

HW3

Name: Zhou Liu

Q1, a

Gradient descent is an algorithm for finding the minimizer(hence the minimum value) of a scalar valued multivariate function.

It is reasonable because the gradient of a function at a given point indicates the direction of the steepest increase of the function, thus, by taking steps in the opposite direction of the gradient, gradient descent naturally moves towards the minimum of the function.

Q1, b

No.

There are some situations where a step of gradient descent cannot result in a decrease of the objective function. For example, if the learning rate is too large, the algorithm may overshoot the minimum and end up with a higher function value in subsequent steps.

Q1, c

Under the statistical learning framework, we use maximum likelihood estimation to estimate the unknown parameters of a statistical model based on samples of the model, which is to find the values of the parameters that make the observed data as likely as possible.

For a regression problem, when calculating the likelihood of data, we found that to maximize such a function is equivalent to minimize the square loss, thus, square loss is a reasonable choice.

Question 2, a.

Let $\theta \in \mathbb{R}^d$, $x \in \mathbb{R}^d$.

Model: $y_i = x_i^t \theta + \epsilon_i$ w/ $\epsilon_i \sim \text{Laplace}(0, b)$.

Data: $D = \{(x_i, y_i)\}_{i=1, 2, \dots, n}$.

Estimate by maximum likelihood

pdf of ϵ_i is $\frac{1}{2b} e^{-\frac{|z|}{b}}$ over $z \in \mathbb{R}$.

likelihood of data using $\epsilon_i = y_i - x_i^t \theta$.

$$L(\theta) = \prod_{i=1}^n \frac{1}{2b} e^{-\frac{|y_i - x_i^t \theta|}{b}} \quad \text{by independence of data.}$$

$$\log L(\theta) = -\sum_{i=1}^n \frac{|y_i - x_i^t \theta|}{2b^2} + \text{terms constant in } \theta$$

maximizing data likelihood \Leftrightarrow

$$\text{minimizing } \sum_{i=1}^n \frac{|y_i - x_i^t \theta|}{2b^2}$$

$$\text{i.e. } \min_{\theta} \frac{1}{2b^2} \sum_{i=1}^n |y_i - x_i^t \theta|.$$

$$\Rightarrow \min_{\theta} \sum_{i=1}^n |y_i - x_i^t \theta| \text{ as } 2b^2 > 0. \quad \frac{1}{2b^2} > 0.$$

Above shown the maximum likelihood estimate of θ is given by $\min_{\theta} \sum_{i=1}^n |y_i - x_i^t \theta|$.

The loss function is

$$l(\hat{y}, y) = |\hat{y} - y|.$$

while we can still use square loss $|\hat{y} - y|^2$ as loss function because $|a| = |b| \Rightarrow a^2 = b^2$.

And square loss is more convenient for differentiation.

Question 2, b.

Estimate θ by maximum likelihood.

$$L(\theta) = \prod_{i=1}^n P(y=y_i | X^{(i)}; \theta)$$

$$\begin{aligned} & P(y=y_i | X^{(i)}; \theta) \\ &= f(X^{(i)}; \theta)_{y_i} \\ &= f(X^{(i)}; \theta)_1^{y_i=1} f(X^{(i)}; \theta)_2^{y_i=2} \dots f(X^{(i)}; \theta)_{y_i}^{y_i=y_i} \dots \quad (*) \end{aligned}$$

As $y_i \in \{1, 2, \dots, K\}$.

$$(*) = \prod_{k=1}^K f(X^{(i)}; \theta)_k^{1_{y_i=k}}$$

$$\therefore L(\theta) = \prod_{i=1}^n \prod_{k=1}^K f(X^{(i)}; \theta)_k^{1_{y_i=k}}$$

$$\begin{aligned} \log L(\theta) &= \sum_{i=1}^n \log \prod_{k=1}^K f(X^{(i)}; \theta)_k^{1_{y_i=k}} \\ &= \sum_{i=1}^n \sum_{k=1}^K 1_{y_i=k} \log f(X^{(i)}; \theta)_k \end{aligned}$$

to maximize $\log L(\theta)$ is to minimize

$$-\sum_{i=1}^n \sum_{k=1}^K 1_{y_i=k} \log f(X^{(i)}; \theta)_k. \quad \text{Q.E.D.}$$

Question 3, a

$$f(z) = \frac{e^z}{e^z + 1} = 1 - \frac{1}{e^z + 1}$$

$$= 1 - (e^z + 1)^{-1}$$

$$f'(z) = (e^z + 1)^{-2} \cdot (e^z + 1)'$$

$$= (e^z + 1)^{-2} \cdot e^z$$

$$= \frac{e^z}{(e^z + 1)^2}$$

$$= \frac{e^z}{e^z + 1} \cdot \frac{1}{e^z + 1}$$

$$= f(z) (1 - f(z))$$

Q.E.D.

Question 3, b

$$f(\theta) = \sum_{i=1}^n -y_i \log \hat{y}_i - (1-y_i) \log (1-\hat{y}_i)$$

$$\hat{y}_i = b(\theta x_i) \quad b'(z) = b(z)(1-b(z))$$

$$\frac{\partial f}{\partial \theta_j} = \sum_{i=1}^n (-y_i \log \hat{y}_i - (1-y_i) \log (1-\hat{y}_i))$$

$$= \sum_{i=1}^n \left(-y_i \frac{\partial \log \hat{y}_i}{\partial \theta_j} - (1-y_i) \frac{\partial \log (1-\hat{y}_i)}{\partial \theta_j} \right)$$

$$= \sum_{i=1}^n \left(-y_i \frac{\partial \log \hat{y}_i}{\partial \hat{y}_i} \cdot \frac{\partial \hat{y}_i}{\partial \theta_j} - (1-y_i) \frac{\partial \log (1-\hat{y}_i)}{\partial \hat{y}_i} \cdot \frac{\partial \hat{y}_i}{\partial \theta_j} \right)$$

$$= \sum_{i=1}^n \frac{\partial \hat{y}_i}{\partial \theta_j} \left(-y_i \frac{1}{\hat{y}_i} - (1-y_i) \frac{-1}{1-\hat{y}_i} \right)$$

$$= \sum_{i=1}^n \frac{\partial \hat{y}_i}{\partial \theta_j} \cdot \frac{\hat{y}_i - y_i}{\hat{y}_i (1-\hat{y}_i)}$$

$$\begin{aligned}
 \frac{\partial \hat{y}_i}{\partial \theta_j} &= \frac{\partial b(\theta x_i)}{\partial \theta_j} = \frac{\partial b(\theta x_i)}{\partial \theta x_i} \cdot \frac{\partial \theta x_i}{\partial \theta_j} \\
 &= b(\theta x_i) (1 - b(\theta x_i)) \cdot \frac{\partial \sum_{j=1}^d \theta_j x_{ij}}{\partial \theta_j} \\
 &= \hat{y}_i (1 - \hat{y}_i) \cdot x_{ij}
 \end{aligned}$$

$$\begin{aligned}
 \therefore \frac{\partial f}{\partial \theta_j} &= \sum_{i=1}^n \hat{y}_i (1 - \hat{y}_i) \cdot x_{ij} \cdot \frac{\hat{y}_i - y_i}{\hat{y}_i (1 - \hat{y}_i)} \\
 &= \sum_{i=1}^n x_{ij} (\hat{y}_i - y_i)
 \end{aligned}$$

$$\begin{aligned}
 \therefore \frac{\partial f}{\partial \theta} &= \sum_{i=1}^n x_i (\hat{y}_i - y_i) \\
 &= x \cdot (\hat{y} - y)
 \end{aligned}$$