

1. (a)

$$\begin{aligned}
 \nabla_{\theta} J(\theta) &= \\
 -\frac{1}{m} \sum_{i=1}^m &\left(\frac{y^{(i)} g(\theta^T x^{(i)}) (1-g(\theta^T x^{(i)})) x^{(i)}}{h_{\theta}(x^{(i)})} \right. \\
 &\left. - \frac{(1-y^{(i)}) g(\theta^T x^{(i)}) (1-g(\theta^T x^{(i)})) X^{(i)}}{1-h_{\theta}(x^{(i)})} \right) \\
 = -\frac{1}{m} \sum_{i=1}^m &(y^{(i)} (1-h_{\theta}(x^{(i)})) X - (1-y^{(i)}) h_{\theta}(x^{(i)}) X) \\
 = -\frac{1}{m} \sum_{i=1}^m &(y^{(i)}(X^{(i)} - h_{\theta}(x^{(i)}) X^{(i)}) = -\frac{1}{m} X^T (y - h_{\theta}(X)) \\
 \text{where } X = &\begin{pmatrix} x^{(1)^T} \\ x^{(2)^T} \\ \vdots \\ x^{(m)^T} \end{pmatrix}, h_{\theta}(X) g(X\theta) \\
 \nabla_{\theta} J(\theta) &= \left[\frac{\partial P_{\theta} J(\theta)}{\partial \theta_1}, \dots, \frac{\partial P_{\theta} J(\theta)}{\partial \theta_n} \right] \\
 \frac{\partial P_{\theta} J(\theta)}{\partial \theta_j} &= \frac{1}{m} \sum_{i=1}^m g(\theta^T x^{(i)}) x^{(i)} \\
 = \frac{1}{m} \sum_{i=1}^m &g(\theta^T x^{(i)}) (1-g(\theta^T x^{(i)})) x_j^{(i)} X^{(i)} \\
 \text{so } H = \nabla_{\theta}^2 J(\theta) &= \frac{1}{m} \sum_{i=1}^m h_{\theta}(x^{(i)}) (1-h_{\theta}(x^{(i)})) X^{(i)} X^{(i)T} \\
 &= \frac{1}{m} [X^T h_{\theta}(X)] \left[(1-h_{\theta}(X)) X \right]
 \end{aligned}$$

$$\begin{aligned}
 Z^T H Z &= \frac{1}{m} \sum_{i=1}^m h_{\theta}(x^{(i)}) (1-h_{\theta}(x^{(i)})) Z^T X^{(i)} X^{(i)T} Z \\
 &= \frac{1}{m} \sum_{i=1}^m h_{\theta}(x^{(i)}) (1-h_{\theta}(x^{(i)})) (X^{(i)T} Z)^2 \geq 0
 \end{aligned}$$

$$\begin{aligned}
 (c) P(y=1|X) &= \frac{P(X|Y=1) P(Y=1)}{P(X)} \\
 &= \frac{P(X|Y=1) P(Y=1)}{P(X|Y=1) P(Y=1) + P(X|Y=0) P(Y=0)} \\
 &= \frac{\phi \exp(-\frac{1}{2}(X-\mu_0)^T \Sigma^{-1}(X-\mu_0))}{\phi \exp(-\frac{1}{2}(X-\mu_1)^T \Sigma^{-1}(X-\mu_1)) + (1-\phi)} \\
 &= \frac{\exp(-\frac{1}{2}(X-\mu_0)^T \Sigma^{-1}(X-\mu_0))}{1 + \exp(-\frac{1}{2}(X-\mu_0)^T \Sigma^{-1}(X-\mu_0) + \frac{1}{2}(X-\mu_1)^T \Sigma^{-1}(X-\mu_1))} \\
 &= \frac{1}{1 + \exp(-(\frac{1}{2}\mu_0^T \Sigma^{\frac{1}{2}} \mu_0 + \frac{1}{2}\mu_1^T \Sigma^{\frac{1}{2}} \mu_1 - \frac{1}{2} X^T \Sigma^{\frac{1}{2}} (\mu_0 + \mu_1) \\
 &\quad - \frac{1}{2} (\mu_0 - \mu_1)^T \Sigma^{-1} X + \ln \frac{\phi}{1-\phi}))} \\
 &= \frac{1}{1 + \exp(- (W_0 \cdot \mu_0)^T \Sigma^{\frac{1}{2}} X + \frac{1}{2} (W_0^T \Sigma^{\frac{1}{2}} \mu_0 - W_1^T \Sigma^{\frac{1}{2}} \mu_1) + b_0) }
 \end{aligned}$$

$$\theta = (-(\mu_0 - \mu_1)^T \Sigma^{-1})^T = \Sigma^{-1}(\mu_1 - \mu_0)$$

$$\theta_0 = \frac{1}{2} (\mu_0^T \Sigma^{-1} \mu_0 - \mu_1^T \Sigma^{-1} \mu_1) + \ln \frac{\phi}{1-\phi}$$

(d)

$$f = \sum_{i=1}^m \log \frac{1}{(2\pi)^{\frac{m}{2}} |\Sigma|^{\frac{1}{2}}} \exp(-\frac{1}{2}(x - M_{y^{(i)}})^T \Sigma^{-1} (x - M_{y^{(i)}}) + \sum_{i=1}^m \log \phi^{y^{(i)}} (1-\phi)^{1-y^{(i)}})$$

$$M_{y^{(i)}} = y^{(i)} \mu_1 + (1-y^{(i)}) \mu_0$$

$$\ell = -\frac{mn}{2} \log 2\pi - \frac{m}{2} \log |\Sigma| + \sum_{i=1}^m \frac{1}{2} (x - (y^{(i)} \mu_1 + (1-y^{(i)}) \mu_0))^T \Sigma^{-1} (x - (y^{(i)} \mu_1 + (1-y^{(i)}) \mu_0)) + y^{(i)} \log \phi + (1-y^{(i)}) \log (1-\phi)$$

$$\frac{\partial \ell}{\partial \phi} = \sum_{i=1}^m \left(\frac{y^{(i)}}{\phi} - \frac{1-y^{(i)}}{1-\phi} \right) = 0$$

$$\sum_{i=1}^m (y^{(i)} - y^{(i)} \phi) = \sum_{i=1}^m (\phi - y^{(i)} \phi)$$

$$\phi = \frac{1}{m} \sum_{i=1}^m y^{(i)} = \frac{1}{m} \sum_{i=1}^m I\{y^{(i)} = 1\}$$

$$\frac{\partial \ell}{\partial \mu_0} = \sum_{i=1}^m \frac{1}{2} (1-y^{(i)}) \sum_{i=1}^m (x - (y^{(i)} \mu_1 + (1-y^{(i)}) \mu_0))^T \Sigma^{-1} (x - (y^{(i)} \mu_1 + (1-y^{(i)}) \mu_0)) \Sigma^{-1} (1-y^{(i)}) = 0$$

$$\sum_{i=1}^m (1-y^{(i)}) (x - (y^{(i)} \mu_1 + (1-y^{(i)}) \mu_0))^T \Sigma^{-1} (x - (y^{(i)} \mu_1 + (1-y^{(i)}) \mu_0)) = 0$$

the same reasoning leads to

$$\sum_{i=1}^m y^{(i)} (x - (y^{(i)} \mu_1 + (1-y^{(i)}) \mu_0)) = 0$$

$$\text{so } \sum_{i=1}^m (1-y^{(i)}) (x^{(i)} - \mu_0) = 0 \quad (\text{when } y^{(i)} = 1, \text{ it goes to 0})$$

$$\sum_{i=1}^m y^{(i)} (x^{(i)} - \mu_1) = 0$$

$$\mu_0 = \frac{\sum_{i=1}^m (1-y^{(i)}) x^{(i)}}{\sum_{i=1}^m (1-y^{(i)})}$$

$$\mu_1 = \frac{\sum_{i=1}^m y^{(i)} x^{(i)}}{\sum_{i=1}^m y^{(i)}}$$

$$\nabla_{\Sigma} \ell = \frac{m}{2} \sum_{i=1}^m \sum_{j=1}^m \nabla (U_i^T \Sigma^{-1} U_j)$$

= 0

$$\sum_{i=1}^m \nabla (U_i^T \Sigma^{-1} U_i) = -m \Sigma^{-1}$$

$$(\Sigma + h)^{-1} = (\Sigma (I + \Sigma^{-1} h))^{-1}$$

$$= (I + \Sigma^{-1} h)^{-1} \Sigma^{-1}$$

$$= (I - \Sigma^{-1} h) \Sigma^{-1}$$

$$= \Sigma^{-1} - \Sigma^{-1} h \Sigma^{-1}$$

$$= U^T (\Sigma + h)^{-1} U$$

$$= U^T \Sigma^{-1} U - U^T \Sigma^{-1} h \Sigma^{-1} U$$

$$= U^T \Sigma^{-1} U + \nabla(U^T \Sigma^{-1} U) \cdot h^T > + d(h)$$

$$\nabla(U^T \Sigma^{-1} U) h = -U^T \Sigma^{-1} h \Sigma^{-1} U$$

$$= -\nabla(U^T \Sigma^{-1} h \Sigma^{-1} U) = -\nabla(\Sigma^T U^T \Sigma^{-1} h)$$

$$\text{so } \nabla U^T \Sigma^{-1} U = -\Sigma^{-1} U U^T \Sigma^{-1}$$

$$\sum_{i=1}^m \sum_{j=1}^m u_i u_j U_i^T \Sigma^{-1} = m \Sigma^{-1}$$

$$\Sigma = \frac{1}{m} \sum_{i=1}^m u_i u_i^T$$

where $u_i = x^{(i)} - M_{y^{(i)}}$

g) dataset 1. because the distribution of the $p(x|y)$ is not multi-gaussians distribution

h) the box-cox transformation

We should add 1 to every data in x to avoid the negative or 0 that makes the boxcox trasformation disabled

$$y^{(\lambda)} = \begin{cases} \frac{y^\lambda - 1}{\lambda} & \text{if } \lambda \neq 0 \\ \ln(y) & \text{if } \lambda = 0 \end{cases}$$

The function statsboxcox() can help us to find the optimised λ .

$$2. \quad (a) \quad P(t|y) =$$

$$\frac{P(y+t)}{P(t+\lambda)} = \frac{P(y+t)}{P(t)}$$

$$P(t)P(y+t) = P(y+t)P(t+x)$$

$$\frac{P(t+x)}{P(y+x)} = \frac{P(t)P(y+x)}{P(y+t)P(y+x)}$$

$$P(y+t) = P(t|\theta) P(y) = P(y)$$

$$P(y+t+x) = P(y+t|x) P(x)$$

$$= P(y|x) P(x) = P(xy)$$

$$\text{so } \frac{P(t+x)}{P(y+x)} = \frac{P(t)P(x|y)}{P(y)P(y|x)} = \frac{P(t)}{P(y)} = \frac{1}{\alpha}$$

$$(b) h(x^{(i)}) \sim p(y^{(i)}|x^{(i)}) \Rightarrow p(t^{(i)}|x^{(i)})$$

$$=\alpha$$

(Q) for the boundary

$$\theta^T x = 0$$

now we need correction

$$\exp(\theta^T x) + 1 = \frac{\alpha}{\alpha}$$

$$\theta^T x = -\ln\left(\frac{\alpha}{\alpha} - 1\right) \text{ contribute to } \theta_0$$

$$\theta'_0 = \theta_0 + \ln\left(\frac{\alpha}{\alpha} - 1\right) = \alpha \cdot \theta_0$$

$$\alpha = \frac{-\ln\left(\frac{\alpha}{\alpha} - 1\right)}{\theta_0}$$

$$3. (a) P(y|\lambda) = \frac{1}{y!} \exp(-\lambda + y \ln \lambda)$$

$$b(y) = \frac{1}{y!}, \eta = (\ln \lambda, T(y)-y), a(\eta) = e^\eta$$

$$(b) E[T(y); \eta] = \lambda = e^\eta = e^{\theta^T x}$$

$$(c) \frac{\partial}{\partial \theta_j} \log P(y^{(i)}|x^{(i)}; \theta)$$

$$= \frac{\partial}{\partial \theta_j} \left(\log \frac{1}{y!} + (-e^{\theta^T x^{(i)}} + y^{(i)} e^{\theta^T x^{(i)}}) \right)$$

$$= -e^{\theta^T x^{(i)}} \cdot x_j^{(i)} + y^{(i)} x_j^{(i)}$$

$$= (y^{(i)} - e^{\theta^T x^{(i)}}) x_j^{(i)}$$

$$\theta_j := \theta_j + \alpha (y^{(i)} - e^{\theta^T x^{(i)}}) x_j^{(i)}$$

$$\theta := \theta + x^T (Y - \exp(X\theta))$$

4. (a)

$$\frac{\partial}{\partial \eta} \int p(y; \eta) dy = 0$$

$$\frac{\partial}{\partial \eta} \int p(y; \eta) dy$$

$$= \int \frac{\partial}{\partial \eta} p(y; \eta) dy$$

$$= \int b(y) \exp(\eta y - a(\eta)) (y - a'(\eta)) dy$$

$$= \int p(y; \eta) y dy - a'(\eta) \int p(y; \eta) dy$$

$$= E[Y|X; \theta] - a'(\eta) = 0$$

$$\text{so } E[Y|X; \theta] = \frac{da(\eta)}{d\eta}$$

$$(b) \frac{\partial^2}{\partial \eta^2} \int p(y; \eta) dy = 0$$

$$\int \frac{\partial}{\partial \eta} p(y; \eta) y dy - a''(\eta) = 0$$

$$\int p(y; \eta) (y - a(\eta)) y dy - a''(\eta) = 0$$

$$\int p(y; \eta) (y - a(\eta))^2 + a'(y) \underbrace{\int p(y; \eta) (y - a(\eta)) dy}_{=0} - a''(\eta) = 0$$

$$\text{Var}(Y|X; \theta) = a''(\eta)$$

(c)

$$-\log p(y; \eta) = -\log b(y) - \eta y + a(\theta^T x)$$

$$NLL = -\log b(y) - \theta^T x y + a(\theta^T x)$$

$$(\theta + h)^T x = \theta^T x + \langle \nabla (\theta^T x), h \rangle + o(\|h\|)$$

$$x^T h = \nabla (\theta^T x)^T h$$

$$\nabla (\theta^T x) = x$$

$$\begin{aligned} \nabla_{\theta} NLL &= -xy + a'(\theta^T x) x \\ &= x(a'(\theta^T x) - y) \end{aligned}$$

Hessian NLL =

$$\left[\frac{\partial^2 NLL}{\partial \theta_1}, \frac{\partial^2 NLL}{\partial \theta_2}, \dots, \frac{\partial^2 NLL}{\partial \theta_m} \right]$$

$$= x \left[a''(\theta^T x) x_1, a''(\theta^T x) x_2, \dots \right]$$

$$= a''(\theta^T x) xx^T$$

h^T Hessian NLL h

$$= a''(\theta^T x) h^T x x^T h$$

$$= a''(\theta^T x) (h^T x)^2$$

$$= \text{Var}(Y|X; \theta) (h^T x)^2 \geq 0$$

so Hessian NLL is PSD

so it's convex

5. (a)

(i) suppose $M = X\theta - y = \begin{bmatrix} \theta^T x^{(1)} - y^{(1)} \\ \theta^T x^{(2)} - y^{(2)} \\ \vdots \\ \theta^T x^{(m)} - y^{(m)} \end{bmatrix}$

$$M^T M = \sum_{i=1}^n \sum_{j=1}^n (\theta^T x^{(i)}) w_{ij} (\theta^T x^{(j)})$$

$$= \frac{1}{2} \sum_{i=1}^n w^{(i)} (\theta^T x^{(i)} - y^{(i)})^2$$

$$\text{so } w_{ij} = \begin{cases} \frac{1}{2} w^{(i)} & (i=j) \\ 0 & (i \neq j) \end{cases}$$

(ii) $J(\theta) = \theta^T X^T W X \theta - \theta^T X^T W y - y^T W y$

$$\nabla_\theta J(\theta) = \nabla_\theta \text{tr} J(\theta) = \nabla_\theta (\text{tr} (\theta^T X^T W X \theta))$$

$$-2 \text{tr} (\theta^T X^T W y) = (\nabla_\theta \text{tr} (\theta^T X^T W X \theta))^T$$

$$-2 X^T W y = (\theta^T X^T W X + \theta^T X^T W X)^T - 2 X^T W y$$

$$= 2 (X^T W X \theta - X^T W y)$$

$$\text{so } \nabla_\theta J(\theta) = 0, \quad \theta = (X^T W X)^{-1} X^T W y$$

$$(iii)-(iv) = \log \prod_{i=1}^m p(y^{(i)} | x^{(i)}; \theta) = \sum_{i=1}^m \log p(y^{(i)} | x^{(i)}; \theta)$$

$$= \sum_{i=1}^m \log \frac{1}{\sqrt{2\pi \sigma^{(i)}}} - \frac{(y^{(i)} - \theta^T x^{(i)})^2}{2\sigma^{(i)2}}$$

$$= \sum_{i=1}^m \log \frac{1}{\sqrt{2\pi \sigma^{(i)}}} - \sum_{i=1}^m \frac{(y^{(i)} - \theta^T x^{(i)})^2}{2(\sigma^{(i)})^2}$$

to maximum $-L(\theta)$, that is to

$$\text{minimum } \sum_{i=1}^m \frac{(y^{(i)} - \theta^T x^{(i)})^2}{2(\sigma^{(i)})^2}$$

which is a weighted linear regression problem

$w^{(i)}$ is $\frac{1}{(\sigma^{(i)})^2}$ in this case