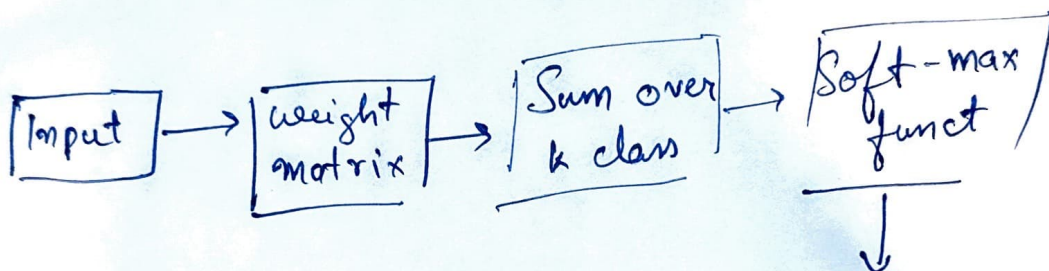


Ques 2

$$\hat{p}_k = \delta(s_k(x))_k = \frac{\exp(s_k(x))}{\sum_{j=1}^k \exp(s_j(x))}$$

$$(a) \quad s_k(x) = x^T (\Theta)^k$$

we need $\Theta^{(k)}$ parameters. All these are stored in parameter matrix.



The input vect is fed into the weight matrix and the bias vector is added to the weighted sum. The result is then passed through the Softmax function, which produces the output probability distribution over the k class.

$$(b) \quad J(\Theta) = -\frac{1}{m} \sum_{i=1}^m \sum_{k=1}^K y_k^{(i)} \log(\hat{p}_k^{(i)})$$

we have to derive gradient of $J(\Theta)$

$$\partial(J(\Theta)) = -\frac{1}{m} \sum_{i=1}^m x_i \cdot \frac{1(e^{\Theta^T x_i})}{\sum_{j=1}^K e^{\Theta^T x_j}}$$

$$= \frac{1}{m} \sum_{i=1}^m (\hat{p}(x)^{(i)} - y_k^{(i)}) x^{(i)}$$

$$-2x^T y + 2(x^T x + \lambda I)w = 0$$

$$-2x^T y = -2(x^T x + \lambda I)w$$

$$x^T y = (x^T x + \lambda I)w$$

$$(x^T x + \lambda I)w = x^T y$$

multiply both sides by

$$(x^T x + \lambda I)^{-1}$$

$$(x^T x + \lambda I)^{-1} (x^T x + \lambda I)w$$

$$\Rightarrow (x^T x + \lambda I)^{-1} x^T y$$

$$-2x^T y + 2(x^T x + \lambda I)w = 0$$

$$x^T y = (x^T x + \lambda I)w$$

now multiply by $(x^T x + \lambda I)^{-1}$

$$(x^T x + \lambda I)^{-1} (x^T x + \lambda I)w$$

$$= (x^T x + \lambda I)^{-1} x^T y$$

$$w = (x^T x + \lambda I)^{-1} x^T y$$

which is ridge regression.