Q.A) Back propogation involves computing of gradients in a multiple layer neural network. This is typically done after the forward propogation step.

F Propogation

$$z^1 = w^1 x + b^1$$

$$h(z) = (\sigma(z_1), \sigma(z_2) \cdots \sigma(z_n)) \; R^2 \to R^2$$

$$a^2 = h(z_1)$$

$$z^2 = w^2 a^2 + b^2$$

$$\vdots$$

$$z^m = w^m a^m + b^m$$

$$a^{m+1} = h(z^m)$$

$$\tilde{a}_j^{m+1} = r_j^{m+1} a^{m+1}$$

$$\xrightarrow{\quad} \text{Bernoullirv. } w \to r_j^{m+1}$$

$$\therefore z^{m_0+1} = w^{m_0+1} \underline{a^{m_0+1}} + b^{m_0+1} \xrightarrow{\quad} r_j^{m+1} a^{m+1}$$

$$\therefore z_0(x,y) = (g(z^{m_0+1})))y$$

$$z^{m_0+1} = w^{m_0+1} a^{m_0+1} + b^{m_0+1}$$

$$\Rightarrow z_0(x,y) = (g(z^{m_0+1}))y$$

$$\Theta \Rightarrow \{w^1, w^{m_0+1}, b^1 \cdots b^{m_0+1}\} \text{ where } w^1 \in R^{\{x a}, w^{m_0+1} \in R^{kx}$$

$$b^m \in R^2, b^{m_0+1} \in R^k \quad 1 \leq m \leq m_0 \quad 1 < m < m_2 \qquad w^m \in R^{LxL}$$

Back propogation

$$z^m = w^m a^m + b^m \qquad \text{Multiply by Bernoull rv}$$

$$a^{m+1} = h(z^m) \qquad \tilde{a}_j^{m+1} = r_j^{m+1} a_j^{m+1}$$

$$\text{where } r_j^{m+1} \sim \text{Bernoulli}(P) \text{ and } r_j^m \text{ and } r_j^{m'} \text{ are independent.}$$

Prediction $\Rightarrow z^m = p w^m a^m + b^m \qquad a^{m+1} = h(z^m)$ [Output of dropout]
[while training]

$$E[\tilde{a}_j^m] = a_j^m p + (1-p)0 = a_j^m p$$

A good optimizer $\Rightarrow$ SGD w momentu.

$\Rightarrow$ SGD has learning rate $\eta$ and initial paramter $\theta$

while not finished :—

    sample a minibatch of size m from ha

set $\{x^{(i)} \ldots x^{(n)}\}$ and $\hat{g} = 0$

for $i = 1 \ldots m$ do

    Computed Gradient Descent.

$$\hat{g} \leftarrow \hat{g} + \nabla_\theta L(f(x^{(i)}, \theta))$$

    Apply update $\theta = \theta_x - \eta \hat{g}$

SGD w momentu has the difference when

Computing gradient descent $\Rightarrow \hat{g} \leftarrow g + \nabla_\theta L(f(x^{(i)}, \theta)$

$$r = \alpha r - zg$$
$$\Rightarrow \theta \leftarrow \theta + r$$

Thus we have derived the back propagation
Algorithm for the multi layer feedforward
neural network. We have also stated
SGD Algorithm for this neural network.

Hence proved.  QED.