P1

(a)

$$l_i(\theta) = \frac{1}{2}\left(X_i^T\theta - Y_i\right)^2 = \frac{1}{2}\left(\sum_j X_{ij}\theta_j - Y_i\right)^2$$

$$\frac{\partial}{\partial\theta_k}l_i(\theta) = \frac{\partial}{\partial\theta_k}\left(\frac{1}{2}\left(\sum_j X_{ij}\theta_j - Y_i\right)^2\right)$$

$$= \left(\sum_j X_{ij}\theta_j - Y_i\right)X_{ik} = \left(X_i^T\theta - Y_i\right)X_{ik}$$

$$\therefore \quad \nabla_\theta l_i(\theta) = \left(X_i^T\theta - Y_i\right)\left(X_{i1} \quad X_{i2} \quad \cdots \quad X_{ip}\right)^T = \left(X_i^T\theta - Y_i\right)X_i \qquad ✗$$

(b)

$$\mathcal{L}(\theta) = \frac{1}{2}\|X\theta - Y\|^2 = \frac{1}{2}\sum_{i=1}^{N}(X\theta - Y)_i^2 = \frac{1}{2}\sum_{i=1}^{N}\left(\sum_{j=1}^{p}X_{ij}\theta_j - Y_i\right)^2 = \sum_{i=1}^{N}\sum_{j=1}^{p}\frac{1}{2}\left(X_{ij}\theta_j - Y_i\right)^2$$

$$\therefore \quad \frac{\partial}{\partial\theta_k}\mathcal{L}(\theta) = \sum_{i=1}^{N}\left(\sum_{j=1}^{p}X_{ij}\theta_j - Y_i\right)X_{ik} = \sum_{i=1}^{N}(X\theta - Y)_i X_{ik}$$

$$\therefore \quad \nabla_\theta \mathcal{L}(\theta) = \left((X\theta - Y)^T X\right)^T = X^T(X\theta - Y) \qquad ✗$$

P2

$$\theta_{k+1} = \theta_k - \alpha f'(\theta_k) = \theta_k - \alpha\theta_k = (1-\alpha)\theta_k$$

$$= (1-\alpha)^{k+1}\theta_0$$

if $|1-\alpha| > 1$, it diverges $\quad \therefore \alpha > 2$ makes GD diverge.

P3

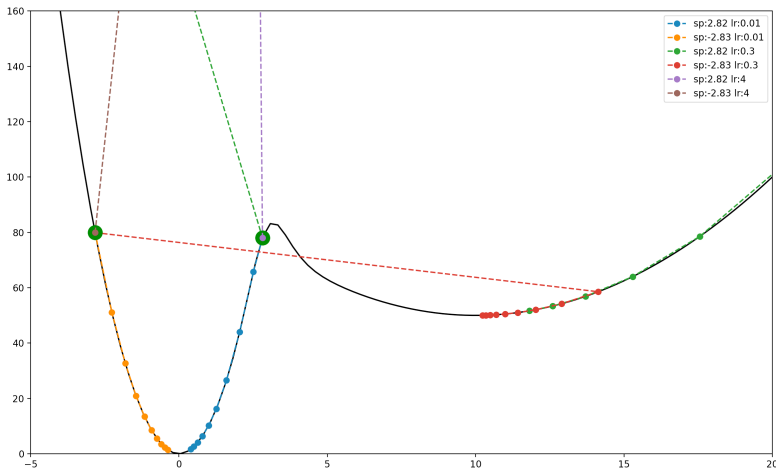$$\min_{\theta \in \mathbb{R}^p} f(\theta), \quad f(\theta) = \frac{1}{2}\|X\theta - Y\|^2$$

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

$$\theta_{k+1} - \theta_* = \theta_k - \alpha X^T(X\theta_k - Y) - \theta_*$$

$$= \theta_k - \theta_* - \alpha X^T X(\theta_k - \theta_*) \qquad (\because X\theta_* = Y)$$

$$= \theta_k - \theta_* - \alpha\rho(\theta_k - \theta_*)$$

$$= (1-\alpha\rho)(\theta_k - \theta_*)$$

$$\therefore \quad \|\theta_{k+1} - \theta_*\| = |1-\alpha\rho|\;\|\theta_k - \theta_*\|$$

Thus, using the result of problem 2, $\quad \alpha\rho > 2 \iff \alpha > \frac{2}{\rho}$ makes $\{\theta_k\}$ diverge.

# P4



You can see with learning rate (lr) 0.01, it approaches to sharp minima as iteration goes on. Meanwhile, with lr 0.3, they converges to wide minima. On the other hand, with lr 4, gradient descent diverges, thus aren't shown in the figure even after their first iteration.

# P5

```
› python3 -u "/Users/yongjae/Desktop/SNU/Mathematical-Foundations-of-DNN/hw/week1/conv1D.py"
5.098441152369599
```