## Q1

a) Regression is used to predict a value using data already provided to us by the user.

With the given data, our model makes adjustments in values of parameter ($\beta$ randomly set at first) such that the squared error is least and comes up with an equation of curve which helps in our prediction.

- Minimizing sq. error helps us computationally.
- It makes the function differentiable, making it easier to find grad. of cost fun$^c$.
- It also makes all values of error positive. Direct summation may lead to zero error, telling our model has 100% accuracy, which isn't true.

b)
$$J = \frac{1}{2} \|y - X\beta\|^2$$

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

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

$$y = \begin{bmatrix} y_1 \\ y_2 \\ \vdots \\ \vdots \\ y_n \end{bmatrix}_{n \times 1}$$

$$x = \begin{bmatrix} \\ \\ \end{bmatrix}_{n \times (d+1)}$$

$$\beta = \begin{bmatrix} \beta_0 \\ \vdots \\ \vdots \\ \beta_d \end{bmatrix}_{(d+1) \times 1}$$

$y^T X\beta \rightarrow (1 \times n)\,(n \times (d+1))\,((d+1) \times 1)$
$\Rightarrow (1 \times 1)$
$\hookrightarrow$ Scalar

$$\nabla\left(\frac{1}{2}(y^T y - y^T X\beta - (y^T X\beta)^T + \beta^T X^T X\beta)\right)$$
$$\underbrace{\hookrightarrow \text{Scalar} \therefore \text{equal.}}$$

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

$$= \frac{1}{2}(-X^T y + X^T X\beta) = 0$$

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

$$\beta y = x$$

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

c) Multiplication and taking inverse of a high dimensional matrix is very time taking and also computationally expensive.

In gradient descent, ~~we don't requ018 require~~ we don't have to calculate inverse of a matrix, this makes the process very fast, also at the same time taking less memory.

## Q2

a) In back propogation, we compute the gradient of our loss func and adjust our parameters and biases accordingly so that our loss func minimises.

Chain Rule makes our computation cost low ~~as~~ as it is much easier and faster.

b) $L = -y \log a_2 + (1-y) \log (1-a_2)$

$$\frac{\partial L}{\partial z_2} = a_2 - y \qquad\qquad z_2 = \omega_2 a_1 + b_2$$

$$\frac{\partial L}{\partial \omega_2} = \frac{\partial L}{\partial z_2} \frac{\partial z_2}{\partial \omega_2}$$

$$= (a_2 - y)(a_1)$$

$$\frac{\partial L}{\partial b_2} = \frac{\partial L}{\partial z_2} \frac{\partial z_2}{\partial b_2}$$

$$\frac{\partial L}{\partial b_2} = \frac{\partial L}{\partial z_2} \frac{\partial z_2}{\partial b_2}$$

$$= (a_2 - y)(1)$$

$$\frac{\partial L}{\partial w_1} = \left( \frac{\partial L}{\partial z_1} \frac{\partial z_2}{\partial a_1} \frac{\partial z_1}{\partial a_1} \frac{\partial a_1}{\partial w_1} \right)$$

$$= (a_2 - y)(w_2)(a_1)(1 - a_1)$$

$$\frac{\partial L}{\partial b_1} = \left( \frac{\partial L}{\partial z_1} \right) \frac{\partial z_1}{\partial a_1}$$

$$\frac{\partial L}{\partial w_1} = \frac{\partial L}{\partial z_1} \frac{\partial z_1}{\partial w_1}$$

$$= \frac{\partial L}{\partial z_2} \frac{\partial z_2}{\partial a_1} \frac{\partial a_1}{\partial z_1} \frac{\partial z_1}{\partial w_1}$$

$$= (a_2 - y) w_2 \, a_1 (1 - a_1) \, x$$

$$\frac{\partial L}{\partial b_1} = \frac{\partial L}{\partial z_1} \frac{\partial z_1}{\partial b_1}$$

$$= (a_2 - y) w_2 a_1 (1 - a_1) \, (1)$$

c) $\quad w_1' := w_1 - \alpha \dfrac{\partial L}{\partial w_1}$

$w_2' := w_2 - \alpha \dfrac{\partial L}{\partial w_2}$

$b_1' := b_1 - \alpha \dfrac{\partial L}{\partial b_1}$

$b_2' := b_2 - \alpha \dfrac{\partial L}{\partial b_2}$

$\alpha$ (learning rate)
Controls step size ie,
if the step is too large, it
may completely miss the
local minima or if the step
is too small it converges
very slowly.

## Q3

a) In AANN, the input is processed independently whereas in RNN is designed for sequential data.

RNN have a hidden layer, whereas ANN doesn't

a) RNNs process input sequentially, by using an hidden state that serves as a memory of past inputs.

ANN processes input by passing data through layers of interconnected neurons.

b) RNN struggle with long term dependencies due to vanishing grad. problem.

c) Role of gates in LSTMs is to control flow of information. which information to forget/input/output.
The gates allow LSTMs to retain selective information, helping in long term memory

d) With the help of gates. in LSTMs, gradients vanishes very slowly addressing to the problem of vanishing grad.

e) ANN → Price of house prediction
RNN → Word predictor / Sentence Completer.
LSTM → language translation.

## Q4

Eg: The ~~boy~~ man asked his girlfriend to marry him.

To correctly predict the last word (him) the model must remember that we are talking about ~~them~~ the man and not his girlfriend.

Yes, a standard RNN would struggle because of the vanishing grad, as a result the model will forget about the man ~~so~~ by the time it reaches to ~~RNN~~ 'marry'.

b) The LSTM gates allow LSTMs to retain relevant info for long time.
The gates, forget gate, input gate, output gate, help in this ~~tasks~~ tasks.
It forgets the irrelevant info and only retains the important info.

Eg ~~I went to store as~~ I went for a walk as I needed some fresh air.

The model must remember that the ~~its~~ sentence is talking about me to predict the ~~its~~ second I in the sentence.