

MathDNN Homework 1

Department of Computer Science and Engineering
2021-16988 Jaewan Park

Problem 1

(a) Let $X_i = \begin{bmatrix} x_{i1} \\ \vdots \\ x_{ip} \end{bmatrix}$ and $\theta = \begin{bmatrix} \theta_1 \\ \vdots \\ \theta_p \end{bmatrix}$, then for $j = 1, \dots, p$,

$$\begin{aligned} \frac{\partial}{\partial \theta_j} l_i(\theta) &= \frac{\partial}{\partial \theta_j} \left[\frac{1}{2} (X_i^\top \theta - Y_i)^2 \right] \\ &= \frac{\partial}{\partial \theta_j} \left[\frac{1}{2} (x_{i1}\theta_1 + \dots + x_{ip}\theta_p - Y_i)^2 \right] \\ &= (x_{i1}\theta_1 + \dots + x_{ip}\theta_p - Y_i) x_{ij} \\ &= (X_i^\top \theta - Y_i) x_{ij}. \end{aligned}$$

Therefore

$$\begin{aligned} \nabla_\theta l_i(\theta) &= \begin{bmatrix} \frac{\partial}{\partial \theta_1} l_i(\theta) \\ \vdots \\ \frac{\partial}{\partial \theta_p} l_i(\theta) \end{bmatrix} \\ &= \begin{bmatrix} (X_i^\top \theta - Y_i) x_{i1} \\ \vdots \\ (X_i^\top \theta - Y_i) x_{ip} \end{bmatrix} = (X_i^\top \theta - Y_i) \begin{bmatrix} x_{i1} \\ \vdots \\ x_{ip} \end{bmatrix} \\ &= (X_i^\top \theta - Y_i) X_i. \end{aligned}$$

(b) Let $X_i = \begin{bmatrix} x_{i1} \\ \vdots \\ x_{ip} \end{bmatrix}$ and $\theta = \begin{bmatrix} \theta_1 \\ \vdots \\ \theta_p \end{bmatrix}$, then for $j = 1, \dots, p$,

$$\begin{aligned} \frac{\partial}{\partial \theta_j} \mathcal{L}(\theta) &= \frac{\partial}{\partial \theta_j} \left[\frac{1}{2} \|X\theta - Y\|^2 \right] \\ &= \frac{\partial}{\partial \theta_j} \left[\frac{1}{2} \left\| \begin{bmatrix} X_1^\top \theta - Y_1 \\ \vdots \\ X_N^\top \theta - Y_N \end{bmatrix} \right\|^2 \right] = \frac{\partial}{\partial \theta_j} \left[\sum_{i=1}^N \frac{1}{2} (X_i^\top \theta - Y_i)^2 \right] = \frac{\partial}{\partial \theta_j} \left[\sum_{i=1}^N l_i(\theta) \right] \\ &= \sum_{i=1}^N \frac{\partial}{\partial \theta_j} l_i(\theta) = \sum_{i=1}^N (X_i^\top \theta - Y_i) x_{ij}. \end{aligned}$$

Therefore

$$\begin{aligned}
\nabla_{\theta} \mathcal{L}(\theta) &= \begin{bmatrix} \sum_{i=1}^N (X_i^{\top} \theta - Y_i) x_{i1} \\ \vdots \\ \sum_{i=1}^N (X_i^{\top} \theta - Y_i) x_{ip} \end{bmatrix} \\
&= \begin{bmatrix} (X_1^{\top} \theta - Y_1) x_{11} + (X_2^{\top} \theta - Y_2) x_{21} + \cdots + (X_N^{\top} \theta - Y_N) x_{N1} \\ \vdots \\ (X_1^{\top} \theta - Y_1) x_{1p} + (X_2^{\top} \theta - Y_2) x_{2p} + \cdots + (X_N^{\top} \theta - Y_N) x_{Np} \end{bmatrix} \\
&= \sum_{i=1}^N X_i (X_i^{\top} \theta - Y_i) = \sum_{i=1}^N (X^{\top})_{:,i} (X\theta - Y)_i \\
&= X^{\top} (X\theta - Y).
\end{aligned}$$

Problem 2

Since $f(\theta) = \theta^2/2$, we can get $f'(\theta) = \theta$ and the iteration equation can be written as

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

Thus the n th iteration gives

$$\theta^n = (1 - \alpha)\theta^{n-1} = \cdots = (1 - \alpha)^n \theta^0.$$

Therefore if $\alpha > 2$ and $\theta^0 \neq 0$, θ^n diverges since it is a geometric sequence where its ratio is less than -1 .

Problem 3

Since $f(\theta) = \frac{1}{2} \|X\theta - Y\|^2$, we can get $\nabla f(\theta) = X^{\top}(X\theta - Y)$ and the iteration equation can be written as

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

Assume $X^{\top}X$ is invertible and let $\theta^* = (X^{\top}X)^{-1}X^{\top}Y$, then

$$\begin{aligned}
\theta^{k+1} - \theta^* &= \theta^k - \alpha X^{\top}(X\theta^k - Y) - (X^{\top}X)^{-1}X^{\top}Y \\
&= \theta^k - \alpha X^{\top}X\theta^k + \alpha X^{\top}Y - (X^{\top}X)^{-1}X^{\top}Y \\
&= (I - \alpha X^{\top}X)\theta^k + (\alpha X^{\top}X - I)(X^{\top}X)^{-1}X^{\top}Y = (I - \alpha X^{\top}X)\theta^k + (\alpha X^{\top}X - I)\theta^* \\
&= (I - \alpha X^{\top}X)(\theta^k - \theta^*)
\end{aligned}$$

where $I \in \mathbb{R}^{p \times p}$ is an identity matrix. Thus the n th iteration gives

$$\theta^n - \theta^* = (I - \alpha X^{\top}X)(\theta^{n-1} - \theta^*) = \cdots = (I - \alpha X^{\top}X)^n(\theta^0 - \theta^*).$$

Since $X^\top X$ is a symmetric matrix and a positive matrix, it has p real positive eigenvalues and is diagonalizable by them. Let $\lambda_1 \geq \lambda_2 \geq \dots \geq \lambda_p \geq 0$ the eigenvalues and M the matrix with the eigenvectors as its columns, then

$$M^{-1}X^\top XM = \text{diag}(\lambda_1, \dots, \lambda_p), \quad X^\top X \sim \text{diag}(\lambda_1, \dots, \lambda_p)$$

and gradually

$$\begin{aligned} I - \alpha X^\top X &= I - \alpha M \text{diag}(\lambda_1, \dots, \lambda_p) M^{-1} \\ &= M I M^{-1} - M(\alpha \text{diag}(\lambda_1, \dots, \lambda_p)) M^{-1} \\ &= M(I - \alpha \text{diag}(\lambda_1, \dots, \lambda_p)) M^{-1} = M \text{diag}(1 - \alpha \lambda_1, \dots, 1 - \alpha \lambda_p) M^{-1} \end{aligned}$$

so $I - \alpha X^\top X \sim \text{diag}(1 - \alpha \lambda_1, \dots, 1 - \alpha \lambda_p)$. Therefore

$$\begin{aligned} (I - \alpha X^\top X)^n &= M \text{diag}(1 - \alpha \lambda_1, \dots, 1 - \alpha \lambda_p)^n M^{-1} \\ &= M \text{diag}((1 - \alpha \lambda_1)^n, \dots, (1 - \alpha \lambda_p)^n) M^{-1}. \end{aligned}$$

If $\alpha > 2/\rho(X^\top X)$, we can obtain $1 - \alpha \lambda_1 < -1$ since $\rho(X^\top X) = \lambda_1 \geq 0$. Therefore $(1 - \alpha \lambda_1)^n$ diverges, and consequently $(I - \alpha X^\top X)^n$, $\theta^n - \theta^*$, and θ^n also diverge.