(Q2)

$$P(y=1|x,w) = \frac{1}{2}\left(1 + \frac{w^T x}{\sqrt{1+(w^T x)^2}}\right)$$

$$P(y=0|x,w) = 1 - P(y=1|x,w)$$

$J(w) = $ Cost function

$$= -\left\{\sum_{i=1}^{n} y^i \log P(y=1|x_i,w) + (1-y^i)\log P(y=0|x_i,w)\right\}$$

* Neglecting the -ve sign for now and forming a generalized equation without $\sum$, will add $\sum$ later

$$J(w)^+ = y^i \log P(y=1|x_i,w) + (1-y^i)\log P(y=0|x_i,w)$$

$$= y^i \log \frac{1}{2}\left(1 + \frac{w^T x_i}{\sqrt{1+(w^T x_i)^2}}\right) + (1-y^i)\log\left(1 - \frac{1}{2}\left(1 + \frac{w^T x_i}{\sqrt{1+(w^T x_i)^2}}\right)\right)$$

$$= y^i \log \frac{1}{2}\left(1 + \frac{w^T x_i}{\sqrt{1+(w^T x_i)^2}}\right) + (1-y^i)\log\frac{1}{2}\left(1 - \frac{w^T x_i}{\sqrt{1+(w^T x_i)^2}}\right)$$

$\Rightarrow \log a \times b = \log a + \log b \qquad a = \frac{1}{2} ; b = \left(1 + \frac{w^T x_i}{\sqrt{1+(w^T x_i)^2}}\right)$ or $\left(1 - \frac{w^T x_i}{\sqrt{1+(w^T x_i)^2}}\right)$

Since $y^i \log a$ and $(1-y^i)\log a$ are gonna be constant when we derive/differentiate $J(w)^+$ with respect to $w_j$ there is no need to solve it and we neglect the terms.

①

$$\frac{\partial J(\omega)^k}{\partial \omega_{ij}} = \frac{\partial \, y^i \log\left(1 + \frac{\omega^T x_i}{\sqrt{1+(\omega^T x_i)^2}}\right)}{\partial w_{ij}} + \frac{\partial \, (1-y^i) \log\left(1 - \frac{\omega^T x_i}{\sqrt{1+(\omega^T x_i)^2}}\right)}{\partial w_{ij}}$$

$$\Rightarrow \frac{\partial \log x}{\partial x'} = \frac{1}{x} \times \frac{\partial x}{\partial x'} \quad ; \quad \frac{\partial a \times b}{\partial x} = b \frac{\partial a}{x} + \frac{a \partial b}{x} \quad \text{where } a \text{ and } b \text{ are functions of } x$$

$$\underbrace{\qquad\qquad\qquad\qquad\qquad}_{\text{Using the rules}}$$

$$= \frac{y^i \sqrt{1+(\omega^T x_i)^2}}{\sqrt{1+(\omega^T x_i)^2} + \omega^T x_i} \left( \frac{x_{ij}}{\sqrt{1+(\omega^T x_i)^2}} - \frac{\cancel{2}}{\cancel{2}} \frac{(\omega^T x)^2 \times x_{ij}}{\left(1+(\omega^T x_i)^2\right)^{3/2}} \right)$$

$$+ \frac{(1-y_i)\sqrt{1+(\omega^T x_i)^2}}{\sqrt{1+(\omega^T x_i)^2} - \omega^T x_i} \left( -\frac{x_{ij}}{\sqrt{1+(\omega^T x_i)^2}} + \frac{\cancel{2}}{\cancel{2}} \frac{(\omega^T x)^2 \times x_{ij}}{\left(1+(\omega^T x_i)^2\right)^{3/2}} \right)$$

$$= \frac{y^i \sqrt{1+(\omega^T x_i)^2}}{\left\{\sqrt{1+(\omega^T x_i)^2} + (\omega^T x_i)\right\} \times \sqrt{1+(\omega^T x_i)^2}} \left( x_{ij} - \frac{(\omega^T x_i)^2 x_{ij}}{1+(\omega^T x_i)^2} \right)$$

$$+ \frac{(1-y_i)\sqrt{1+(\omega^T x_i)^2}}{\left\{\sqrt{1+(\omega^T x_i)^2} + (\omega^T x_i)\right\} \times \sqrt{1+(\omega^T x_i)^2}} \left( -x_{ij} + \frac{(\omega^T x_i)^2 x_{ij}}{1+(\omega^T x_i)^2} \right)$$

②

$$= \frac{y^i x_{ij}}{\sqrt{1+(\omega^T x_i)^2}+(\omega^T x_i)} \left( \frac{1 + (\omega^T x_i)^2 - (\omega^T x_i)^2}{1+(\omega^T x_i)^2} \right)$$

$$+ \frac{(1-y^i) x_{ij}}{\sqrt{1+(\omega^T x)^2}-(\omega^T x_i)} \left( \frac{-1 - (\omega^T x_i)^2 + (\omega^T x_i)^2}{1+(\omega^T x_i)^2} \right)$$

$$= \frac{x_{ij}}{((\omega^T x_i)^2 + 1)} \left\{ \frac{y_i}{(\sqrt{1+(\omega^T x_i)^2}+(\omega^T x_i)) \times 1} + \frac{(y_i - 1)}{(\sqrt{1+(\omega^T x_i)^2}-(\omega^T x_i))} \right\}$$

$$= \frac{x_{ij}}{(\omega^T x_i)^2 + 1} \left( \frac{y_i(\sqrt{1+(\omega^T x_i)^2} - \omega^T x_i) + (y_i - 1)(\sqrt{1+(\omega^T x_i)^2} + \omega^T x_i)}{(\sqrt{1+(\omega^T x_i)^2})^2 - (\omega^T x_i)^2} \right)$$

$$\Rightarrow (a+b)(a-b) = a^2 - b^2$$

$$= \frac{x_{ij}}{(\omega^T x_i)^2 + 1} \left( 2y_i \sqrt{1+(\omega^T x_i)^2} - \sqrt{1+(\omega^T x_i)^2} - \omega^T x_i \right)$$

$$= x_{ij} \left( \frac{2y_i}{\sqrt{1+(\omega^T x_i)^2}} - \frac{1}{\sqrt{1+(\omega^T x_i)^2}} - \frac{\omega^T x_i}{1+(\omega^T x_i)^2} \right)$$

So $\dfrac{\delta J(\omega)}{\delta \omega_{ij}} = x_{ij}\left(\dfrac{2y_i}{\sqrt{P}} - \dfrac{1}{\sqrt{P}} - \dfrac{\omega^T x}{P}\right)$

where $P = \left(1 + \omega^T x_i\right)^2$

This is the gradient descent with respect to weight $j$. Since we have $n$ samples $i \in n$. We can use batch or stochastic gradient descent to calculate the optimal weights.

$$\omega_{k+1}^j = \omega_k^j - \alpha \dfrac{\delta J(\omega)}{\delta \omega_j} \qquad \text{(iterative approach)}$$

$\dfrac{\delta J(\omega)}{\delta \omega_j} = \left(\displaystyle\sum_{i=1}^{n} \dfrac{\delta J(\omega)}{\delta \omega_{ij}}\right) / n$ ; Since earlier I had omitted the $-ve$ sign from cost function which directly impacts the gradient descent function since it is a derivation of $J(\omega)$ cost function. I will add the $-ve$ sign direct to step rule of iterative approach of my model. So updated rule looks like

$$\omega_{k+1}^j = \omega_k^j + \alpha \dfrac{\delta J(\omega)}{\delta \omega_j}$$

where $\omega^j$ is the $j^{th}$ feature $\omega_k$ is the $k$ the iteration & $\alpha$ is the step size.

④