Mohamed Ehab Fathy        2002597

## Assignment 3

1. 1) $h_1 = w_x x_1 + w_h \cdot h_0 = 0.1 \times 10 + 1 \times 1$

$$= 2$$

$h_2 = w_h \cdot h_1 + w_x \cdot x_2 = 1 \times 2 + 0.1 \times 10$

$$= 3$$

$\hat{y}_2 = w_y \cdot h_2 = 2 \times 3 = 6$ #

2) $L_t \Longrightarrow$ For 2 outputs

$\hat{y}_2 = 6$

$\hat{y}_1 = w_h \cdot h_1 = 2$

$L_t = \sum_{i=1}^{2} (\hat{y}_i - y_i) = 1 + 9 = 10$

3) $\dfrac{\partial L_t}{\partial h_1} = \dfrac{\partial L_1}{\partial h_1} + \dfrac{\partial L_2}{\partial h_1}$

* $\dfrac{\partial L_1}{\partial h_1} = \dfrac{\partial L_1}{\partial y_1} * \dfrac{\partial y_1}{\partial h_1}$

$= 2(\hat{y}_1 - y_1) * w_y$

$= -6 * 2 = -12$

* $\dfrac{\partial L_2}{\partial h_1} = \dfrac{\partial L_2}{\partial y_2} * \dfrac{\partial y_2}{\partial h_2} * \dfrac{\partial h_2}{\partial h_1}$

$= 2(\hat{y}_2 - y_2) * w_y * w_h$

$= 2 * 2 * 1 = 4$

$\therefore \dfrac{\partial L_t}{\partial h_1} = -8$

4) $\dfrac{\partial L_t}{\partial w_h} = \dfrac{\partial L_1}{\partial w_h} + \dfrac{\partial L_2}{\partial w_h}$

$* \quad \dfrac{\partial L_1}{\partial w_h} = \dfrac{\partial L_1}{\partial y_1} * \dfrac{\partial y_1}{\partial h_1} * \dfrac{\partial h_1}{\partial w_h}$

$\qquad = 2(\hat{y}_1 - y_1) * w_y * h_o$

$\qquad = -6 * 2 * 1 = -12$

$* \quad \dfrac{\partial L_2}{\partial w_h} = \overset{o-}{\underset{i=2}{\sum}} \dfrac{\partial L_2}{\partial y_2} * \dfrac{\partial y_2}{\partial h_2} \dfrac{\partial h_2}{\partial h_i} * \dfrac{\partial h_i}{\partial w_h}$

$\qquad = 2(\hat{y}_2 - y_2) * w_y \left[ \dfrac{\partial h_2}{\partial w_h} + \dfrac{\partial h_2}{\partial h_1} * \dfrac{\partial h_1}{w_h} + \dfrac{\partial h_2}{\partial h_o} * \dfrac{\partial h_o}{\partial w_h} \right]$

$\qquad = 2(\hat{y}_2 - y_2) * w_y \left[ h_1 + w_h * h_o \right]$

$\qquad = 2 * 2 * [2 + 1] = 8$

$$\therefore \frac{\partial L_t}{\partial W_n} = -4$$

2    * to observe   this , we must check

   the impact  of  the  First  input  on   the

   last  outpt

$$\frac{\partial J_n}{\partial w} = \sum_{i=0}^{n} \frac{\partial J_n}{\partial y_n} * \frac{\partial y_n}{\partial h_n} * \left( \frac{\partial h_n}{\partial h_i} * \frac{\partial h_i}{\partial w} \right)$$

@ $i=0$ } ⟶ First input

$$\frac{\partial h_n}{\partial h_o} = \frac{\partial h_n}{\partial h_{n-1}} * \frac{\partial h_{n-1}}{\partial_{n-2}} * \cdots * \frac{\partial h_1}{\partial h_o}$$

and $\dfrac{\partial h_n}{\partial_{n-1}} = W_{hh} * \tanh'(W_{hh} h_{n-1}, W_{xh} x_n)$

* where $W_{hh}$  is  probably  sampled  From  unit

* Hi .

gaussian $(<1)$, and derivative of $tanh \leq 1$

, so we keep multiplying small numbers

, so $\frac{\partial h_n}{\partial h_0}$ becomes very small no., hence

$h_0$ has no impact on $y_n$ (Vanishing gradients)

(short memory)

3] For long sequences, where Vanishing gradients must be avoided, so we can keep track of long term info.

4] Adv: using as training algorithm which is applied to sequence data, where RNN is used (since weights are shared for each time step, then errors are calculated and accumulated for each time step to update single

weight)

* disadv: high Computational Cost for Single parameter update

5) a) feed forward NN will assume that prediction of each letter depends only on that letter, which isn't the Case, since we have to learn from the whole sequence before predicting one letter

b) word embedding where each letter is encoded to hot vector

c) many to many (encoder-decoder) since we need to watch the whole sequence before output the First letter
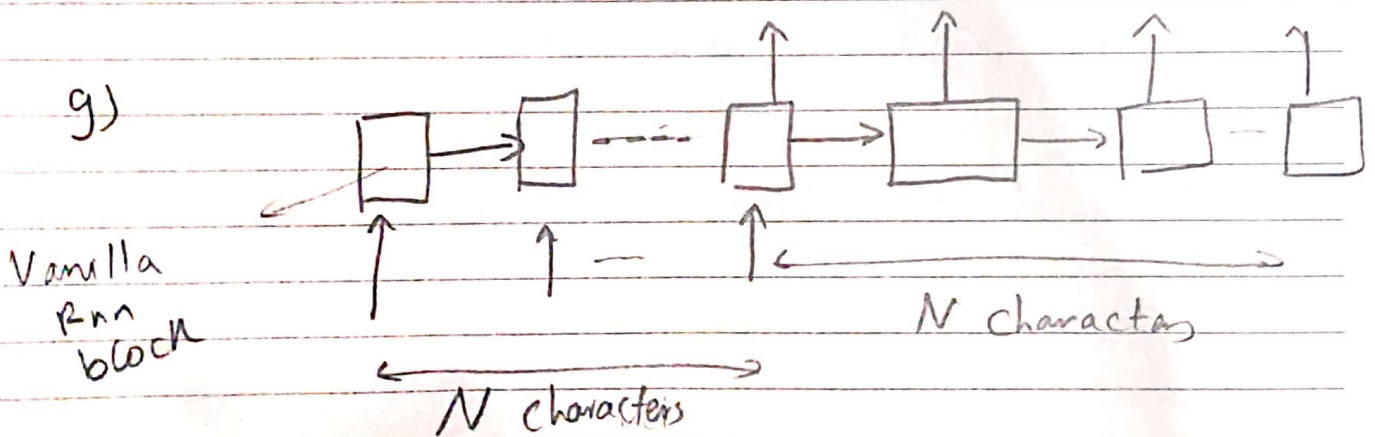
e) * choose similar length texts to be in the same patch

* use padding

* or choose patch size of 1

f) Convert each charchter into hot vector

* divide training data into patches, same text lengths per patch

g)

Vanilla
Rnn
block

N characters

N characters

* enCoder deCoder arch

h. → decode the hot vectors into characters

the concatanate them