/\partial by$ .

```
#Compute yhat - y.
d_yy_cap = np.copy(ycap[t])
d_yy_cap[targets[t]]-= 1

# Add this line to calculate gradient with repect to dby.
dL_dby+= d_yy_cap
```

• To calculate  $\partial L/\partial bh$  we will need derivation of tanh (activation function).

```
# Compute W^Tyh(y-ycap)
dL_dh = np.dot(self.W_yh.T,d_yy_cap) + dhnext
# Diffrentiaion of tanh(x) is (1-tanh^2(x))
dL_dh_dtanh = (1 -hs[t]*hs[t])*dL_dh

#dL_dbh is W^Tyh(y-ycap) * derivation of tanh(hs) so add below line to compute dL_dbh.
dL_dbh+= dL_dh_dtanh
```

Comparison of result with same configuration but without and with biases:

**Input:** "Amazon is river."

Output without biases: It did converged at the last 10,000 iterations.

```
iter num 9000 [0.41933032]
Amazon in river.
None
```

Output with biases: It converged at 8000 iterations and loss decreased to 0.3035.

iter num 9000 [0.30353884] Amazon is river. **Problem 2:** Replace basic sgd technique used in the function update\_model() with any other sophisticated gradient update technique popular in literature. Record relevant observations during training after modifying gradient update method.

#### **Solution:**

I have tried 2 gradient update technique namely AdaGrad and RMSprop.

### AdaGrad:

AdaGrad (Adaptive Gradient) does not update directly based on learning rate but it actually maintains the running sum of previous gradients and at time of update it divides the learning rate by this running sum. And to avoid division by zero a small offset is used.

Gradient update equation:

$$w_{t+1} = w_t - \frac{\alpha}{\sqrt{v_t + \epsilon}} \cdot \frac{\partial L}{\partial w_t}$$

Where vt is,

$$v_t = v_{t-1} + \left[\frac{\partial L}{\partial w_t}\right]^2$$

### **Code for AdaGrad:**

First declare variables to store running sum of gradients.

```
# memory variables
self.mW_xh = np.zeros_like(self.W_xh)
self.mW_hh = np.zeros_like(self.W_hh)
self.mW_yh = np.zeros_like(self.W_yh)
self.mbh = np.zeros_like(self.bh)
self.mby = np.zeros_like(self.by)
```

Call adaGrade() function from update\_method.

## RMSprop:

In RMSprop(Root mean square prop) is similar to adaGrad() but the difference is in RMSprop weights are assigned to current gradient and running sum of gradients.

Gradient update equations:

$$w_{t+1} = w_t - \frac{\alpha}{\sqrt{v_t + \epsilon}} \cdot \frac{\partial L}{\partial w_t}$$

Where,

$$v_t = \beta v_{t-1} + (1 - \beta) \left[ \frac{\partial L}{\partial w_t} \right]^2$$

# **Code for RMSprop:**

Call RMSprop() method from update\_method.

# Comparison of SGD, AdaGrad and RMSprop:

**Input:** "All the faith he had had had had no effect on the outcome of his life."

**Output with SGD:** Even after 10000 iterations, the loss remained very high and it didn't converged.

```
iter num 9000 [128.40807327]
Alemfeud hlf e d ihd nefeelh ta d hetiad nttc ihefiad h i ilhfohe ust
```

**Output with AdaGrad:** For AdaGrad too, the loss remained very high after 10000 iterations.

```
iter num 9000 [142.44594829]
Adi nooltoe s ttfahh ofsh dhhfhutahl uheAioec efnadhna otifls dcsie
```

**Output with RMSprop:** It converged very early i.e at 3000 iteration it started converging and till last iteration loss decreased to 0.07.

```
iter num 2000 [0.85262771]
All the faith he had had had had no effect on the outcome of his life
iter num 3000 [0.28077867]
All the faith he had had had had no effect on the outcome of his life
```

**Input:** "Amazon is river."

**Output with SGD:** It converged at 10000 iterations and loss decreased to 0.3035.

```
iter num 9000 [0.30353884]
Amazon is river.
```

**Output with AdaGrad:** For AdaGrad, the loss remained very high after 10000 iterations and it did not converged.

```
iter num 9000 [25.3674136]
AvoozAvos. iv vov
```

**Output with RMSprop:** It converged very early i.e at 2000 iteration with very less loss compared to SGD.

```
iter num 1000 [0.09690172]
Amazon is river.
```

# **Observation:**

From the above results it is clear that performance of RMSprop is best compared to SGD and AdaGrad. Whereas performance of SGD is better compared to AdaGrad.

**Problem 3:** Experiment with various hidden vector sizes and record your observation.

### **Solution:**

**Input:** "Amazon is river."

# Parameter configuration:

• Learning rate: 0.001

• Iterations: 10000

Gradient update method: SGD

• Hidden vector size: 10

## **Output:**

```
iter num 9000 [0.30353884]
Amazon is river.
```

With same input and configuration but only hidden vector size increased to 20, the result obtained is as follows:

## **Output:**

```
iter num 5000 [0.30995026]
Amazon is river.
iter num 6000 [0.24641345]
Amazon is river.
```

# Now increasing hidden vector size to 30,

```
iter num 2000 [0.72503865]
Amazon is river.
```

It started predicting the output correctly very early but with high loss.

# **Trying with other input:**

**Input:** "All the faith he had had had no effect on the outcome of his life"

**Parameters:** Same parameters just hidden vector size is increased to 30.

### **Output:**

iter num 2000 [1.64148478]
All the faith he had had had had no effect on the outcome of his life
iter num 3000 [0.883951]
All the faith he had had had had no effect on the outcome of his life

As we saw above, earlier with hidden vector size as 10 and gradient update method as SGD, it was not converging even after 10000 iterations but with hidden vector size increased to 30, it easily converged at very early iteration and till 10000 iteration the loss decreased to 0.21.

### **Observation:**

From above results we can see that, with increase in size of hidden vectors the result obtained are getting better and this is because as we increase hidden vector size, the RNN is able to store more contextual data and that helps in predicting the next character.