

컴퓨터비전특론 3rd week summary

[Lec 04. Loss]

Analytic Gradient

- Hinge loss

1) Binary hinge loss

$$L = \max(0, 1 - y \cdot s) \quad \left(s = w^T x + b \right)$$

$$= \begin{cases} 1 - y \cdot s & \text{if } 1 - y \cdot s > 0 \\ 0 & \text{otherwise} \end{cases} \quad \left(y = \pm 1 \right)$$

$$\Rightarrow \frac{\partial L}{\partial s} = \begin{cases} -y & \text{if } 1 - y \cdot s > 0 \\ 0 & \text{otherwise} \end{cases}$$

2) hinge loss (multi-class)

$$L = \sum_{j=1, j \neq y}^n \max(0, s_j - s_y + 1)$$

$$\Rightarrow \frac{\partial L}{\partial s_y} = - \sum_{j=1, j \neq y}^n \mathbb{1}(s_j - s_y + 1 > 0) \quad \text{for } j=y$$

$$\frac{\partial L}{\partial s_j} = \mathbb{1}(s_j - s_y + 1 > 0) \quad \text{for } j \neq y$$

$$s = w x + b, \quad w = \begin{pmatrix} w_1^T \\ \vdots \\ w_n^T \end{pmatrix}, \quad s = \begin{pmatrix} s_1 \\ \vdots \\ s_n \end{pmatrix}$$

x : image, y : class label ($1 \leq y \leq n$)

$$\mathbb{1}(F) = \begin{cases} 1 & \text{if } F \text{ is true} \\ 0 & \text{otherwise} \end{cases}$$

- Log likelihood loss

$$L = -\log p_y \quad \text{where } z_y = 1$$

$p = \begin{pmatrix} p_1 \\ \vdots \\ p_n \end{pmatrix}$: Prob. for i^{th} image

y : class label for i^{th} image, $1 \leq y \leq n$

z : class probability = $(z_1, \dots, z_n)^T$, $z_y = 1$, $z_k \neq y = 0$

$$\Rightarrow \frac{\partial L}{\partial p} = - \begin{bmatrix} 0 \\ \vdots \\ 1/p_y \\ \vdots \\ 0 \end{bmatrix}$$

→ 정답 label만 0이 아님

- Cross-Entropy loss

$$L = - \sum_{j=1}^n (z_j \log p_j + (1 - z_j) \log (1 - p_j))$$

$p = \begin{pmatrix} p_1 \\ \vdots \\ p_n \end{pmatrix}$: Prob. for i^{th} image

y : class label for i^{th} image, $1 \leq y \leq n$

z : class probability = $(z_1, \dots, z_n)^T$, $z_y = 1$, $z_k \neq y = 0$

$$\Rightarrow \frac{\partial L}{\partial p} = \begin{bmatrix} 1/(1-p_1) \\ \vdots \\ -1/p_y \\ \vdots \\ 1/(1-p_n) \end{bmatrix}$$

→ 정답 label

- Regression loss

$$y = \begin{pmatrix} y_1 \\ \vdots \\ y_n \end{pmatrix}, \quad s = \begin{pmatrix} s_1 \\ \vdots \\ s_n \end{pmatrix} \quad \text{일 때}$$

$$L = (y - s)^2 \Rightarrow \frac{\partial L}{\partial s} = -2(y - s)$$

$$L = |y - s| \Rightarrow \frac{\partial L}{\partial s} = \begin{cases} -1 & \text{for } y - s > 0 \\ 1 & \text{for } y - s < 0 \end{cases}$$

* 0일 때는 미분 불가능!

Image features vs ConvNet

- Image features : SIFT, HoG, BoW 등의 feature transformation 방법으로 이미지를 vector로 표현

Handcrafted feature : 변형 사이즈 지정, edge 분할 등

인간이 정해진 규칙에 따라 feature extraction

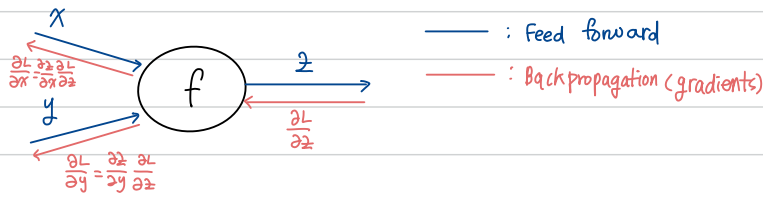
→ Feature extraction & classifier training (SVM, RF 등)

- ConvNets : Conv를 통한 feature representation 과정과 softmax + log likelihood 과정이 모두 training에 포함됨

⇒ 답러닝이 더 좋은 성능을 보인다

[Lec 05. Backpropagation and Neural Networks]

Backpropagation



$\star \star$

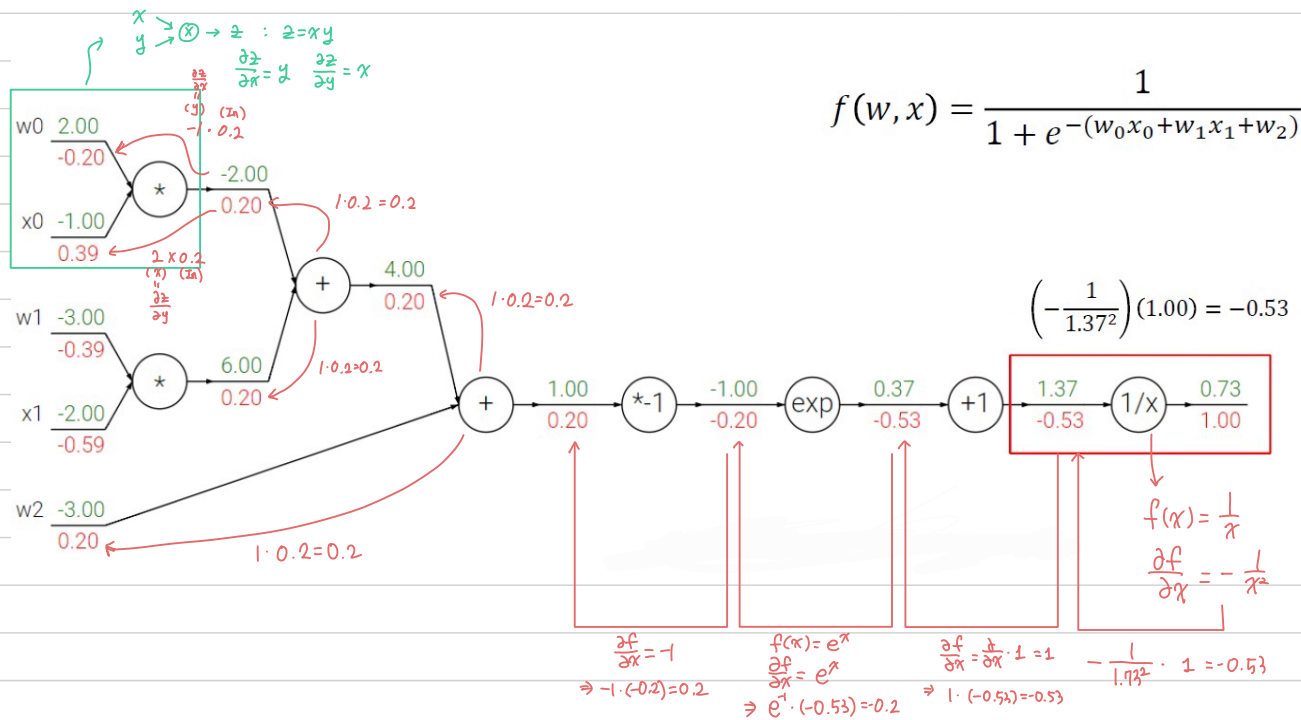
Chain rule

$p = \begin{pmatrix} x \\ y \end{pmatrix}, \frac{\partial L}{\partial z}$: 이전 단계의 outgoing gradient일 때

$\frac{\partial L}{\partial p} = \frac{\partial L}{\partial p} \cdot \frac{\partial z}{\partial p}$

Outgoing gradient Jacobian matrix Incoming gradient

- scalar example

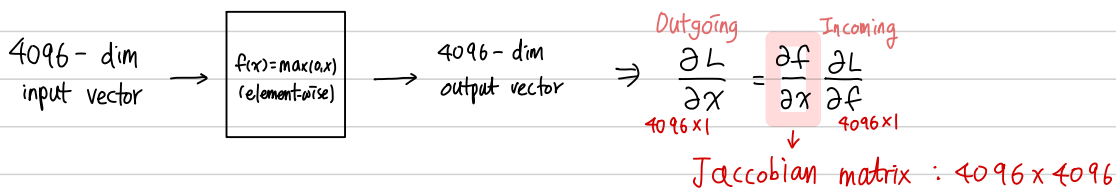


* Derivative of sigmoid function

Sigmoid function $\sigma(x) = \frac{1}{1+e^{-x}}$

$\rightarrow \frac{\partial \sigma(x)}{\partial x} = \frac{1+e^{-x}-1}{(1+e^{-x})^2} = \frac{e^{-x}}{1+e^{-x}} \cdot \frac{1}{1+e^{-x}} = (1-\sigma(x))\sigma(x)$ \rightarrow Sigmoid function의 backpropagation을 한번에 계산가능

- vectorized operations



* Minibatch = 100 일 경우

: 4096 짜리 vector가 100번 들어옴

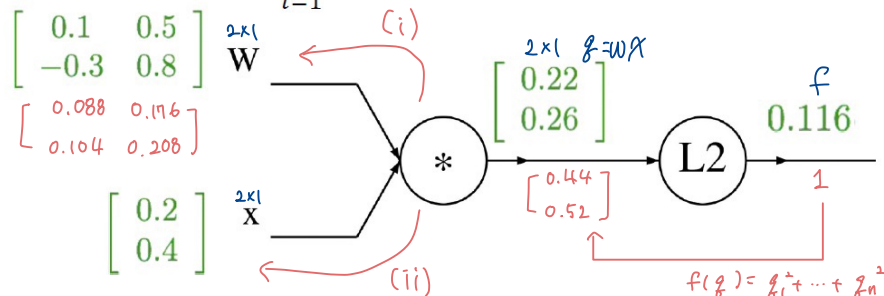
\Rightarrow 409600 x 409600 matrix

(= 4096 x 4096 x 100 3D Tensor)

- vectorized example

$$f(x, W) = |Wx|^2 = \sum_{i=1}^n (Wx)_i^2 \quad x \in \mathbb{R}^d, W \in \mathbb{R}^{n \times d}$$

$q_i = (Wx)_i$: i^{th} element
 q : $n \times 1$ vector



$$q = Wx = \begin{pmatrix} W_{1,1}x_1 + \dots + W_{1,d}x_d \\ \vdots \\ W_{n,1}x_1 + \dots + W_{n,d}x_d \end{pmatrix} = \begin{pmatrix} q_1 \\ \vdots \\ q_n \end{pmatrix}$$

$$f(q) = |q|^2 = q_1^2 + \dots + q_n^2$$

$$(ii) \frac{\partial q_k}{\partial x_j} = W_{k,j}$$

$$\frac{\partial f}{\partial x_j} = \sum_k \frac{\partial f}{\partial q_k} \frac{\partial q_k}{\partial x_j} = \sum_k 2q_k \cdot W_{k,j}$$

$$\rightarrow \frac{\partial f}{\partial x} = 2W^T q = 2W^T W x$$

$$\frac{\partial f}{\partial x}, \frac{\partial f}{\partial W} = ?$$

$$(i) \frac{\partial f}{\partial w_{ij}} = \sum_k \frac{\partial f}{\partial q_k} \frac{\partial q_k}{\partial w_{ij}} \quad q_k = W_{k,1}x_1 + \dots + W_{k,d}x_d$$

$$\rightarrow \frac{\partial q_k}{\partial w_{ij}} = 1 \quad \text{if } k=i, j=j \quad \text{otherwise } 0$$

$$\rightarrow \frac{\partial f}{\partial w_{ij}} = \sum_k \frac{\partial f}{\partial q_k} \frac{\partial q_k}{\partial w_{ij}} = \sum_k 2q_k (1_{k=i} x_j) = 2q_i x_j$$

$$\frac{\partial f}{\partial W} = 2q x^T = 2W x x^T$$

$n \times d$ $1 \times d$ $n \times d$ $d \times d$

$$\begin{aligned} \frac{\partial f}{\partial w_i} &= \frac{\partial f}{\partial w_i} \frac{\partial f}{\partial q} = 2(W^T x)_i x \rightarrow \frac{\partial f}{\partial w} = 2W x x^T \\ \frac{\partial f}{\partial x} &= \frac{\partial f}{\partial x} \frac{\partial f}{\partial q} = 2W^T q \rightarrow \frac{\partial f}{\partial x} = 2W^T W x \end{aligned}$$

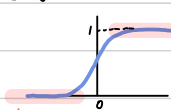
Neural Networks

- Linear vs NN

Linear: layer가 여러개 있어도 한번 쓴 것과 동일한 결과
NN: Nonlinear operation을 사용 → Layer별로 다른 weight가 채워짐

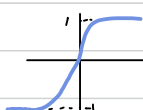
- Activation functions

[Sigmoid] $\sigma(x) = \frac{1}{1+e^{-x}}$

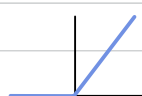


Saturation (Gradient Vanishing 발생 가능)

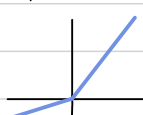
[tanh] $\tanh(x) = \frac{e^x - e^{-x}}{e^x + e^{-x}}$



[ReLU] $\max(0, x)$

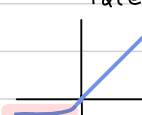


[Leaky ReLU] $\max(0, \alpha x)$

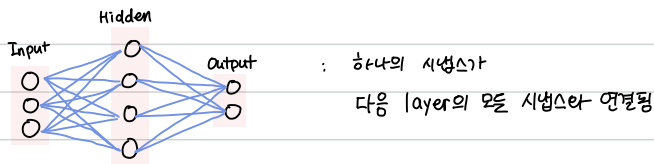


[Maxout] $\max(w_1^T x + b_1, w_2^T x + b_2)$

[ELU] $\begin{cases} x, & x \geq 0 \\ \alpha(e^x - 1), & x < 0 \end{cases}$ * α : hyperparameter

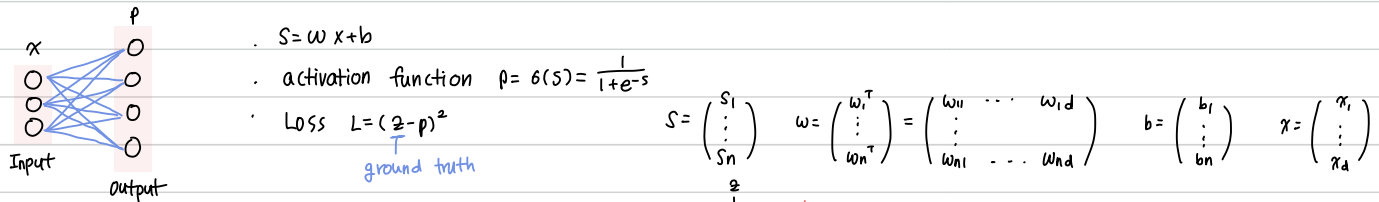


- Fully connected layers (= Multi layer Perceptron)



- Derivative of NN using chain rules

1) 1-layer NN (L2 regression loss)



$$x \xrightarrow{d \times 1} W \xrightarrow{n \times d} S \xrightarrow{n \times 1} \text{sigmoid} \xrightarrow{n \times 1} p \xrightarrow{n \times 1} \text{L2 loss} \Rightarrow \frac{\partial L}{\partial p}, \frac{\partial L}{\partial s}, \frac{\partial L}{\partial w}, \frac{\partial L}{\partial b} \text{ 을 계산해야 함}$$

$$\frac{\partial L}{\partial p} = (z - p)^T \cdot \frac{\partial}{\partial p} = -2(z - p)$$

$$\frac{\partial L}{\partial s} = \frac{\partial p}{\partial s} \frac{\partial L}{\partial p} = \text{diag}((1 - \sigma(s_j)) \sigma(s_j)) \cdot \frac{\partial L}{\partial p} = -2 \begin{bmatrix} (1 - \sigma(s_1)) \sigma(s_1) (z_1 - p_1) \\ \vdots \\ (1 - \sigma(s_n)) \sigma(s_n) (z_n - p_n) \end{bmatrix} = (1 - \sigma(s)) \otimes \sigma(s) \otimes \frac{\partial L}{\partial p}$$

Element-wise multiplication

$$\frac{\partial L}{\partial w_j} = \frac{\partial s}{\partial w_j} \frac{\partial L}{\partial s} = x_j \frac{\partial L}{\partial s} = [0 \ 0 \ \dots \ x_j \ \dots \ 0] \frac{\partial L}{\partial s} = \left(\frac{\partial L}{\partial s} \right)_j x$$

jth column

$$\Rightarrow \frac{\partial L}{\partial w} = \frac{\partial L}{\partial s} x^T$$

$$\frac{\partial L}{\partial b} = \frac{\partial s}{\partial b} \frac{\partial L}{\partial s} = \frac{\partial L}{\partial s}$$

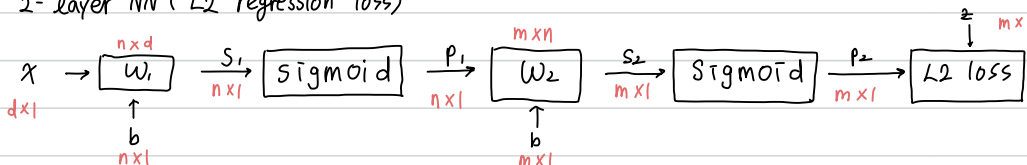
identity

$$* S = Wx + b$$

$$\Rightarrow \frac{\partial L}{\partial w} = \frac{\partial L}{\partial s} \cdot x^T \quad \frac{\partial L}{\partial b} = \frac{\partial L}{\partial s} \cdot 1 = \frac{\partial L}{\partial s}$$

$$\frac{\partial L}{\partial x} = \frac{\partial s}{\partial x} \frac{\partial L}{\partial s} = W^T \frac{\partial L}{\partial s}$$

2) 2-layer NN (L2 regression loss)



$$\frac{\partial L}{\partial p_2} = -2(z - p_2) \quad \frac{\partial L}{\partial s_2} = \frac{\partial p_2}{\partial s_2} \frac{\partial L}{\partial p_2} = \text{diag}((1 - \sigma(s_{2,j})) \sigma(s_{2,j})) \frac{\partial L}{\partial p_2}$$

$$\frac{\partial L}{\partial w_2} = \frac{\partial L}{\partial s_2} p_1^T \quad \frac{\partial L}{\partial b_2} = \frac{\partial L}{\partial s_2}$$

$$\frac{\partial L}{\partial p_1} = \frac{\partial s_2}{\partial p_1} \frac{\partial L}{\partial s_2} = w_2^T \frac{\partial L}{\partial s_2} \quad \frac{\partial L}{\partial s_1} = \text{diag}((1 - \sigma(s_1)) \sigma(s_1)) \frac{\partial L}{\partial p_1} \rightarrow \begin{bmatrix} \frac{\partial L}{\partial w_1} = \frac{\partial L}{\partial s_1} x^T \\ \frac{\partial L}{\partial b_1} = \frac{\partial L}{\partial s_1} \end{bmatrix}$$

두번째 layer의 input data와 같음!

3) 1-layer NN (Softmax classifier)



$$s = Wx + b, \quad s_j = w_j^T x + b_j$$

$$\text{Activation function } p = \frac{e^{s_j}}{\sum_{i=1}^n e^{s_i}}$$

$$\text{Loss } L = -\log p_y \quad \text{where } z_y = 1, \quad z = (z_1, \dots, z_n)^T, \quad z_y = 1, \quad z_{k \neq y} = 0$$

$$S = \begin{pmatrix} s_1 \\ \vdots \\ s_n \end{pmatrix} \quad W = \begin{pmatrix} w_{11} & \dots & w_{1d} \\ \vdots & & \vdots \\ w_{n1} & \dots & w_{nd} \end{pmatrix} \quad b = \begin{pmatrix} b_1 \\ \vdots \\ b_n \end{pmatrix} \quad x = \begin{pmatrix} x_1 \\ \vdots \\ x_d \end{pmatrix}$$

$$\frac{\partial L}{\partial p} = \begin{bmatrix} 0 \\ \vdots \\ -1/p_y \end{bmatrix} \quad \frac{\partial L}{\partial s} = \frac{\partial p}{\partial s} \frac{\partial L}{\partial p} = -\frac{1}{p_y} \begin{bmatrix} 0 \\ \vdots \\ 1 \end{bmatrix} = -\frac{1}{p_y} \begin{bmatrix} -p_y \\ p_y(1-p_y) \\ \vdots \\ -p_y p_y \end{bmatrix} = \begin{bmatrix} p_1 \\ p_{j-1} \\ \vdots \\ p_n \end{bmatrix} = p - z$$

$$\Rightarrow \frac{\partial L}{\partial w_j} = \frac{\partial s}{\partial w_j} \frac{\partial L}{\partial s} = x_j \frac{\partial L}{\partial s} = [0 \ \dots \ x_j \ \dots \ 0] \frac{\partial L}{\partial s} = \left(\frac{\partial L}{\partial s} \right)_j x$$

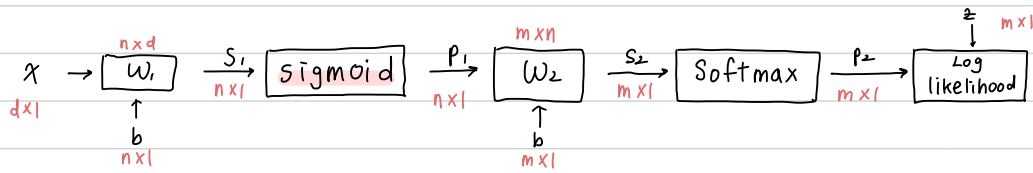
jth column

$$\rightarrow \frac{\partial L}{\partial w} = \left(\frac{\partial L}{\partial w_1} \ \dots \ \frac{\partial L}{\partial w_n} \right)^T = \frac{\partial L}{\partial s} x^T$$

$$\frac{\partial L}{\partial b} = \frac{\partial s}{\partial b} \frac{\partial L}{\partial s} = \frac{\partial L}{\partial s}$$

$$* \frac{\partial L}{\partial x} = \frac{\partial s}{\partial x} \frac{\partial L}{\partial s} = W^T \frac{\partial L}{\partial s}$$

4) 2-layer NN (Softmax classifier)



$$\frac{\partial L}{\partial p_2} = \begin{bmatrix} 0 \\ \vdots \\ -1/p_y \\ 0 \end{bmatrix}$$

$$\frac{\partial L}{\partial s_2} = D \frac{\partial L}{\partial p_2}, \quad D_{ab} = p_a (p_{ab} - p_b), \quad p_{ab} = \int_0^1 \delta_{ab} dw$$

$$\frac{\partial L}{\partial W_2} = \frac{\partial L}{\partial s_2} p_1^T$$

$$\frac{\partial L}{\partial b_2} = \frac{\partial L}{\partial s_2}$$

$$\frac{\partial L}{\partial p_1} = W_2^T \frac{\partial L}{\partial s_2}$$

$$\frac{\partial L}{\partial s_1} = \text{diag}((1 - \sigma(s_{1,j})) \sigma(s_{1,j})) \frac{\partial L}{\partial p_1}$$

$$\frac{\partial L}{\partial W_1} = \frac{\partial L}{\partial s_1} x^T$$

$$\frac{\partial L}{\partial b} = \frac{\partial L}{\partial s_1}$$