

1. Questions

3.

- (a) ANN has no idea of a sequence. It takes input one by one and does not store information about the past. In RNNs input comes as a stream/sequence, it processes them step by step and it also remembers what it saw in the past.
- (b) RNNs struggle with long term dependencies because of the issue of vanishing gradients.
- (c) LSTM's uses gates [forget gate, input gate, output gate] and using these information can flow unchanged ~~time~~.
- (d) RNNs repeatedly multiply gradients by a no. $(\eta) < 1$ so they shrink with time/many inputs. LSTM's using gates have additive updates.
- (e) ANN, simple pre-diction
RNN, sentence completion
- LSTM's talking bots (small scale)

4. (a) The content of the word. Eg, the use of the word model can mean differently in diff contexts, so long term dependency is reqd. No, RNNs struggle for long range dependency, LSTMs are useful.

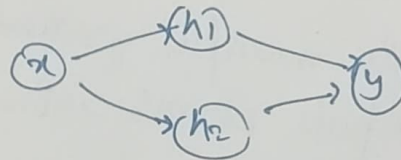
(b) LSTMs have a cell state that flows forward almost unchanged if has gates, which learn controlled values

- 1- forget Gate: if the old info still relevant, if not, forget it.
- 2- Input Gate: is the new info to be remembered or not?
- 3- Output Gate: Does the remembered content need to be applied while generating the next output or not?

forget gate 0 \rightarrow erase memory

when the content of a word changes in the sentence / text it should forget the old one and remember the new one.

2.
$$\sigma(z) = \frac{1}{1+e^{-z}}$$



(c) The idea behind Backprop is to compute the gradients efficiently, and then update the previous weight and biases, so as to reduce the cost function, and make the predictions closer to the actual output.

Each parameter affects the other one in a link, so applying the chain rule helps us to model all these dependencies.

(b) (i) $\frac{\partial L}{\partial w_2}$, using chain rule $\frac{\partial L}{\partial w_2} = \frac{\partial L}{\partial z_2} \cdot \frac{\partial z_2}{\partial w_2}$

$$\frac{\partial L}{\partial w_2} = (a_2 - y)(a_1)$$

$$(i) \frac{\partial L}{\partial b_2} = \frac{\partial L}{\partial z_2} \cdot \frac{\partial z_2}{\partial b_2} = (a_2 - y)(1)$$

$$z_2 = w_2 a_1 + b_2$$

$$z_1 = w_1 x + b_1$$

$$a_1 = \sigma(z_1)$$

$$(i) \frac{\partial L}{\partial w_1} = \frac{\partial L}{\partial z_2} \cdot \frac{\partial z_2}{\partial a_1} \cdot \frac{\partial a_1}{\partial z_1} \cdot \frac{\partial z_1}{\partial w_1}$$

$$\Rightarrow (a_2 - y)(w_2)(a_1)(1 - a_1)(x)$$

$$(i) \frac{\partial L}{\partial b_1} = \frac{\partial L}{\partial z_2} \cdot \frac{\partial z_2}{\partial a_1} \cdot \frac{\partial a_1}{\partial z_1} \cdot \frac{\partial z_1}{\partial b_1}$$

$$\Rightarrow (a_2 - y)(w_2)(a_1)(1 - a_1)(1)$$

- the

(c) learning rate works on any parameter θ as $\theta := \theta - \alpha \frac{\partial L}{\partial \theta}$

we will update each parameter the same way.

→ learning rate controls the growth of the steps convergence time of the algorithm, speed and efficiency of the algo. It shouldn't be too large or either too small.

$$1. \hat{y}_i = \beta_0 + \beta_1 x_{i1} + \beta_2 x_{i2} \dots \beta_d x_{id}$$

$$X = \begin{bmatrix} 1 & x_{11} & x_{12} & \dots & x_{1d} \\ \vdots & x_{21} & \vdots & \ddots & x_{2d} \\ \vdots & \vdots & \vdots & \ddots & \vdots \\ 1 & x_{n1} & \dots & \dots & x_{nd} \end{bmatrix} \quad \hat{y} = X\beta$$

$$J(\beta) = \frac{1}{2} \|y - X\beta\|^2 = \frac{1}{2} \sum_{i=1}^n (y_i - \hat{y}_i)^2$$

(a) Regression is just trying to find the best values of the parameters and fitting a better and better straight line to the data. so that diff b/w ~~expectation~~ and the prediction and the actual output is as close as possible. squared error is a good method, because it penalises large diff. to prediction even more, and also makes $-ve$ and ve errors both positive because both are equally bad. and it also has a ~~single~~ global ~~min~~ min. So the local min is the global

$$(b) \hat{y} = X\beta$$

$$\text{error} = y - X\beta$$

$$J(\beta) = \frac{1}{2} \|y - X\beta\|^2 = \frac{1}{2} (y - X\beta)^T (y - X\beta) \quad (\text{if } A \text{ sym})$$

$$J(\beta) = \frac{1}{2} (y^T y - 2\beta^T X^T y + \beta^T X^T X \beta)$$

take grad. w.r.t β

$$\nabla_{\beta} J(\beta) = \frac{1}{2} (\nabla(y^T y) - \nabla(2\beta^T X^T y) + \nabla(\beta^T X^T X \beta))$$

$$= \frac{1}{2} \begin{bmatrix} 0 & -2X^T y & 2X^T X \beta \end{bmatrix}$$

$$= \boxed{X^T X \beta - X^T y}$$

$$\begin{cases} \nabla_{\beta} (\beta^T A \beta) = 2A\beta \\ \nabla_{\beta} (\beta^T c) = c \end{cases}$$

setting it to 0

$$X^T X \beta = X^T y$$

$$\therefore \beta = (X^T X)^{-1} X^T y$$

(c) matrix inversion is ~~roughly~~ the order of $O(n^3)$

and if the data size is large (> 1000) then doing the matrix computation and inversion will be very time-heavy. so it is not used for large datasets. iterative methods are good because they scale to large data also and are bound to converge if the correct learning rate is chosen.