**1.**

(a)

$$\frac{\partial}{\partial \theta_j} \ell_i(\theta) = \frac{\partial}{\partial \theta_j}\left[\frac{1}{2}\left(x_i^T\theta - Y_i\right)^2\right]$$

$$= \frac{\partial}{\partial \theta_j}\left[\frac{1}{2}\left(x_{i1}\theta_1 + \cdots + x_{ip}\theta_p - Y_i\right)^2\right]$$

$$= \frac{1}{2}\cdot 2\left(x_{i1}\theta_1 + \cdots + x_{ip}\theta_p - Y_i\right)\cdot x_{ij}$$

$$= \left(x_i^T\theta - Y_i\right)x_{ij}$$

$$\nabla_\theta \ell_i(\theta) = \begin{bmatrix} \frac{\partial}{\partial \theta_1}\ell_i(\theta) \\ \vdots \\ \frac{\partial}{\partial \theta_p}\ell_i(\theta) \end{bmatrix} = \begin{bmatrix} (x_i^T\theta - Y_i)x_{i1} \\ \vdots \\ (x_i^T\theta - Y_i)x_{ip} \end{bmatrix} = (x_i^T\theta - Y_i)x_i \qquad \square$$

(b)

$$\nabla_\theta L(\theta) = \nabla_\theta\left[\frac{1}{2}\|X\theta - Y\|^2\right]$$

$$= \nabla_\theta\left[\frac{1}{2}\sum_{i=1}^{N}(x_i\theta - Y_i)^2\right]$$

$$= \sum_{i=1}^{N}\nabla_\theta \ell_i(\theta)$$

$$= \sum_{i=1}^{N}\left(x_i^T\theta - Y_i\right)x_i$$

$$= \sum_{i=1}^{N}x_i\left(x_i^T\theta - Y_i\right)$$

$$= \begin{bmatrix} x_1, & \cdots, & x_N \end{bmatrix}\begin{bmatrix} x_1^T\theta - Y_1 \\ \vdots \\ x_N^T\theta - Y_N \end{bmatrix}$$

$$= X^T\left(X\theta - Y\right) \qquad \square$$

2.

$$f(\theta) = \frac{\theta^2}{2}$$

$$f'(\theta) = \theta$$

$$\theta^{k+1} = \theta^k - \alpha f'(\theta^k) = \theta^k - \alpha\theta^k = (1-\alpha)\theta^k$$

$$\therefore \theta^n = (1-\alpha)^n \theta^0, \quad k = 0, 1, 2, \cdots$$

$$\alpha > 2 \Rightarrow |1-\alpha| > 1$$

$\therefore$ If $\theta^0 \neq 0$, $\alpha > 2$,

then $\theta^n$ diverges as $n \to \infty$. $\square$

3.

$$\nabla f(\theta) = X^T(X\theta - Y)$$

$$\theta^{k+1} = \theta^k - \alpha \nabla f(\theta^k)$$

$$= \theta^k - \alpha X^T(X\theta^k - Y)$$

$$= (I - \alpha X^T X)\theta^k + \alpha X^T Y$$

$$\theta^{k+1} - (X^T X)^{-1} X^T Y = (I - \alpha X^T X)\theta^k + \alpha X^T Y - (X^T X)^{-1} X^T Y$$

$$= (I - \alpha X^T X)(\theta^k - (X^T X)^{-1} X^T Y)$$

By letting $(X^T X)^{-1} X^T Y = \theta^*$,

$$\theta^{k+1} - \theta^* = (I - \alpha X^T X)(\theta^k - \theta^*)$$

$$\theta^n - \theta^* = (I - \alpha X^T X)^n (\theta^0 - \theta^*)$$

$X^T X$ is symmetric psd and so diagonalizable.

$$X^T X = Q^T \Lambda Q, \quad Q \in \mathbb{R}^{p \times p} : \text{orthogonal}, \quad Q = \begin{bmatrix} q_1 \\ \vdots \\ q_p \end{bmatrix}$$

$$\Lambda = \begin{bmatrix} \lambda_1 & & 0 \\ & \ddots & \\ 0 & & \lambda_p \end{bmatrix}, \quad \lambda_1 \geq \lambda_2 \geq \cdots \geq \lambda_p \geq 0 : \text{eigenvalues of } X^T X$$

$$(I - \alpha X^T X)^n = (Q^T I Q - \alpha Q^T \Lambda Q)^n$$

$$= [Q^T(I - \alpha \Lambda)Q]^n$$

$$= Q^T \begin{bmatrix} (1 - \alpha \lambda_1)^n & & 0 \\ & \ddots & \\ 0 & & (1 - \alpha \lambda_p)^n \end{bmatrix} Q = \sum_{i=1}^{p} q_i^T (1 - \alpha \lambda_i)^n q_i$$

Since $\alpha > \dfrac{2}{\rho(X^T X)} = \dfrac{2}{\lambda_1} \Rightarrow 1 - \alpha \lambda_1 < -1$

$(1 - \alpha \lambda_1)^n$ diverges as $n \to \infty$

$\therefore \theta^n = (I - \alpha X^T X)^n (\theta^0 - \theta^*) + \theta^*$ diverges

if first element of $\theta^0 - \theta^*$ is non zero.

$\therefore \langle \theta^n \rangle$ diverges for most starting points $\theta_0$.

□