

### Assignment 1

Shourabh Payal (MT2020054)

1. Derive expressions for  $\frac{\partial L}{\partial b_h}$  and  $\frac{\partial L}{\partial b_y}$  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 word. Record relevant observations during training after adding bias terms.

We have

$$L = -y_t \log(\hat{y}_t)$$

$$\hat{y}_t = \text{softmax}(z_t)$$

$$z_t = W_{yh} h_t + b_y$$

(a) Derivative wrt  $b_y$

$$\begin{aligned} \frac{\partial L}{\partial b_y} &= \sum_{i=1}^T \frac{\partial L_i}{\partial b_y} \\ &= \sum_{i=1}^T \frac{\partial L_i}{\partial \hat{y}_i} \frac{\partial \hat{y}_i}{\partial z_i} \frac{\partial z_i}{\partial b_y} \end{aligned}$$

$$\bullet \frac{\partial L_i}{\partial \hat{y}_i} = -\frac{y_i}{\hat{y}_i}$$

$$\begin{aligned} \bullet \frac{\partial \hat{y}_i}{\partial z_i} &= \hat{y}_i(1 - \hat{y}_i) \text{ when } i = k \\ &= -\hat{y}_i \hat{y}_k \text{ when } i \neq k \end{aligned}$$

$$\bullet \frac{\partial z_i}{\partial b_y} = 1$$

- Combining and solving we get

$$\frac{\partial L}{\partial b_y} = \sum_{i=1}^T (\hat{y}_i - y_i)$$

---

(b) Derivative wrt  $b_h$

Let's consider for time  $t+1$  first

$$\frac{\partial L_{t+1}}{\partial b_h} = \frac{\partial L_{t+1}}{\partial h_{t+1}} \cdot \frac{\partial h_{t+1}}{\partial b_h}$$

We know from derivations of  $\frac{\partial L_{t+1}}{\partial W_{xh}}$  and  $\frac{\partial L_{t+1}}{\partial W_{hh}}$  that

$$\frac{\partial L_{t+1}}{\partial h_{t+1}} = -W_{yh}^T (y_t - \hat{y}_t)$$

Now  $\frac{\partial h_{t+1}}{\partial b_h} = ?$

Let  $z_{t+1} = W_{xh}X_{t+1} + W_{hh}h_t + b_h$

We know  $h_{t+1} = \phi_h(z_{t+1})$

$$\therefore \frac{\partial h_{t+1}}{\partial b_h} = \phi'_h(z_{t+1}) \cdot \frac{\partial z_{t+1}}{\partial b_h}$$

$$\frac{\partial z_{t+1}}{\partial b_h} = W_{hh} \frac{\partial h_t}{\partial b_h} + 1$$

$$\frac{\partial h_{t+1}}{\partial b_h} = \phi'_h(z_{t+1}) \cdot (W_{hh} \frac{\partial h_t}{\partial b_h} + 1) \quad (\text{Reccursion !})$$

$$\frac{\partial h_t}{\partial b_h} = \phi'_h(z_t) \cdot (W_{hh} \frac{\partial h_{t-1}}{\partial b_h} + 1)$$

and so on...

Total gradient on all losses wrt  $b_h = \sum_{k=1}^T \frac{\partial L_k}{\partial b_h}$

---

(c) Observations after adding bias term

```
for bias in biases:  
    plt.plot(bias_dic[bias]['iter'],bias_dic[bias]['loss'], label=bias)  
    plt.xlabel('iterations')  
    plt.ylabel('loss')  
    plt.legend(loc="upper right")
```

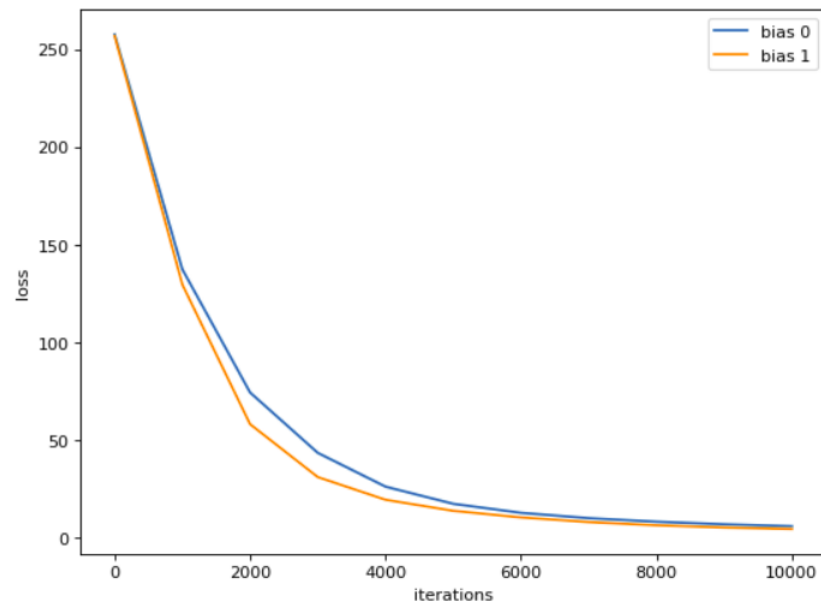


Figure 1: Bias comparison

Did not observe significant difference before and after adding bias.

- 
2. Replace the 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.

### AdaGrad, AdaDelta, Adam

Let's consider a simple high level comparison between Basic SGD and Adagrad using their respective update equations.

#### Basic SGD:

$$W_t = W_{t-1} - \eta \frac{\partial L}{\partial W_{t-1}}$$

Here  $\eta$  is learning rate and it is same for all weights and is constant.

#### AdaGrad:

$$W_t = W_{t-1} - \eta'_t \frac{\partial L}{\partial W_{t-1}}$$

Here  $\eta'_t$  is adaptive learning rate and it is different for each weight and changes each iteration.

$$\eta'_t = \frac{\eta}{\sqrt{\alpha_{t-1} + \epsilon}}$$

$\epsilon$  is a small positive value to avoid division by zero problem.

$$\alpha_{t-1} = \sum_{i=1}^{t-1} \left( \frac{\partial L}{\partial W_{i-1}} \right)^2$$

#### Advantages of Adagrad

- No need of manually tuning  $\eta$
- Takes care of right learning rate for sparse and dense features

#### Disadvantages of Adagrad

- $\alpha_{t-1}$  can become very large as  $t$  increases as a result convergence can be slowed down.

---

## AdaDelta

To overcome the disadvantage of adagrad we imply a simple technique of exponential decaying averages.

Instead of  $\alpha_{t-1}$  we have  $EDA_{t-1}$

$$EDA_{t-1} = (\gamma)EDA_{t-2} + (1 - \gamma) \left( \frac{\partial L}{\partial W_{t-1}} \right)^2$$

Growth of the denominator term in  $\eta'_t$  is therefore controlled.

## Adam (Adaptive moment estimation)

Adam tries to bring together the above methods along with the concept of moments.

- Mean : First order moment
- Variance : Second order moment

Let's define  $g_t = \frac{\partial L}{\partial W_{t-1}}$

$$m_t = \beta_1(m_{t-1}) + (1 - \beta_1)g_t$$

$$v_t = \beta_2(v_{t-1}) + (1 - \beta_2)g_t^2$$

$$\hat{m}_t = \frac{m_t}{1 - (\beta_1)^2}$$

$$\hat{v}_t = \frac{v_t}{1 - (\beta_2)^2}$$

Here  $0 \leq \beta_1, \beta_2 \leq 1$

Typical good values are  $\beta_1 = 0.9$  and  $\beta_2 = 0.999$

Our final update equation is

$$W_t = W_{t-1} - \eta \left( \frac{\hat{m}_t}{\sqrt{\hat{v}_t + \epsilon}} \right)$$

## Results on code

We will be implementing Adam optimizer. Below are the results of experiment.

```
***** Optimizer = sgd *****
iter num: 0   loss: 256.32902761937254
acJKZS6jZBKLthj6i0VHmpXZFMT97CsSV8RezCa4HBUswjHZYNE

iter num: 1000 loss: 125.49305610208975
aH90FfMQqjgFKopeRVmmsgdTnkhmT 7qr9UNjSz z1J5RFxWSSh

iter num: 2000 loss: 44.99013560354285
aq9BNfkptudhKVdA0Zhm0kJRnktmWmpn5gtr6tFTP64z3J89EQ0

iter num: 3000 loss: 21.38854662109456
aKcdKughijcJKVveKcdefghmn5Iqbkhijclujkjmjdpmsgsfpgpb

iter num: 4000 loss: 13.041387544632776
a9cdefcheqghtfvwefghifghzjMeurhePghi0kpqrohK9MNOPLi

iter num: 5000 loss: 9.071217496881347
ajghijklmnopqrst5y 67z9UXyX3W56781 3456709 e056789

iter num: 6000 loss: 6.8450964052604695
a0khefghijklmnopqrstuv5xyzABCDEQGHijklmnopQ7STUVsX3

iter num: 7000 loss: 5.454772212054938
abcdefghijklmnohQHWtjkmndpqrhtfghijklmnopqrstuvw

iter num: 8000 loss: 4.515745324216776
abcdiRcdefgXirghULgtbndefghijklmnopqrstuvwxyzABCDEL

iter num: 9000 loss: 3.8435387761963145
aKcdefghijklmnopqrstuvwxyzABCDEFGHijklmnopQRSTUVMXY

iter num: 10000 loss: 3.3405925170849944
abchefghijklmnopqrstuvwxyzABCDEFGHijklmnOSQRSTUVWXY

***** Optimizer = adam *****
iter num: 0   loss: 256.92660123981443
aBJX9qhle0 q8FZBB63aFRZeigHqyqDG9CQqtUjx2CzwZnXB2oFQbdZa00zi6e

iter num: 1000 loss: 17.032066685685322
abcdeghijklCDAB0DJKLnopVRRSTUkXx9Z01X9bcdffqBCDY6BH9JKLt9Z0wWXYZ

iter num: 2000 loss: 0.129939943308596
abcdeghijklmnopqrstuvwxyABCDEFHijklmnopQRSTUVWXYZ0123456789

iter num: 3000 loss: 0.005504777963686187
abcdeghijklmnopqrstuvwxyABCDEFHijklmnopQRSTUVWXYZ0123456789

iter num: 4000 loss: 0.0017524094885176364
abcdeghijklmnopqrstuvwxyABCDEFHijklmnopQRSTUVWXYZ0123456789

iter num: 5000 loss: 0.0009855137710198826
abcdeghijklmnopqrstuvwxyABCDEFHijklmnopQRSTUVWXYZ0123456789

iter num: 6000 loss: 0.0006776010023463483
abcdeghijklmnopqrstuvwxyABCDEFHijklmnopQRSTUVWXYZ0123456789

iter num: 7000 loss: 0.0005142756096922056
abcdeghijklmnopqrstuvwxyABCDEFHijklmnopQRSTUVWXYZ0123456789

iter num: 8000 loss: 0.0004136866833997377
abcdeghijklmnopqrstuvwxyABCDEFHijklmnopQRSTUVWXYZ0123456789

iter num: 9000 loss: 0.0003457091330526938
abcdeghijklmnopqrstuvwxyABCDEFHijklmnopQRSTUVWXYZ0123456789

iter num: 10000 loss: 0.0002967732773813102
abcdeghijklmnopqrstuvwxyABCDEFHijklmnopQRSTUVWXYZ0123456789
```

(a) Normal SGD

(b) Adam

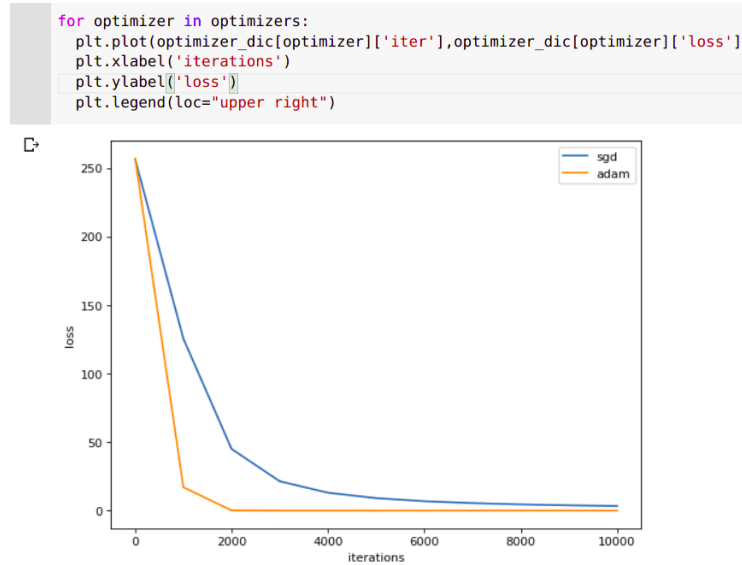


Figure 3: Optimizer comparison

The Adam optimizer converges a lot quickly compared to the regular SGD for this problem.

### 3. Experiment with various hidden vector sizes and record your observations

```
for layer in layers:
    plt.plot(layer_dic[layer]['iter'], layer_dic[layer]['loss'], label=
    plt.xlabel('iterations')
    plt.ylabel('loss')
    plt.legend(loc="upper right")
```

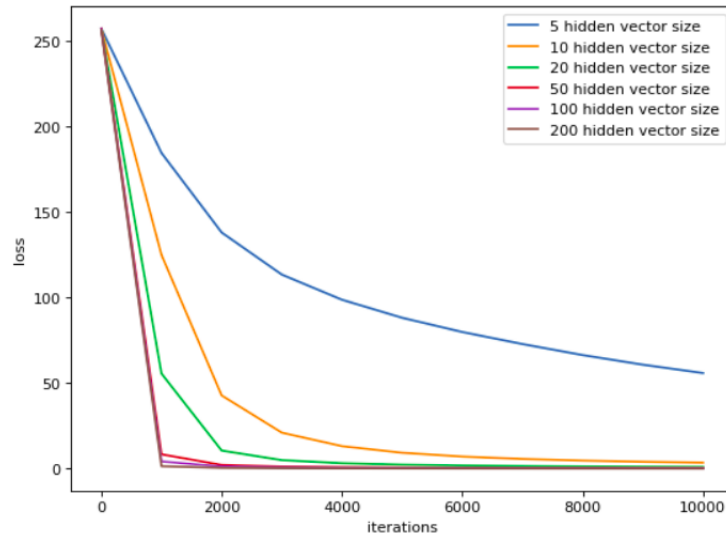


Figure 4: Hidden layer vector size comparison using SGD optimizer

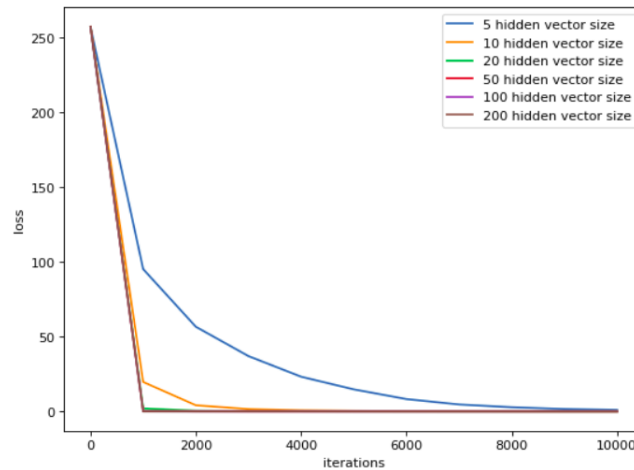


Figure 5: Hidden layer vector size comparison using Adam optimizer