

Problem 1

Represent the derivative of the following scalar functions with respect to $\mathbf{X} \in \mathbb{R}^{D \times D}$.

(a) $f(\mathbf{X}) = \text{tr}(\mathbf{X}^2)$. Here, $\text{tr}(A)$ is the trace of a square matrix A .

(b) $g(\mathbf{X}) = \text{tr}(\mathbf{X}^3)$.

(c) $h(\mathbf{X}) = \text{tr}(\mathbf{X}^k)$ for $k \in \mathbb{N}$.

$$(a) \quad f(\mathbf{X}) = \text{tr}(\mathbf{X}^2) = \sum_{i=1}^D \sum_{j=1}^D X_{ij} X_{ji}$$

$$\frac{df(\mathbf{X})}{d\mathbf{X}} = \frac{\partial f}{\partial X_{kl}} X_{ji} + X_{ij} \frac{\partial f}{\partial X_{kl}} \quad \text{ } l=i \text{ \& \& } k=j \text{인 항만 살아남음}$$

$$= X_{lk} + X_{lk} = 2X_{lk} = 2(X_{kl})^T = 2\mathbf{X}^T$$

$$(b) \quad g(\mathbf{X}) = \text{tr}(\mathbf{X}^3) = \sum_{i=1}^D \sum_{j=1}^D \sum_{k=1}^D X_{ij} X_{jk} X_{ki}$$

$$\frac{dg(\mathbf{X})}{d\mathbf{X}} = \frac{\partial g}{\partial X_{lm}} \sum_{k=1}^D X_{jk} X_{ki} + \frac{\partial g}{\partial X_{lm}} \sum_{i=1}^D X_{ij} X_{ki} + \frac{\partial g}{\partial X_{lm}} \sum_{i=1}^D X_{ij} X_{jk}$$

$$= \sum_{k=1}^D X_{mk} X_{kl} + \sum_{i=1}^D X_{il} X_{mi} + \sum_{j=1}^D X_{mj} X_{jl}$$

$$= \sum_{k=1}^D (X_{lk})^T \cdot (X_{km})^T + \sum_{i=1}^D (X_{li})^T \cdot (X_{im})^T + \sum_{j=1}^D (X_{lj})^T \cdot (X_{jm})^T$$

$$= (\mathbf{X}^2)^T + (\mathbf{X}^2)^T + (\mathbf{X}^2)^T = 3(\mathbf{X}^2)^T$$

$$(c) \quad h(\mathbf{X}) = \text{tr}(\mathbf{X}^k)$$

$$\text{tr}[\mathbf{X}^k] = \sum_{a_1=1}^D \sum_{a_2=1}^D \cdots \sum_{a_k=1}^D \overbrace{(X_{a_1 a_2} X_{a_2 a_3} X_{a_3 a_4} \cdots X_{a_k a_1})}^{k\text{TH}}$$

$$\frac{dh(\mathbf{X})}{d\mathbf{X}} = \sum_{a_3=1}^D \cdots \sum_{a_k=1}^D \mathbf{X}^{k-1} + \sum_{a_1=1}^D \sum_{a_2=1}^D \cdots \sum_{a_k=1}^D \mathbf{X}^{k-1} + \cdots$$

$$\underbrace{\hspace{10em}}_{k\text{TH}}$$

$$= k(\mathbf{X}^{k-1})^T$$

Problem 2

In order to alleviate overfitting in logistic regression, regularization technique can be used with the L_2 -norm of the weight parameter, $\|\mathbf{w}\|^2 = \sum_{i=1}^D w_i^2$. When we are given $\mathbf{x}_1, \dots, \mathbf{x}_N \in \mathbb{R}^D$ and $y_1, \dots, y_N \in \{0, 1\}$ for training set, we want to derive the update rule for $\mathbf{w} \in \mathbb{R}^D$ in order to minimize the following loss function $L(\mathbf{w})$ having L_2 -norm,

$$L(\mathbf{w}) = \sum_{i=1}^N \left(-y_i \ln f(\mathbf{x}_i; \mathbf{w}) - (1 - y_i) \ln(1 - f(\mathbf{x}_i; \mathbf{w})) \right) + \boxed{\|\mathbf{w}\|^2} \quad (1)$$

$$f(\mathbf{x}; \mathbf{w}) = \frac{1}{1 + \exp(-\mathbf{w}^\top \mathbf{x})} \quad \text{L2-norm} \quad (2)$$

(1) Derive the update rule for \mathbf{w} when we use gradient descent.

(2) Discuss the effect of L_2 -norm regularization.

$$\begin{aligned} \therefore L(\mathbf{w}) &= \sum_{i=1}^N \left(-y_i \ln(f(\mathbf{x}_i; \mathbf{w})) - (1 - y_i) \ln(1 - f(\mathbf{x}_i; \mathbf{w})) \right) + \|\mathbf{w}\|^2 \\ &= -y \ln f - (1 - y) \ln(1 - f) + \mathbf{w}^\top \mathbf{w} \end{aligned}$$

$$\begin{aligned} \left[\frac{d}{d\mathbf{w}} (\mathbf{w}^\top \mathbf{w}) \right]_i &= \frac{d}{d\mathbf{w}_i} (\sum w_i^2) = \frac{d}{d\mathbf{w}_i} (w_1^2 + w_2^2 + \dots + w_i^2 + \dots + w_p^2) \\ &= \frac{d}{d\mathbf{w}_i} (w_i^2) = 2w_i \quad \therefore \frac{d}{d\mathbf{w}} (\mathbf{w}^\top \mathbf{w}) = 2\mathbf{w} \text{ --- (A)} \end{aligned}$$

$$\begin{aligned} \frac{df}{d\mathbf{w}} &= \frac{-\exp(-\mathbf{w}^\top \mathbf{x})(-\mathbf{x})}{(1 + \exp(-\mathbf{w}^\top \mathbf{x}))^2} = \frac{\mathbf{x} \cdot \exp(-\mathbf{w}^\top \mathbf{x})}{(1 + \exp(-\mathbf{w}^\top \mathbf{x}))^2} \\ &= \frac{1}{1 + \exp(-\mathbf{w}^\top \mathbf{x})} \cdot \frac{\exp(-\mathbf{w}^\top \mathbf{x})}{1 + \exp(-\mathbf{w}^\top \mathbf{x})} \cdot \mathbf{x} \end{aligned}$$

$$= f(1 - f) \cdot \mathbf{x} \text{ --- (B)} \quad 2$$

$$\frac{d}{dw} L(w) = -y \cdot \frac{1}{f} \cdot \frac{df}{dw} - (1-y) \cdot \frac{1}{1-f} \left(-\frac{df}{dw} \right) + 2w - \epsilon$$

①과 ②를 ③에 대입하면,

$$\begin{aligned} \frac{d}{dw} L(w) &= \sum_{i=1}^N \left[-y_i \cdot \frac{1}{f(x_i; w)} \cdot \cancel{f(x_i; w)} \cdot (1 - \cancel{f(x_i; w)}) \cdot x_i \right. \\ &\quad \left. + (1-y_i) \cdot \frac{1}{1 - \cancel{f(x_i; w)}} \cdot \cancel{f(x_i; w)} (\cancel{f(x_i; w)} - 1) \cdot x_i \right] + 2w \end{aligned}$$

$$= \sum_{i=1}^N (-y_i (1 - f(x_i; w)) \cdot x_i + (1-y_i) \cdot f(x_i; w) \cdot x_i) + 2w$$

$$\therefore w_{t+1} = w_t - \eta \cdot \frac{dL}{dw}$$

$$= w_t - \eta \cdot \left[\sum_{i=1}^N (-y(1-f)x + (1-y)f \cdot x) + 2w \right]$$

(2) 기본적으로 regularization을 한다는 것의 의미는
weight가 작아지도록 학습을 하게 만드는 것
의미한다.

$$L_1 \quad \text{loss} = \sum |w|$$

$$L_2 \quad \text{loss} = \sqrt{\sum |w|^2}$$

$$\text{ex) } a = \{0.3, -0.3, 0.4\}, \quad b = \{0.5, -0.5, 0\}$$

$$\text{i) } L_1 \quad \text{loss} : \sum |a| = 0.3 + 0.3 + 0.4 = 1$$

$$\sum |b| = 0.5 + 0.5 + 0 = 1$$

$$\text{ii) } L_2 \quad \text{loss} : \sum |a|^2 = 0.09 + 0.09 + 0.16 = 0.34$$

$$\sum |b|^2 = 0.25 + 0.25 + 0 = 0.50$$

반대로, weight가 너무 큰 값을 가지면
model이 학습한 data를 외운다. (\rightarrow 일반화 성능 저하)

L_2 norm regularization은 모델의 학습에서

local minimum의 영향을 덜 받게 (특이점의 영향 \downarrow)

학습을 하여 일반화 성능을 높인다고 할 수 있다.

Problem 3

Consider the following function,

$$f(\mathbf{x}; \mathbf{w}) = \frac{1 - \exp(-\mathbf{w}^T \mathbf{x})}{1 + \exp(-\mathbf{w}^T \mathbf{x})}, \quad (3)$$

having the shape in Fig. 1.

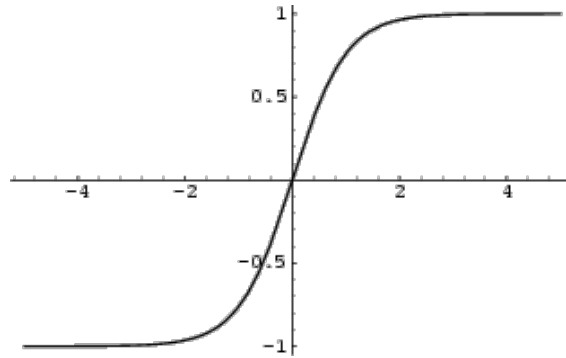


Figure 1: $f(\mathbf{x}; \mathbf{w})$

Find the gradient descent update rule for \mathbf{w} to minimize the loss

$$L = \frac{1}{2} \sum_{i=1}^N (f(\mathbf{x}_i; \mathbf{w}) - y_i)^2. \quad (4)$$

Here, we use N number of $\mathbf{x}_i \in \mathbb{R}^D$ and $y_i \in \{-1, 1\}$ for $i \in \{1, \dots, N\}$ which are given in advance.

$$f(x; w) = \frac{1 - \exp(-w^T x)}{1 + \exp(-w^T x)}$$

$$L = \frac{1}{2} \sum_{i=1}^N (f(x_i; w) - y_i)^2$$

$$\frac{dL}{dw} = \sum_{i=1}^N (f(x_i; w) - y_i) \cdot \frac{df(x_i; w)}{dw} \quad \text{--- ①}$$

$$\frac{df(x_i; w)}{dw} = \frac{d}{dw} \left(\frac{1 - \exp(-w^T x)}{1 + \exp(-w^T x)} \right)$$

$$= \frac{\{ -\exp(-w^T x) \cdot (-x) \} \{ 1 + \exp(-w^T x) \} - \{ 1 - \exp(-w^T x) \} \{ \exp(-w^T x) \cdot (-x) \}}{\{ 1 + \exp(-w^T x) \}^2}$$

$$\downarrow \left\{ \frac{f(x)}{g(x)} \right\}' = \frac{f'(x)g(x) - f(x)g'(x)}{\{g(x)\}^2}$$

$$= \frac{x \cdot \exp(-w^T x) + x \cdot \{\exp(-w^T x)\}^2 + x \cdot \exp(-w^T x) - x \{\exp(-w^T x)\}^2}{\{1 + \exp(-w^T x)\}^2}$$

$$= \frac{2x \cdot \exp(-w^T x)}{\{1 + \exp(-w^T x)\}^2} \quad \text{--- ②}$$

② → ① 을 하면,

$$\frac{dL}{dw} = \left[\sum_{i=1}^N \left\{ \frac{1 - \exp(-w^T x_i)}{1 + \exp(-w^T x_i)} - y_i \right\} \right] \left[\frac{2x \cdot \exp(-w^T x)}{\{1 + \exp(-w^T x)\}^2} \right]$$

$$= \sum_{i=1}^N (f - y_i) \{ 2x_i \cdot \sigma(1 - \sigma) \}$$

$$\therefore w_{t+1} = w_t - \eta \cdot \frac{dL}{dw}$$

$$= w_t - \eta \cdot \left[\sum_{i=1}^N (f - y_i) \{ 2x_i \cdot \sigma(1 - \sigma) \} \right]$$

$$\left(\text{단, } \sigma = \frac{1}{1 + \exp(-w^T x)} \right)$$

$$f = \frac{1 - \exp(-w^T x)}{1 + \exp(-w^T x)}$$

Problem 4

For a given data set $\mathcal{D} = \{\mathbf{x}_i, y_i\}_{i=1}^N$, derive the closed-form solution for \mathbf{w} that maximizes the following probability P .

$$P = \prod_{i=1}^N \frac{1}{\sqrt{2\pi\sigma^2}} \exp\left(-\frac{1}{2\sigma^2}(y_i - f(\mathbf{x}_i; \mathbf{w}))^2\right), \quad (5)$$

with $f(\mathbf{x}; \mathbf{w}) = \mathbf{w}^\top \mathbf{x}$ for $\mathbf{x}, \mathbf{w} \in \mathbb{R}^D$.

exponential 부분이 포함된 항이 \prod 를 통해 곱의 합으로 전개되는 형태의 식이다. 이에 따라, 지수를 좀 더 간단하게 다루기 위해 양변에 로그를 취하는 방식을 고려하였다.

이때 로그함수는 monotonically increasing의 특성을 지니고 있다. 따라서, 임의의 함수 $f(x)$ 에 대하여

$$\operatorname{argmin}_x f(x) = \operatorname{argmin}_x \ln[f(x)]$$

가 성립한다.

$$\ln P = - \sum_{i=1}^N \frac{1}{\sqrt{2\pi\sigma^2}} \left(\frac{1}{2\sigma^2} (y_i - f(\mathbf{x}_i; \mathbf{w}))^2 \right)$$

위의 식에서, $\frac{1}{\sqrt{2\pi\sigma^2}}$ 부분을 시그마 앞으로 빼고, y 의 index를 생략하여 y 로 표현한다.

$$\ln P = - \frac{1}{\sqrt{2\pi\sigma^2}} \left(\frac{1}{2\sigma^2} (y - \mathbf{w}^\top \mathbf{x})^2 \right)$$

$$\begin{aligned}
\frac{d \ln P}{d w} &= -\frac{1}{\sqrt{2\pi\sigma^2}} \cdot \frac{1}{2\sigma^2} \cdot \frac{d}{dw} (y - w^T x)^2 \\
&= -\frac{1}{\sqrt{2\pi\sigma^2}} \cdot \frac{1}{2\sigma^2} \cdot 2(y - w^T x) \cdot \left[\frac{d}{dw} (y - w^T x) \right] \\
&= -\frac{1}{\sqrt{2\pi\sigma^2}} \cdot \frac{1}{2\sigma^2} \cdot 2(y - w^T x) \cdot (-x) \\
&= \boxed{\frac{1}{\sqrt{2\pi\sigma^2}} \cdot \frac{1}{2\sigma^2}} \cdot 2(y - w^T x) \cdot x = 0 \\
&\quad \text{constant}
\end{aligned}$$

$$(y - w^T x)x = 0 \rightarrow w^T x^2 = yx$$

$$w^T x x = y x \rightarrow w^T = y x x^{-1} x^{-1} = y I x^{-1} = y x^{-1}$$

$$\therefore w^T = y x^{-1}$$

$$w = (y x^{-1})^T = (x^{-1})^T y^T$$

이때, $x \in \mathbb{R}^D : D \text{ by } 1$ $w \in \mathbb{R}^D : D \text{ by } 1$

$$\therefore w^T x : (1 \text{ by } D) \times (D \text{ by } 1) = \text{scalar}$$

$$\therefore \boxed{w = (x^{-1})^T y^T}$$

Problem 5

The Frobenius dot product $\langle \mathbf{A}, \mathbf{B} \rangle$ is defined as

$$\langle \mathbf{A}, \mathbf{B} \rangle = \text{tr}(\mathbf{A}^\top \mathbf{B}), \quad (6)$$

for $\mathbf{A}, \mathbf{B} \in \mathbb{R}^{N \times N}$ using the scalar function $\text{tr}(\cdot)$ for the trace of a matrix,

$$\text{tr}(\mathbf{M}) = \sum_i^N M_{ii}, \quad \text{for } \mathbf{M} \in \mathbb{R}^{N \times N}. \quad (7)$$

Find the derivative of $\langle \mathbf{A}, \mathbf{B} \rangle$ with respect to \mathbf{A} using the definition of the matrix derivative:

$$\left[\frac{d}{d\mathbf{A}} \langle \mathbf{A}, \mathbf{B} \rangle \right]_{ij} = \frac{\partial}{\partial A_{ij}} \langle \mathbf{A}, \mathbf{B} \rangle. \quad (8)$$

$$\langle \mathbf{A}, \mathbf{B} \rangle = \text{tr}(\mathbf{A}^\top \mathbf{B}) : \text{Frobenius dot product}$$

$$\left[\frac{d}{d\mathbf{A}} \langle \mathbf{A}, \mathbf{B} \rangle \right]_{ij} = \frac{\partial}{\partial A_{ij}} \langle \mathbf{A}, \mathbf{B} \rangle = \frac{\partial}{\partial A_{ij}} \text{tr}(\mathbf{A}^\top \mathbf{B})$$

각각의 부분을 element-wise notation으로 바꾼다.

$$[\mathbf{AB}]_{pq} = \sum_{k=1}^D A_{pk} B_{kq}$$

$$[\mathbf{A}^\top \mathbf{B}]_{pq} = \sum_{k=1}^D A_{kp} B_{kq}$$

$$\text{tr}[\mathbf{A}^\top \mathbf{B}] = \sum_{l=1}^D \sum_{k=1}^D A_{kl} B_{kl}$$

$$\left[\frac{d}{d\mathbf{A}} \langle \mathbf{A}, \mathbf{B} \rangle \right]_{pq} = \frac{\partial}{\partial A_{pq}} \langle \mathbf{A}, \mathbf{B} \rangle$$

$$= \frac{\partial}{\partial A_{pq}} \sum_{l=1}^D \sum_{k=1}^D a_{kl} b_{kl} = b_{pq}$$

$$\therefore \frac{d}{d\mathbf{A}} \langle \mathbf{A}, \mathbf{B} \rangle = \mathbf{B}$$

아래의 상황을 가정

$$A : \begin{bmatrix} a_{11} & a_{12} \\ a_{21} & a_{22} \\ a_{31} & a_{32} \end{bmatrix}$$

$$B : \begin{bmatrix} b_{11} & b_{12} \\ b_{21} & b_{22} \\ b_{31} & b_{32} \end{bmatrix}$$

$$\langle A, B \rangle = \text{tr}(A^T B)$$

$$A^T B = \begin{bmatrix} a_{11} & a_{21} & a_{31} \\ a_{12} & a_{22} & a_{32} \end{bmatrix} \begin{bmatrix} b_{11} & b_{12} \\ b_{21} & b_{22} \\ b_{31} & b_{32} \end{bmatrix}$$

— : $a_{11}b_{11} + a_{21}b_{21} + a_{31}b_{31}$
— : $a_{12}b_{12} + a_{22}b_{22} + a_{32}b_{32}$

이 둘을 더하면 $\text{tr}(A^T B) = \langle A, B \rangle$

따라서, Frobenius inner product는 matrix A, B 의 element-wise multiplication의 합이라고 정성적으로 해석할 수 있다.

만약 특수한 상황으로 $\langle A, A \rangle = \text{tr}(A^T A)$ 를

가정하면 이 값은 Frobenius norm으로

각각의 element를 제곱한 것의 합을 의미한다.

Problem 6

We are given a scalar function $f(\mathbf{X}) = \mathbf{w}^\top \mathbf{X} \mathbf{w}$ for $\mathbf{w} \in \mathbb{R}^D$, $\mathbf{X} \in \mathbb{R}^{D \times D}$. Find the derivative of $f(\mathbf{X})$ with respect to the vector \mathbf{w} . Use the definition of the vector derivative $\left[\frac{df}{d\mathbf{w}} \right]_k = \frac{\partial f}{\partial w_k}$, where w_k is the k -th element of \mathbf{w} .

$$f(x) = \mathbf{w}^\top \mathbf{X} \mathbf{w} \quad \begin{array}{l} \mathbf{w} \in \mathbb{R}^D \\ \mathbf{X} \in \mathbb{R}^{D \times D} \end{array}$$

$$\hookrightarrow (1 \text{ by } D) \times (D \text{ by } D) \times (D \text{ by } 1)$$

$$= 1 \text{ by } 1 \Rightarrow \text{a scalar}$$

element-wise notation을 이용하여 표현하면 아래와 같다.

$$f(\mathbf{X}) = \begin{bmatrix} w_1 & w_2 & \dots & w_D \end{bmatrix} \begin{bmatrix} X_{11} & X_{12} & \dots & X_{1D} \\ X_{21} & X_{22} & \dots & X_{2D} \\ \vdots & \vdots & \ddots & \vdots \\ X_{D1} & X_{D2} & \dots & X_{DD} \end{bmatrix} \begin{bmatrix} w_1 \\ w_2 \\ \vdots \\ w_D \end{bmatrix}$$

$$= \begin{bmatrix} \sum_{i=1}^D w_i X_{i1} & \sum_{i=1}^D w_i X_{i2} & \dots & \sum_{i=1}^D w_i X_{iD} \end{bmatrix} \begin{bmatrix} w_1 \\ w_2 \\ \vdots \\ w_D \end{bmatrix}$$

$$= w_1 \sum_{i=1}^D w_i X_{i1} + w_2 \sum_{i=1}^D w_i X_{i2} + \dots + w_D \sum_{i=1}^D w_i X_{iD}$$

$$= \sum_{i=1}^D w_1 w_i X_{i1} + \sum_{i=1}^D w_2 w_i X_{i2} + \dots + \sum_{i=1}^D w_D w_i X_{iD}$$

$$= \sum_{j=1}^D \sum_{i=1}^D w_j X_{ij} w_i$$

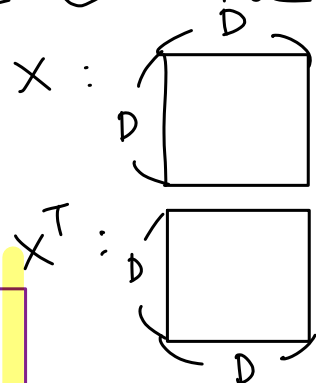
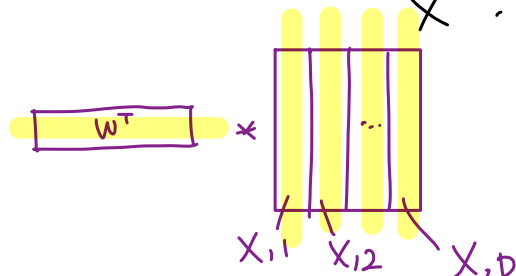
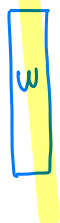
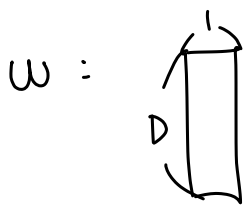
$$\left[\frac{df}{dw} \right]_k = \frac{\partial f}{\partial w_k} = \delta_k^i \sum_{j=1}^D w_j x_{ij} + \delta_k^j \sum_{i=1}^D w_i x_{ij}$$

(라블 시, 해당 element 외에는 모두 0이 됨)
 (δ_b^a 는 $a=b$ 인 항만 남게 함)
 ($\delta_b^a = 1$ if $a=b$)

$$\frac{df}{dw} = \begin{bmatrix} \sum_{j=1}^D x_{1j} w_j + \sum_{i=1}^D w_i x_{i1} \\ \sum_{j=1}^D x_{2j} w_j + \sum_{i=1}^D w_i x_{i2} \\ \vdots \\ \sum_{j=1}^D x_{Dj} w_j + \sum_{i=1}^D w_i x_{iD} \end{bmatrix} = \begin{bmatrix} \sum_{j=1}^D x_{1j} w_j \\ \sum_{j=1}^D x_{2j} w_j \\ \vdots \\ \sum_{j=1}^D x_{Dj} w_j \end{bmatrix} + \begin{bmatrix} \sum_{i=1}^D w_i x_{i1} \\ \sum_{i=1}^D w_i x_{i2} \\ \vdots \\ \sum_{i=1}^D w_i x_{iD} \end{bmatrix}$$

$$= \textcircled{1} \begin{bmatrix} w^T x_1 \\ w^T x_2 \\ \vdots \\ w^T x_D \end{bmatrix} + \textcircled{2} \begin{bmatrix} w x_{,1} \\ w x_{,2} \\ \vdots \\ w x_{,D} \end{bmatrix} = Xw + w^T X$$

↑ 위의 shape, element가 만들어지는 과정을 시각화



Problem 7

For a dataset $\mathcal{D} = \{(\mathbf{x}_i, y_i)\}_{i=1}^N$ with $\mathbf{x}_i \in \mathbb{R}^D$ and $y_i \in \{0, 1\}$, Suppose that we have a two-layer neural network with residual connections shown in Figure 2. Each component is given as follows:

$$L(W) = \frac{1}{2} \sum_{i=1}^N (g(\mathbf{x}_i) - y_i)^2 \quad (9)$$

$$g(\mathbf{x}) = \sigma \left(\sum_{d=1}^D w_{2,d} \cdot h_d(\mathbf{x}) \right) \quad (10)$$

$$h_i(\mathbf{x}) = x_i + \sigma \left(\sum_{m=1}^M w_{1,i,m} \cdot z_m(\mathbf{x}) \right) \quad (11)$$

$$z_i(\mathbf{x}) = \sigma \left(\sum_{d=1}^D w_{0,i,d} \cdot x_d \right) \quad (12)$$

$$\sigma(x) = \frac{1}{1 + \exp(-x)}. \quad (13)$$

Here, the residual connection from each x_i to the h_i node in the second layer is represented in Eq. (11). The weight $w_{i,j,k}$ connects the j -th node of i -th layer to the k -th node of $(i+1)$ -th layer.

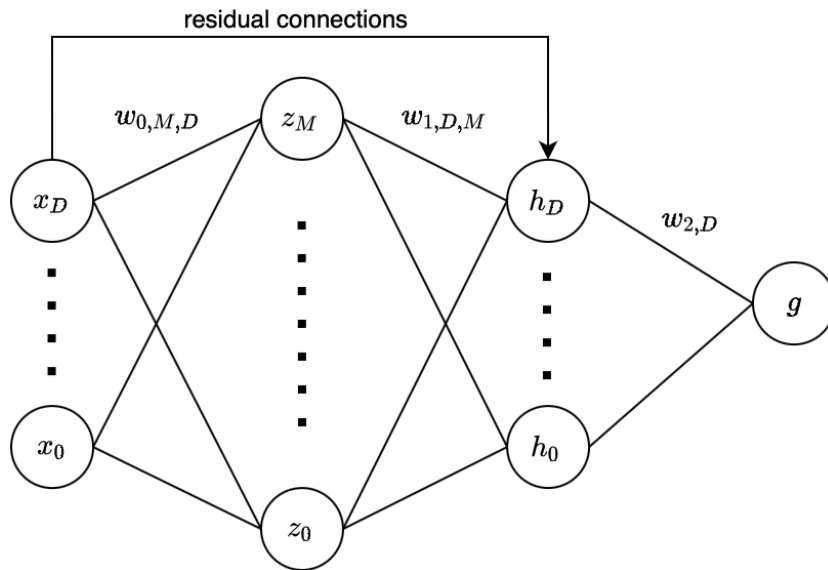


Figure 2: A two-layer network

- (a) Calculate derivative of L with respect to h_j .
- (b) Calculate derivative of L with respect to $w_{1,d,m}$.

먼저 각각을 matrix format으로 바꾸고

(1), (2)를 풀 것이다.

$$L(w) = \frac{1}{2} \sum_{i=1}^N (g(x_i) - y_i)^2$$

$$y_i \in \{0, 1\}$$

$$g(x_i) \in \mathbb{R} \rightarrow \text{scalar}$$

$$g(x) = \sigma \left(\sum_{d=1}^D w_{2,d} \cdot h_d(x) \right)$$

$$w_2 \in \mathbb{R}^D : \begin{pmatrix} w_{2,1} \\ \vdots \\ w_{2,D} \end{pmatrix}$$

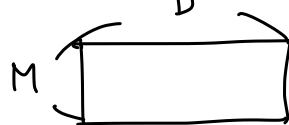
$$h \in \mathbb{R}^D : \begin{pmatrix} h_1(x) \\ h_2(x) \\ \vdots \\ h_D(x) \end{pmatrix}$$

$$\therefore g(x) = \sigma(w_2^T h(x))$$

$$h_i(x) = x_i + \sigma \left(\sum_{m=1}^M w_{1,i,m} \cdot z_m(x) \right)$$

$$x \in \mathbb{R}^D$$

$$w_1 \in \mathbb{R}^{M \times D}$$



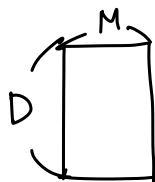
$$z(x) \in \mathbb{R}^M : \begin{pmatrix} z_1(x) \\ z_2(x) \\ \vdots \\ z_M(x) \end{pmatrix}$$

$$\therefore h(x) = x + \sigma(w_1^T z)$$

$$z_i(x) = \sigma \left(\sum_{d=1}^D w_{0,i,d} \cdot x_d \right)$$

$$x \in \mathbb{R}^D$$

$$w_0 \in \mathbb{R}^{D \times M}$$



$$\therefore z(x) = \sigma(w_0^T x)$$

$$\begin{aligned}
 (1) \left[\frac{\partial L}{\partial h} \right]_j &= \frac{\partial L}{\partial g} \cdot \left[\frac{\partial g}{\partial h} \right]_j \\
 &= \left[\frac{1}{2} \cdot 2(g(x) - y) \right] \cdot \sigma(w_2^T h(x)) \cdot (1 - \sigma(w_2^T h(x))) \cdot \underbrace{w_{2,j}}_{\text{scalar}} \\
 &= (g(x) - y_j) \cdot \sigma(w_{2,j} h(x_j)) \cdot (1 - \sigma(w_{2,j} h(x_j))) \cdot w_{2,j}
 \end{aligned}$$

$$(2) \left[\frac{\partial L}{\partial w_1} \right]_{d,m} = \left[\frac{\partial L}{\partial g} \cdot \frac{\partial g}{\partial h} \right] \cdot \left[\frac{\partial h}{\partial w_1} \right]_{d,m}$$

↳ Delta rule을 고쳐.

$$\begin{aligned}
 &= (g(x_d) - y_d) \cdot \sigma(w_{2,d} \cdot h(x_d)) \cdot (1 - \sigma(w_{2,d} \cdot h(x_d))) \cdot w_{2,d} \\
 &\quad \cdot \sigma(w_{1,d,m} \cdot z_d) \cdot (1 - \sigma(w_{1,d,m} \cdot z_d)) \cdot z_d
 \end{aligned}$$

하위 layer에서 derivative를 구할 때에는

상위 layer에서 구한 derivative 값을 이용하여

recursive한 형태의 도출이 가능하다는 것이

delta rule의 핵심이다.