

# Homework II Yuting Sun.

P1

a. Proof:

Given,

$$G(t) = \frac{e^t}{1+e^t} \quad 1-G(t) = \frac{1}{1+e^t}$$

$$P(y_i=1|\vec{\theta}) = G(\vec{\theta}^T \vec{x}_i)$$

$$\therefore P(y_i=1|\vec{\theta}) + P(y_i=0|\vec{\theta}) = 1$$

$$\therefore P(y_i=0|\vec{\theta}) = 1 - P(y_i=1|\vec{\theta}) = 1 - G(\vec{\theta}^T \vec{x}_i)$$

$$P(y|\vec{\theta}) = P(y_1, \dots, y_n|\vec{\theta}) = \prod_{i=1}^n \left( P(y_i=1|\vec{\theta})^{y_i} \cdot P(y_i=0|\vec{\theta})^{1-y_i} \right)$$

$$= \prod_{i=1}^n G(\vec{\theta}^T \vec{x}_i)^{y_i} (1 - G(\vec{\theta}^T \vec{x}_i))^{1-y_i}$$

$$= \prod_{i=1}^n G(\vec{\theta}^T \vec{x}_i)^{y_i} \frac{1 - G(\vec{\theta}^T \vec{x}_i)}{(1 - G(\vec{\theta}^T \vec{x}_i))^{y_i}}$$

$$= \prod_{i=1}^n \left( \frac{G(\vec{\theta}^T \vec{x}_i)}{1 - G(\vec{\theta}^T \vec{x}_i)} \right)^{y_i} (1 - G(\vec{\theta}^T \vec{x}_i))$$

$$\therefore \frac{G(\vec{\theta}^T \vec{x}_i)}{1 - G(\vec{\theta}^T \vec{x}_i)} = \frac{e^t / (1+e^t)}{1 / (1+e^t)} = e^t \quad \text{when } t = \vec{\theta}^T \vec{x}_i$$

$$= \exp(\vec{\theta}^T \vec{x}_i)$$

$$\therefore L(\vec{\theta}|y) = -\log P(y|\vec{\theta}) = -\sum_{i=1}^n \left( \log(e^{\vec{\theta}^T \vec{x}_i})^{y_i} - \log(1 + e^{\vec{\theta}^T \vec{x}_i}) \right)$$

$$= \sum_{i=1}^n \left[ -y_i \vec{\theta}^T \vec{x}_i + \log(1 + e^{\vec{\theta}^T \vec{x}_i}) \right]$$

b. Proof  $\nabla \log(\vec{\theta}|\vec{y}) = \sum_{i=1}^n -\vec{x}_i y_i + \underbrace{\frac{e^{\vec{\theta}^T \vec{x}_i}}{1 + e^{\vec{\theta}^T \vec{x}_i}}}_{\text{scalar}} \vec{x}_i = -\sum_{i=1}^n \vec{x}_i (y_i - G(\vec{\theta}^T \vec{x}_i))$

$$X = (\vec{x}_1^T, \vec{x}_2^T, \dots, \vec{x}_n^T)^T$$

$$G\vec{\theta} = (G(\vec{\theta}), G(\vec{\theta}), \dots, G(\vec{\theta}))^T$$

$$G(\vec{\theta}^T \vec{x}_i) = G(\vec{\theta})$$

$$\nabla(\vec{\theta}|y) = -X^T \vec{y} + X^T G\vec{\theta} = -X^T (\vec{y} - G\vec{\theta})$$

$$c) \quad \nabla_{\vec{\theta}} \log(\vec{\theta} | \vec{y}) = - \sum_{i=1}^n \vec{x}_i \underbrace{(y_i - \sigma(\vec{\theta}^T \vec{x}_i))}_{\text{const.}}$$

$$\nabla_{\vec{\theta}} \log L(\vec{\theta} | \vec{y}) = + \sum_{i=1}^n \nabla_{\vec{\theta}} (\vec{x}_i \cdot \sigma(\vec{\theta}^T \vec{x}_i))$$

$$b_i(\vec{\theta}) = \sigma(\vec{\theta}^T \vec{x}_i) \quad \nabla \cdot \vec{x}_i \cdot \sigma_i = \vec{x}_i \cdot \nabla \left( \frac{e^{\vec{\theta}^T \vec{x}_i}}{1 + e^{\vec{\theta}^T \vec{x}_i}} \right) \cdot \frac{d \sigma_i}{d \vec{\theta}}$$

$$= \vec{x}_i \cdot \frac{e^{\vec{\theta}^T \vec{x}_i} \vec{\theta} - e^{\vec{\theta}^T \vec{x}_i} \cdot e^{\vec{\theta}^T \vec{x}_i}}{(1 + e^{\vec{\theta}^T \vec{x}_i})^2} \cdot \vec{x}_i^T$$

$$= \vec{x}_i \cdot \frac{e^{\vec{\theta}^T \vec{x}_i}}{(1 + e^{\vec{\theta}^T \vec{x}_i})^2} \cdot \vec{x}_i^T$$

$$= \vec{x}_i \cdot \frac{e^{\vec{\theta}^T \vec{x}_i}}{1 + e^{\vec{\theta}^T \vec{x}_i}} \cdot \frac{1}{1 + e^{\vec{\theta}^T \vec{x}_i}} \vec{x}_i^T$$

$$= \vec{x}_i \cdot \underbrace{b_i(\vec{\theta})(1 - b_i(\vec{\theta}))}_{\text{scalar}} \vec{x}_i^T = \vec{x}_i \vec{x}_i^T b_i(\vec{\theta})(1 - b_i(\vec{\theta}))$$

$$\Rightarrow \nabla_{\vec{\theta}}^2 \log L(\vec{\theta} | \vec{y}) = \sum_{i=1}^n \vec{x}_i \vec{x}_i^T b_i(\vec{\theta})(1 - b_i(\vec{\theta}))$$

$$= X^T D(\vec{\theta}) X$$

$$D_i(\vec{\theta}) = b_i(\vec{\theta}^T \vec{x}_i)(1 - b_i(\vec{\theta}^T \vec{x}_i)) = b_i(\vec{\theta})(1 - b_i(\vec{\theta}))$$

diagonal matrix.

d) <sup>Proof.</sup>  $X$  is full rank.  $-\nabla_{\theta}^2 \log L(\omega|y) = + \sum_{i=1}^n \vec{x}_i \vec{x}_i^T b_i(\vec{\theta}^T \vec{x}_i) (1 - b(\vec{\theta}^T \vec{x}_i))$

Pick up a vector  $\vec{y}$  ( $\vec{y} \neq \vec{0}$ )

$$\begin{aligned} & \vec{y}^T \nabla_{\theta}^2 (\log L(\omega|y)) \vec{y} \\ &= \vec{y}^T \left( \sum_{i=1}^n \vec{x}_i \vec{x}_i^T \frac{e^{\vec{\theta}^T \vec{x}_i}}{(1 + e^{\vec{\theta}^T \vec{x}_i})^2} \right) \vec{y} \\ &= \sum_{i=1}^n (\vec{y}^T \vec{x}_i)^2 \frac{e^{\vec{\theta}^T \vec{x}_i}}{(1 + e^{\vec{\theta}^T \vec{x}_i})^2} \end{aligned}$$

$$\because e^{\vec{\theta}^T \vec{x}_i} > 0 \Rightarrow \frac{e^{\vec{\theta}^T \vec{x}_i}}{1 + e^{\vec{\theta}^T \vec{x}_i}} > 0$$

$\because X$  full rank.  $\therefore (\vec{y}^T \vec{x}_i)^2 = 0$  only when  $\vec{y} = \vec{0}$

However,  $\vec{y} \neq \vec{0}$ , thus  $(\vec{y}^T \vec{x}_i)^2 > 0$

Therefore, if  $X$  is full rank, the Hessian is positive definite. (PD)

If  $N(X)$  is non-trivial.

$$\begin{aligned} & \vec{y}^T \nabla_{\theta}^2 (\log L(\omega|y)) \vec{y} \\ &= \sum_{i=1}^n (\vec{y}^T \vec{x}_i)^2 \frac{e^{\vec{\theta}^T \vec{x}_i}}{(1 + e^{\vec{\theta}^T \vec{x}_i})^2} \end{aligned}$$

$\therefore N(X) \neq \emptyset$

By picking suitable  $\vec{y} \neq \vec{0}$ , it is possible that  $(\vec{y}^T \vec{x}_i)^2 = 0$

$$\Rightarrow \vec{y}^T \nabla_{\theta}^2 (\log L(\omega|y)) \vec{y} \geq 0$$

Therefore, if  $N(X)$  is non-trivial, the Hessian is positive semi-definite (PSD)

$$e) \theta^{(l+1)} = \theta^{(l)} - (\nabla_{\theta}^2 \log L(\theta|y)|_{\theta=\theta^{(l)}})^{-1} \nabla_{\theta} \log L(\theta|y)|_{\theta=\theta^{(l)}}$$

The Newton-Raphson step is:

$$\begin{aligned} \theta^{(l+1)} &= \theta^{(l)} - (X^T D(\theta^{(l)}) X)^{-1} (-X^T (y - b(\theta^{(l)}))) \\ &= \theta^{(l)} + (X^T D(\theta^{(l)}) X)^{-1} X^T (y - b(\theta^{(l)})) \\ &= (X^T D(\theta^{(l)}) X)^{-1} \cdot X^T D(\theta^{(l)}) X \theta^{(l)} + (X^T D(\theta^{(l)}) X)^{-1} X^T (y - b(\theta^{(l)})) \\ &= (X^T D(\theta^{(l)}) X)^{-1} (X^T D(\theta^{(l)}) X \theta^{(l)} + X^T D(\theta^{(l)}) \cdot D(\theta^{(l)})^{-1} (y - b(\theta^{(l)}))) \\ &= (X^T D(\theta^{(l)}) X)^{-1} X^T D(\theta^{(l)}) (X \theta^{(l)} + D(\theta^{(l)})^{-1} (y - b(\theta^{(l)}))) \\ &= (X^T D(\theta^{(l)}) X)^{-1} X^T D(\theta^{(l)}) \underbrace{Z}_{\text{response}} \quad \textcircled{1} \end{aligned}$$

Compared with least-square model,

$$\begin{aligned} \hat{\theta}_{LS} &= \arg \min_{\theta} \frac{1}{2} (\vec{y} - X\vec{\theta})^T W (\vec{y} - X\vec{\theta}) \\ \nabla \hat{\theta}_{LS} &= 0 \Rightarrow \vec{\theta} = (X^T W X)^{-1} X^T W \vec{y} \quad \textcircled{2} \end{aligned}$$

According to ①②  $Z \iff \vec{y}$   
equivalent

$$\text{therefore } \theta^{(l+1)} = \arg \min_{\theta} \frac{1}{2} (\tilde{y}(\theta^{(l)}) - X\theta)^T W(\theta^{(l)}) (\tilde{y}(\theta^{(l)}) - X\theta)$$

$$\left( \begin{array}{l} W(\theta^{(l)}) = D(\theta^{(l)}) \\ \tilde{y}(\theta^{(l)}) = X\theta^{(l)} + D(\theta^{(l)})^{-1} (y - b(\theta^{(l)})) \end{array} \right)$$

P2.

Proof.  $f: \mathbb{R}^n \rightarrow \mathbb{R}$  is a convex function.  $\Rightarrow \nabla^2 f \succeq 0$   
 $A\vec{x} + \vec{b} \in \text{dom } f$ .

$$f \text{ convex} \Rightarrow f(\lambda \vec{x} + \bar{\lambda} \vec{y}) \leq \lambda f(\vec{x}) + \bar{\lambda} f(\vec{y}), \quad \bar{\lambda} = 1 - \lambda$$

$$\left( \begin{array}{l} \text{To prove } g(\vec{x}) = f(A\vec{x} + \vec{b}) \text{ convex, we} \\ \text{hope. } g(\lambda \vec{x} + \bar{\lambda} \vec{y}) \leq \lambda g(\vec{x}) + \bar{\lambda} g(\vec{y}). \end{array} \right)$$

$$g(\vec{x}) = f(A\vec{x} + \vec{b})$$

$$g(\lambda \vec{x} + \bar{\lambda} \vec{y}) = f(A(\lambda \vec{x} + \bar{\lambda} \vec{y}) + \vec{b})$$

$$\underline{\lambda + \bar{\lambda} = 1}$$

$$= f(A(\lambda \vec{x} + \bar{\lambda} \vec{y}) + (\lambda + \bar{\lambda}) \vec{b})$$

$$= \underline{f(\lambda(A\vec{x} + \vec{b}) + \bar{\lambda}(A\vec{y} + \vec{b}))} \quad \textcircled{1}$$

$$f(\vec{x}) \text{ is convex, } \Rightarrow \textcircled{1} \leq \lambda f(A\vec{x} + \vec{b}) + \bar{\lambda} f(A\vec{y} + \vec{b})$$

$$= \lambda g(\vec{x}) + \bar{\lambda} g(\vec{y})$$

Therefore,  $g(\vec{x})$  is also a convex function.



P3. Given:

$x_*$  is a local minimizer of a convex function,

assume  $x_*$  could not be a global minimizer. We can find a point  $y \in \mathbb{R}$  such that  $f(y) < f(x_*)$ .  $x_*, y \in \Omega$

As a local minimizer,  $\exists \epsilon > 0$  s.t.  $f(x_*) \leq f(x)$  for  $\forall x$  s.t.  $\|x - x_*\|_2 < \epsilon$ . We now consider a point  $x$  in the line segment between  $x_*$  and  $y$ .

$$x = \bar{\lambda} x_* + \lambda y \quad (\bar{\lambda} = 1 - \lambda) \quad \lambda \in (0, 1) \quad x \in \Omega.$$

By convexity,  $f(x) = f(\lambda y + \bar{\lambda} x_*) \leq \lambda f(y) + \bar{\lambda} f(x_*)$  ①

$$f(y) < f(x_*) \quad f(x) \leq \lambda f(x_*) + (1 - \lambda) f(x_*) = f(x_*)$$

Since there exists a point  $f(x_*) \leq f(\lambda y + (1 - \lambda) x_*)$  ②

① ②  $f(x) = f(\lambda y + \bar{\lambda} x_*) \leq f(x_*) \leq f(\lambda y + \bar{\lambda} x_*)$

It is not possible, thus  $x_*$  is not a local minimizer any more, and should be a global minimizer.

Assume the minimizer is not unique  $\Rightarrow$  there exists two points  $a, b \in \Omega$  such that  $f(a), f(b)$  are both local minima.  $f(a) \geq f(b)$  WLOG

$$\lambda \in (0, 1) \quad 1 - \lambda > 0 \quad \lambda f(a) + (1 - \lambda) f(b) \leq \lambda f(a) + (1 - \lambda) f(a) = f(a)$$

$$f - \text{strictly convex} \Rightarrow f(\lambda a + (1 - \lambda) b) < \lambda f(a) + (1 - \lambda) f(b)$$

$$\text{Since } x = \lambda a + (1 - \lambda) b \in \Omega \Rightarrow f(x) < \lambda f(a) + (1 - \lambda) f(a) = f(a)$$

$f(a)$  is not a minimizer, which contradicts our assumption.

Therefore, if the function is strictly convex, there is a unique minimizer.

P4.

Given.

$g_i(x|\bar{x})$  is a majorizer of  $f_i(x)$ ,

we know.  $g_i(x|\bar{x}) \geq f_i(x)$  for  $\forall x$ .

$$g_i(\bar{x}|\bar{x}) = f_i(\bar{x})$$

$$\text{When } g(x|\bar{x}) = \sum_{i=1}^n g_i(x|\bar{x})$$

$$= g_1(x|\bar{x}) + g_2(x|\bar{x}) + \dots + g_i(x|\bar{x}) + \dots + g_n(x|\bar{x})$$

(each term is greater than  $f_i(x)$   $i=1, 2, \dots, n$  )

$$\geq f_1(x) + f_2(x) + \dots + f_i(x) + \dots + f_n(x)$$

$$= \sum_{i=1}^n f_i(x) \quad \textcircled{1}$$

$$g(\bar{x}|\bar{x}) = \sum_{i=1}^n g_i(\bar{x}|\bar{x})$$

$$= g_1(\bar{x}|\bar{x}) + \dots + g_i(\bar{x}|\bar{x}) + \dots + g_n(\bar{x}|\bar{x})$$

(each term is equal to  $f_i(\bar{x})$   $i=1, 2, \dots, n$  )

$$= f_1(\bar{x}) + \dots + f_i(\bar{x}) + \dots + f_n(\bar{x})$$

$$= \sum_{i=1}^n f_i(\bar{x}) \quad \textcircled{2}$$

Based on ①②.  $g(x|\bar{x}) = \sum_{i=1}^n g_i(x|\bar{x})$  is a majorizer of

$$f(x) = \sum_{i=1}^n f_i(x).$$

$$a) \ell(\omega) = -y_i \theta^T x_i + \log(H e^{\theta^T x_i})$$

$$g_i(\bar{\omega} | \bar{\theta}) = \ell_i(\bar{\omega}) - (y_i - \frac{e^{\bar{\omega}^T x_i}}{1 + e^{\bar{\omega}^T x_i}}) x_i^T (\bar{\omega} - \bar{\theta}) + \frac{1}{8} \|\bar{\omega} - \bar{\theta}\|^2$$

$$= \ell_i(\bar{\omega}) \quad \textcircled{1}$$

According to Taylor's theorem,

$$f(\omega) \leq f(\bar{\omega}) + \nabla f(\bar{\omega})^T (\omega - \bar{\omega}) + \frac{1}{2} (\omega - \bar{\omega})^T M(\omega) (\omega - \bar{\omega})$$

(M is positive definite matrix)

$$\nabla \ell_i(\omega) = - (y_i - \sigma(\omega^T x_i)) x_i^T$$

$$\ell_i(\bar{\omega}) \leq \ell_i(\bar{\omega}) - (y_i - \sigma(\bar{\omega}^T x_i)) x_i^T (\bar{\omega} - \bar{\theta}) + \frac{1}{2} (\bar{\omega} - \bar{\theta})^T M(\bar{\omega}) (\bar{\omega} - \bar{\theta})$$

Since,  $g_i(\omega | \bar{\theta}) = \ell_i(\omega) - (y_i - \sigma(\bar{\omega}^T x_i)) x_i^T (\omega - \bar{\omega}) + \frac{1}{8} (\omega - \bar{\omega})^T x_i x_i^T (\omega - \bar{\omega})$

We need to compare  $\frac{1}{2} (\omega - \bar{\omega})^T M(\omega) (\omega - \bar{\omega})$  with  $\frac{1}{8} (\omega - \bar{\omega})^T x_i x_i^T (\omega - \bar{\omega})$

$$M = \nabla^2 \ell_i(\omega) = x_i ( \sigma(\omega^T x_i) (1 - \sigma(\omega^T x_i)) ) x_i^T$$

$$\sigma(t) (1 - \sigma(t)) \xrightarrow{t = \bar{\omega}^T x_i} \frac{e^t}{(1 + e^t)^2}$$

$$\frac{e^t}{(1 + e^t)^2} - \frac{1}{4} = \frac{4e^t - (1 + e^t)^2}{4(1 + e^t)^2} = \frac{-(1 - e^t)^2}{4(1 + e^t)^2} \leq 0$$

$$\Rightarrow \frac{e^t}{(1 + e^t)^2} \leq \frac{1}{4} \quad M \mathbf{I} \leq \frac{1}{4} x_i x_i^T \Rightarrow \frac{1}{2} M \mathbf{I} \leq \frac{1}{8} x_i x_i^T$$

$$\Rightarrow \frac{1}{2} (\omega - \bar{\omega})^T M(\omega) (\omega - \bar{\omega}) \leq \frac{1}{8} (\omega - \bar{\omega})^T x_i x_i^T (\omega - \bar{\omega})$$

$$g_i(\omega | \bar{\theta}) \geq \ell_i(\bar{\omega}) \quad \textcircled{2}$$

Based on ① ②  $g_i(\omega | \bar{\theta})$  is a majorizer of  $\ell_i(\omega)$



b) To minimize.  $g_i(\omega|\bar{\omega}) = \ell_i(\bar{\omega}) - (y_i - b\bar{\omega}^T x_i) x_i^T (\omega - \bar{\omega}) + \frac{1}{8} (\omega - \bar{\omega})^T x_i x_i^T (\omega - \bar{\omega})$

equivalent to minimize

Construct a quadratic function

$$\Rightarrow -2 \underbrace{(y_i - b\bar{\omega}^T x_i)}_{\Delta} x_i^T (\omega - \bar{\omega}) + \frac{1}{8} (\omega - \bar{\omega})^T x_i x_i^T (\omega - \bar{\omega}) + \frac{1}{8} \Delta^2$$

$$= \left( \frac{1}{2\sqrt{2}} (\omega - \bar{\omega})^T x_i - \sqrt{2} \Delta \right)^2 = \frac{1}{8} \left( (\omega - \bar{\omega})^T x_i - 4\Delta \right)^2$$

For the sum of all minimized terms,  $\text{const}$

$$\min_{\omega} \sum_{i=1}^n g_i(\omega|\bar{\omega}) \Rightarrow \min_{\omega} \sum_{i=1}^n \frac{1}{8} \left( \bar{\omega}^T x_i - \underbrace{\bar{\omega}^T x_i + 2\Delta}_{\text{const}} \right)^2 \cdot (\Delta = y_i - b\bar{\omega}^T x_i)$$

$$= \frac{1}{2} \cdot (\bar{y}(\bar{\omega}) - x(\bar{\omega}))^T W (\bar{y}(\bar{\omega}) - x(\bar{\omega}))$$

$$\bar{y}(\bar{\omega}) = x\bar{\omega} + 4(y - b\bar{\omega}^T x)$$

$$W = \frac{1}{4} I.$$

c). According to Taylor's Theorem,

$f$  is a convex function and  $\underbrace{\nabla^2 f(\bar{\omega}) \leq \mu I}_{\text{PSD}}$

$$f(\omega) \leq f(\bar{\omega}) + \nabla f^T(\bar{\omega})(\omega - \bar{\omega}) + \frac{1}{2} (\omega - \bar{\omega})^T \mu I (\omega - \bar{\omega})$$

$$= f(\bar{\omega}) + \nabla f^T(\bar{\omega})(\omega - \bar{\omega}) + \frac{1}{2} \mu \|\omega - \bar{\omega}\|^2 = \ell(\omega) \quad \textcircled{1}$$

On the other hand.

$$\ell(\bar{\omega}) = f(\bar{\omega}) + \nabla f^T(\bar{\omega})(\bar{\omega} - \bar{\omega}) + \frac{1}{2} (\bar{\omega} - \bar{\omega})^T \mu I (\bar{\omega} - \bar{\omega}) \quad \textcircled{2}$$

Based on  $\textcircled{1}$  &  $\textcircled{2}$ ,  $\ell(\omega)$  is a majorizer of  $f(\omega)$ .

P6

$$a) \min \sum_{i=1}^n -y_i \beta^T x_i + \log(1 + e^{\beta^T x_i})$$

$$N(x) = \vec{r} = [1 \ 1 \ 1 \ \dots \ 0]^T$$

If we pick a vector.  $\vec{\beta}' = \vec{\beta} + \vec{r}t$  ( $t = \text{const } (\neq 0)$ )

$$\min \sum_{i=1}^n -y_i \beta^T x_i + \log(1 + e^{\beta^T x_i})$$

$$= \min \sum_{i=1}^n -y_i (\vec{\beta} + \vec{r}t)^T \vec{x}_i + \log(1 + e^{(\vec{\beta} + \vec{r}t)^T \vec{x}_i})$$

$$(\because \vec{r} \cdot t \cdot x_i = 0)$$

$$= \min \sum_{i=1}^n -y_i (\vec{\beta}^T \vec{x}_i) + \log(1 + e^{\vec{\beta}^T \vec{x}_i})$$

Therefore.  $\vec{\beta}' = \vec{\beta} + \vec{r}t$  ( $\vec{\beta}' \neq \vec{\beta}$ ) can also minimize the function  $\Rightarrow$  the solution is not unique.

$$b) \beta^{(l+1)} = \arg \min_{\beta} \frac{1}{2} (\tilde{y}(\beta^{(l)}) - X\beta)^T W(\beta^{(l)}) (\tilde{y}(\beta^{(l)}) - X\beta)$$

$$\left( \begin{array}{l} \tilde{y}(\beta^{(l)}) = X\beta^{(l)} + D(\beta^{(l)}) (y - b(\beta^{(l)})) \\ W(\beta^{(l)}) = D(\beta^{(l)}) \\ D_{ii}(\beta^{(l)}) = \frac{b_i(\beta^{(l)}) (1 - b_i(\beta^{(l)}))}{\tilde{y}_i(\beta^{(l)})} \end{array} \right)$$

$$g(\beta^{(l)}) = \frac{1}{2} (\tilde{y}(\beta^{(l)}) - X\beta)^T W(\beta^{(l)}) (\tilde{y}(\beta^{(l)}) - X\beta) + \lambda r^T \beta$$

$$\nabla g(\beta) = -X^T W (\tilde{y}(\beta) - X\beta) + \lambda \vec{r} = \vec{0}$$

$$\nabla g(\lambda) = \vec{r}^T \vec{\beta} = 0$$

$$\begin{bmatrix} X^T W X & \vec{r} \\ \vec{r}^T & 0 \end{bmatrix} \begin{bmatrix} \vec{\beta} \\ \lambda \end{bmatrix} = \begin{bmatrix} X^T W \tilde{y} \\ 0 \end{bmatrix}$$

$$L(\beta) = \frac{1}{2} (\tilde{y}(\beta^{(l)}) - X\beta)^T W(\beta^{(l)}) (\tilde{y}(\beta^{(l)}) - X\beta) + \lambda r^T \beta$$

(see code)

c) Similar to b), only need to change.  $W(\beta^{(l)}) = \frac{1}{4} I$ ,

$$\tilde{y}(\beta^{(l)}) = X\beta^{(l)} + 4(y - b(\vec{\theta}^T X))$$

(see code)

d) see code

P6 (see code)

Comments for part e:

Compared with the win/loss percentage, the ranking estimation using logistic regression basically has the similar results but not exactly the same. As seen from the ranking results, some teams' betas are very close, which leads to a different order among 2 or 3 teams from the win/loss ranking. This makes sense because the score of each team should be independent but our derived non-linear model takes the score differences ( $y = 1,0$ ) between 2 teams into consideration. The underlying assumption is that these 30 teams may have correlations, leading to the differences between our results and the win/loss ranking.

Compared with the results from linear model we applied in the last homework, logistic regression seems to have a better performance. (Logistic regression predicts 12 correct rankings but linear regression predicts 11 correct rankings). It seems that for this NBA ranking, logistic regression performs well since score differences are converted into (0,1) binary number, which makes the classification more robust to numerical factors.