



Lecture 3. Model Optimization



Feedbacks,

열정적인 수업 감사드립니다.

코딩수업 바뀐방식이 더 좋은데 조금만더 천천히 진행해주셨으면 좋겠어요

Regularization이 이해가 잘 안되었습니다 순간적으로 집중을 못해서 놓쳤네요

후기작성
upload하겠습니다.

수업 로드가 어제보다 적어 적당했습니다.

1회차보다 훨씬 어려웠지만 훨씬 좋았습니다.

good

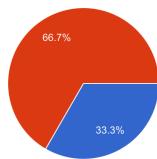
1회차 수업에 비해 강의 난이도가 갑자기 오르긴 했지만, 아직까지는 괜찮은 듯 함.

코딩 실습에서 기본 제공코딩 또한 어떠한 기능을 하는지 설명해주었으면 좋겠습니다.

아마 강의차가 가장 어려울 것 같네요....
수학적 내용은 천천히 이해하셔도 됩니다!

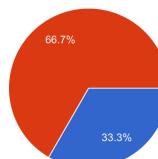
// Loss에 대한 설명이 잘 전달되었나요?

응답 9개



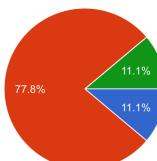
Loss Function에 대한 설명이 잘 전달되었나요?

응답 9개



Regularization에 대한 설명이 잘 전달되었나요?

응답 9개



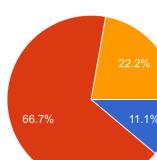
• Github 알기 자료.

• Regularization & Validation 보충설명

• Questions via, Github Issue
(Kakao Talk)

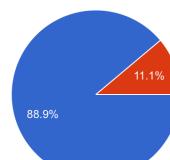
// Linear Classifier 코딩실습의 난이도는 어땠나요?

응답 9개



// 코딩 실습 방식이 바뀐 것에 대해 어떻게 생각하시나요?

응답 9개



• 유익하�, 조금 착향적.

④ Github Repo Star

Review

Questions

I. Train이 잘 되었는지 판단할 수치적 척도 필요

Loss .

: define a **Loss Function** that quantifies our unhappiness with
the scores across the training data

// 2. Parameter를 update하는 algorithm 필요

: come up with a way of efficiently finding the parameters that
minimize the **Loss Function**

//

Review

$$L = \frac{1}{N} \sum_{i=1}^N L_i + \gamma R(\omega).$$

1. Loss
2. **Loss Function**

- Multiclass SVM Loss
- Cross-Entropy Loss



3. **Regularization**

- Why it is needed : to discourage large weight matrix
- Overfitting

$$L = \alpha \| \omega' \|_k \| \omega' \|$$

Review



How is a Model Optimized / Updated?

Loss.



1. Training Set의 Data들을 Linear Classifier에 통과시켜서, 그 결과들을 정답과 비교
- (2. 1에서의 결과를 바탕으로 Parameter의 값들을 Update)
3. 다시 1로

Loss.

Review

Questions

- Train이 잘 되었는지 판단할 수치적 척도 필요

: define a **Loss Function** that quantifies our unhappiness with the scores across the training data

- Parameter를 update하는 algorithm 필요

: come up with a way of efficiently finding the parameters that minimize the **Loss Function**

Optimization and Backpropagation

Loss Function은 Scalar Function

Thus, Loss Function의 input W, b에 대해

Gradient를 구해 주어 W, b에서 빼면,

Loss가 줄어드는 방향으로 학습하지 않을까?

Today's Contents

- 1. Optimization**
- 2. Backpropagation**

Optimization : Idea

(Currently) Incorrect Model → Correct Model 을 위해,

Parameter W와 b를 Update해줘야 함.

How...?

Optimization : Idea

$$\textcircled{1} \quad W=0.5, \quad b=0.3 \rightarrow L=10$$

$$\textcircled{2} \quad W=0.3, \quad b=0.1 \rightarrow L=5$$

IDEA 1. Random Search

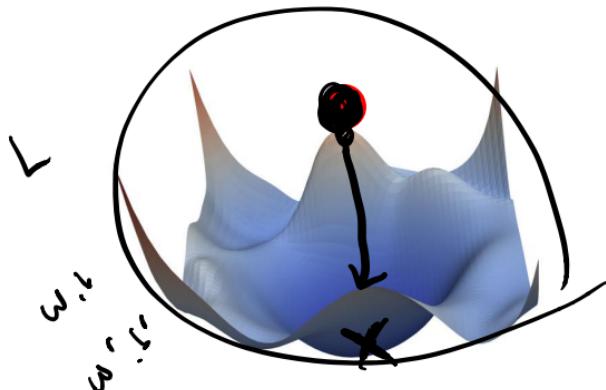
1. Randomly set W and b
2. Select (W, b) that gives the minimum loss in training

Obviously very poor 😞

Optimization : Idea

IDEA 2. Gradient Descent

Loss : W, b, x, y (vector & matrix) $\rightarrow L$ (scalar)



Loss의 Global Minimum으로 가려면 어떻게 해야 할까?

$$L = \frac{1}{N} \sum_{i=1}^N L_i(w_i, b_i, x_i, y_i) + \lambda R(w)$$



Optimization : Idea

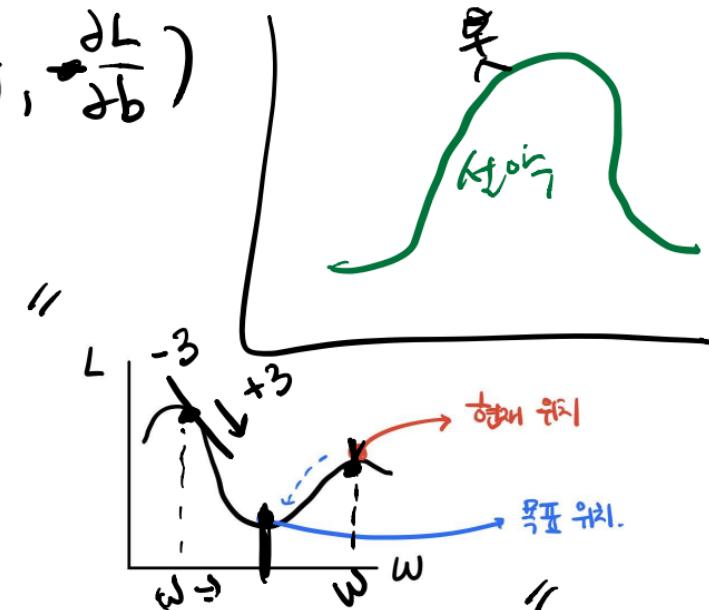
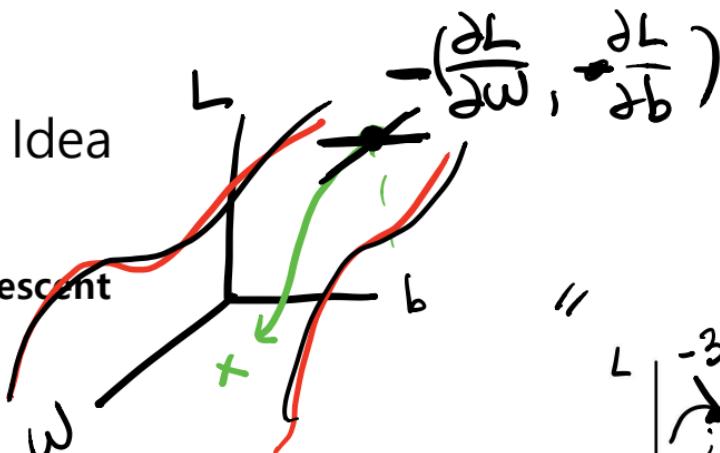
IDEA 2. Gradient Descent

W, b, x, y, L이 모두 Scalar라고 Simplify해 생각

.. ↴

$W' \neq$

$$\underline{W} - \frac{\partial L}{\partial W} \times K$$



- 에서 \bullet 로 가는 방법?
- 에서의 L의 기울기 방향으로 가면 됨.

$$\text{i.e., } \underline{W} - \frac{dL}{dw} \times k$$

(k : some constant).

Optimization : Idea

$$W \begin{matrix} 3 \times 4 \\ b \quad 3 \times 1 \end{matrix}$$

scalar.

$$W = \begin{bmatrix} \text{scalar} \\ | \\ | \end{bmatrix}$$

IDEA 2. Gradient Descent

다시 Vector / Matrix일 때를 생각해 보면,

다차원 공간상의 L 의 Minimum으로 가도록 W, b 를 Update하려면,

We should calculate the GRADIENT ∇L_w

$$\frac{\partial L}{\partial W} = \begin{bmatrix} \frac{\partial L}{\partial w_{00}} & \frac{\partial L}{\partial w_{01}} & \cdots \\ \vdash & \vdash & \vdash \end{bmatrix} \quad 3 \times 4$$

$$b = \begin{bmatrix} \quad \\ \quad \end{bmatrix} \quad b_0, b_1, b_3 \quad 3 \times 1$$

$w_{00}, w_{01}, \dots, w_{23}$ $(2 \times H)$

Optimization : Idea

IDEA 2. Gradient Descent

We can use vector and matrix notation to rewrite things a bit. Define the **gradient** of a scalar-valued function $f: X \subseteq \mathbf{R}^n \rightarrow \mathbf{R}$ to be the *vector*

$$\nabla f(\mathbf{x}) = \left(\frac{\partial f}{\partial x_1}, \frac{\partial f}{\partial x_2}, \dots, \frac{\partial f}{\partial x_n} \right).$$

Consequently,

$$\nabla f(\mathbf{a}) = (f_{x_1}(\mathbf{a}), f_{x_2}(\mathbf{a}), \dots, f_{x_n}(\mathbf{a})).$$

Optimization : Idea

$w_{\text{new}} \leftarrow \text{update } \Sigma \text{ param.}$

IDEA 2. Gradient Descent

//

$$w_{\text{new}} = w_{\text{old}} - lr \cdot \frac{\partial L}{\partial w} \Big|_{w=w_{\text{old}}}$$

$w_{\text{old}} \leftarrow \text{old param.}$

"lr"

$$b_{\text{new}} = b_{\text{old}} - lr \cdot \frac{\partial L}{\partial b} \Big|_{b=b_{\text{old}}} //$$

Hyperparam.

$$\frac{\lambda \cdot R(w)}{\epsilon}$$

L1, L2

Optimization : Calculating Gradient

Gradient Descent

$$\omega_{\text{new}} = \omega_{\text{old}} - lr \cdot \frac{\partial L}{\partial \omega} \quad \omega = \omega_{\text{old}}$$

learning rate
 lr

$$b_{\text{new}} = b_{\text{old}} - lr \cdot \frac{\partial L}{\partial b} \quad b = b_{\text{old}}$$

then, how do we calculate the Gradient?

Optimization : Calculating Gradient

IDEA 1. Numeric

$$\frac{\partial L}{\partial w} = \begin{bmatrix} -2.5 \end{bmatrix}$$

$$\frac{\partial L}{\partial w_{0,0}} = \lim_{h \rightarrow 0} \frac{L(w_{0,0} + h) - L(w_{0,0})}{h}$$

$$f'(x) = \lim_{h \rightarrow 0} \left(\frac{f(x+h) - f(x)}{h} \right)$$

$h = 0.00001$

미분의 정의 활용

$$\frac{\partial L}{\partial w_{0,0}} = \lim_{h \rightarrow 0} \frac{L(w_{0,0} + h) - L(w_{0,0})}{h}$$

$h = 10^{-5}$

예 h 에 10^{-4} 같은 작은 수 대입. $h=0.0001$

ex) $w_{0,0} = 0.34 \rightarrow L(w_{0,0}) = 1.25347$

$w_{0,0} + h = 0.3401 \rightarrow L(w_{0,0} + h) = 1.25322$.

$$\frac{\partial L}{\partial w_{0,0}} = \frac{1.25322 - 1.25347}{0.0001} = -2.5$$

Optimization : Calculating Gradient

IDEA 1. Numeric

poor because,

(1. 근사값만 얻을 수 있음)

(2. Time complexity)

"thus, almost never used in practice"

$$h = 10^{-7} \quad 10^{-5}$$

$$h \leftarrow 0$$

$$w \leftarrow 12\text{개.} \quad b \leftarrow 3\text{개.}$$

y

15 개.

$$\begin{matrix} w_{b+1} + h \\ w_{0,i} + h \end{matrix}$$

Optimization : Calculating Gradient

IDEA 2. Analytic Gradient.

도함수를 직접 구해서 값을 대입해 미분값 구하기

$$L = \underbrace{\left(\frac{1}{N} \sum_{i=1}^N \sum_{j=1}^C \delta_{ij} - \lambda \sum_{k=1}^K \sum_{q=1}^Q w_{kj} \right)}_{\text{Cost Function}}$$

$$\sum_{\substack{j=1 \\ j \neq y_i}}^C \max(0, \delta_j - \delta_{y_i} + \sigma)$$

$$\lim_{h \rightarrow 0} \frac{f(x+h) - f(x)}{h}$$

$$(f'(x) \leftarrow x)$$

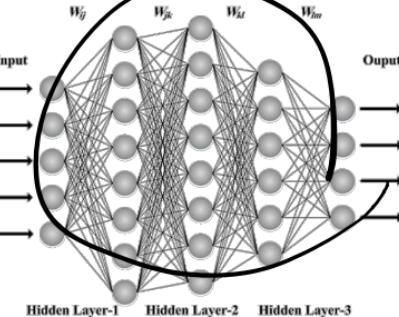
Optimization : Calculating Gradient

IDEA 2. Analytic

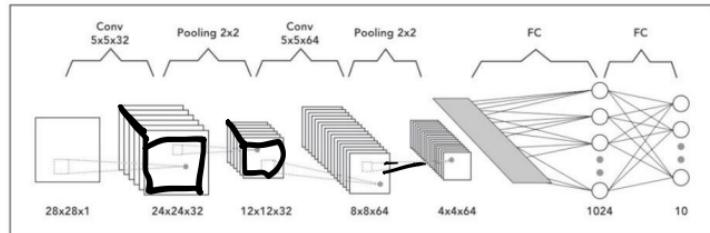
“ 지금은 simple linear classifier만 다루고 있지만,

앞으로 다룰 model은

MLP.

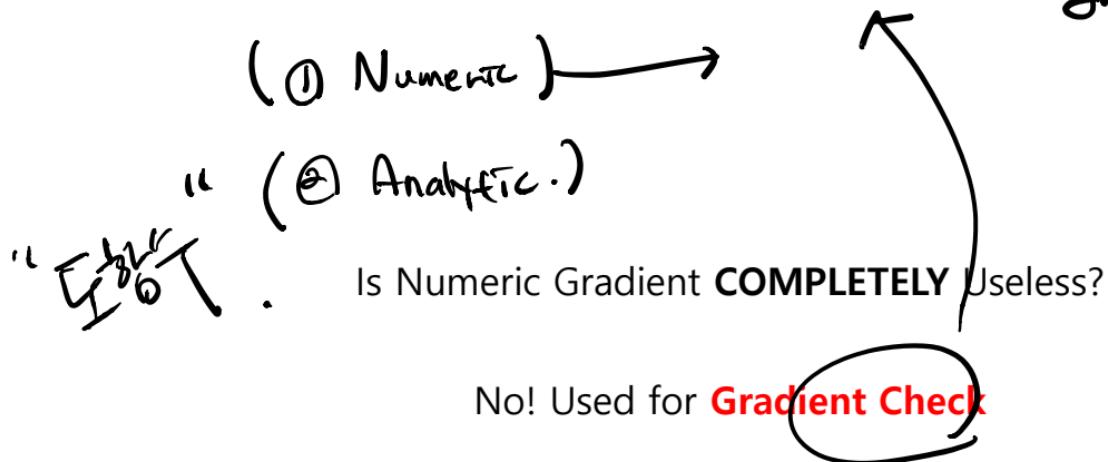


↙
↙
↙ CNN



then, how should we derive the gradient?

Optimization : Calculating Gradient



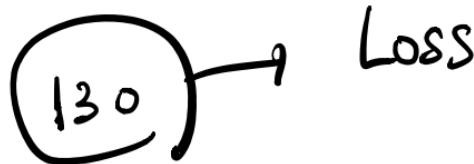
$$\frac{\partial L}{\partial w_j} = g$$

$$g(w')$$

if

$$\lim_{h \rightarrow 0} \frac{f(w'+h) - f(w')}{h}$$

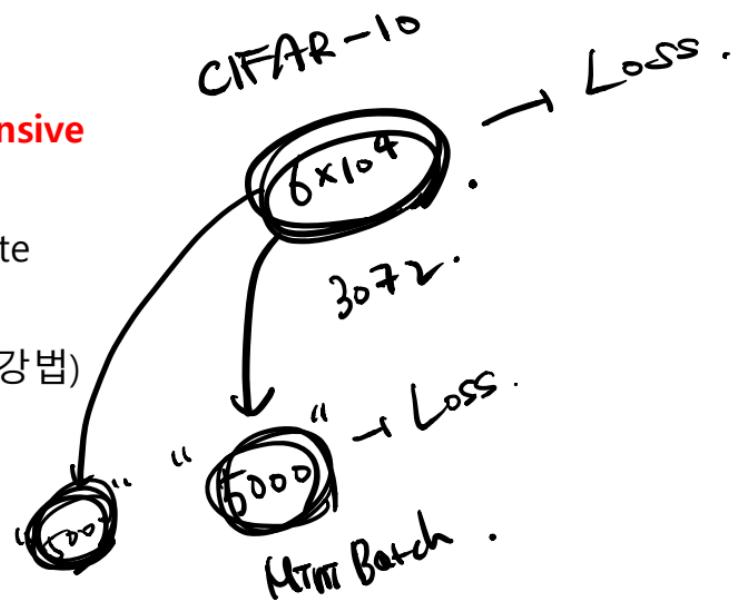
Optimization : Calculating Gradient



but, all dataset에 대해 Loss 계산...? Still Too Expensive

thus, 일부 Dataset(Batch)에 대해 Loss 계산 후 Update

= SGD (Stochastic Gradient Descent, 확률적 경사하강법)
MiniBatch - SGD.



“ Backpropagation : Idea ”

$$L = \frac{1}{N} \sum_{i=1}^N \dots$$

MLP. CNN .

Analytic

We compute the Gradient via, “**Backpropagation**”

(Backpropagation : Idea)

$$g(x) = e^x$$

$$h(x) = \cos x.$$

$$\underline{f = g(h(x))}$$



$$\underline{f(x) = e^{\cos x}}$$

$$L = h(f(g(\dots)))$$

1. 복잡한 함수는 간단한 함수/연산들의 합성함수로 나타낼 수 있다.

2. 간단한 함수들에 Chain Rule을 적용하여, 복잡한 함수의 미분계수를 구한다.

$$\frac{df}{dx} = b'(x) \times g'(h(x)) = \frac{df}{dh} \times \frac{dh}{dx} = -\sin x \times e^x$$

$$\frac{\partial L}{\partial w} = \frac{\partial L}{\partial h} \times \frac{\partial h}{\partial w}$$

Backpropagation : Computational Graph

복잡한 함수를 간단한 함수/연산들의 합성으로 나타낼 때,

“**Computational Graph**”를 활용 (Node : Operator, Leaf : Operand)



$$(1, 2, 3) - \mathbf{lr} \times (3, 3, 3)$$

$$\nabla f = (3, 3, 3)$$

"Backpropagation": Computational Graph

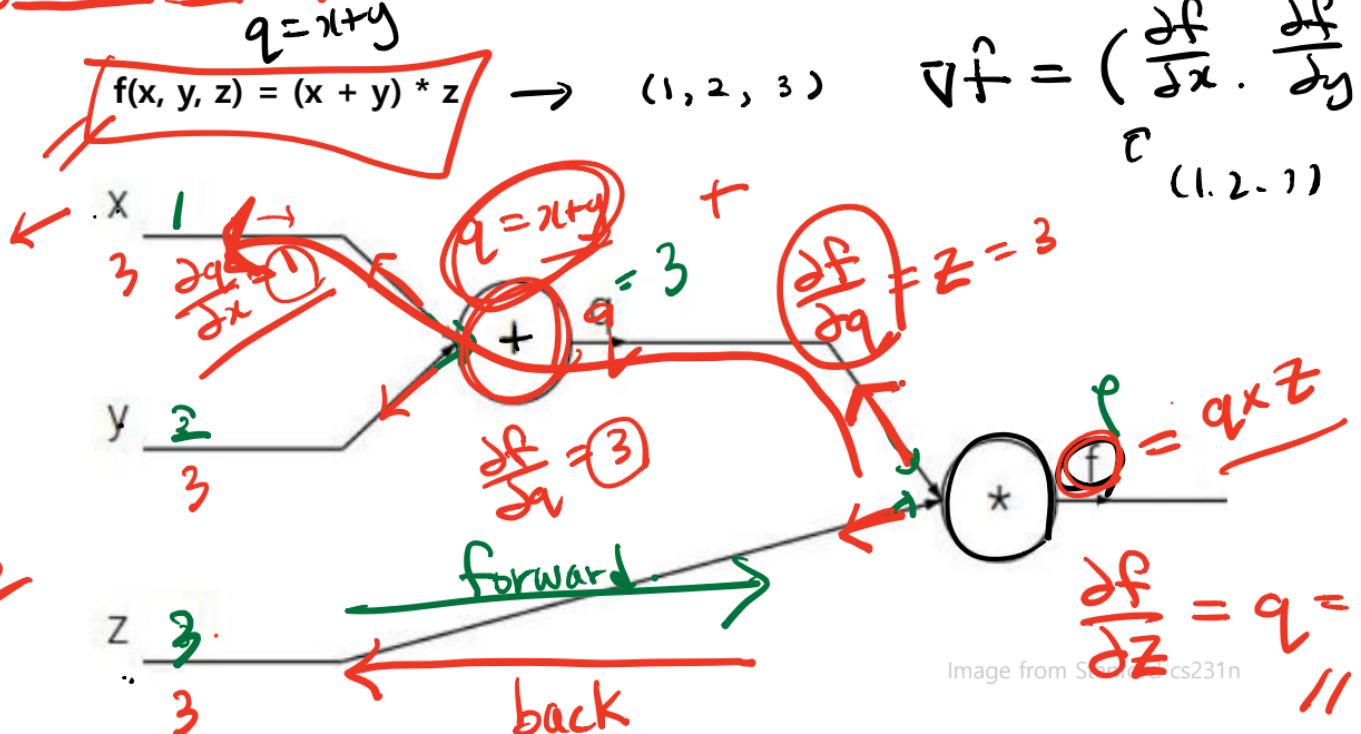


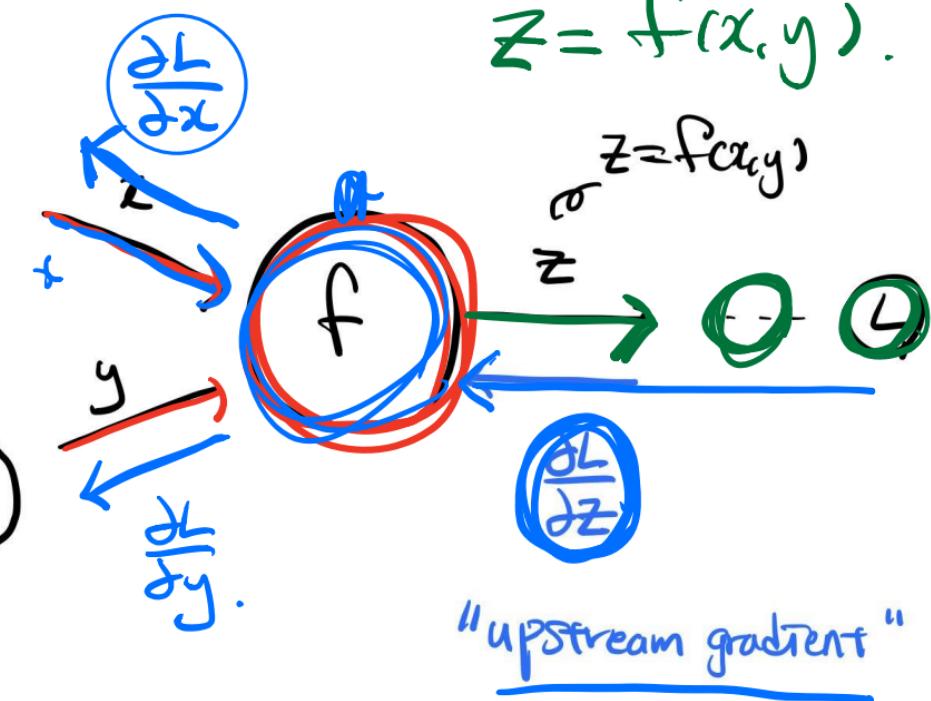
Image from Stanford CS231n

Backpropagation : Computational Graph

Interpreting each node:

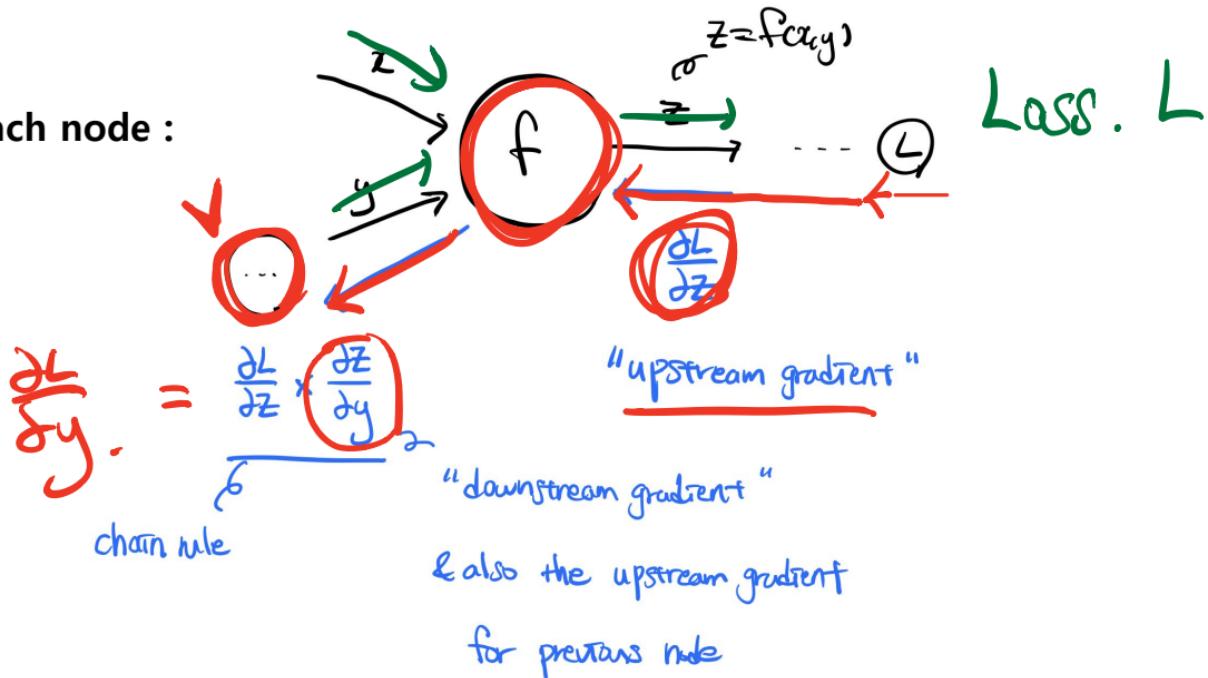
$$\frac{\partial L}{\partial x} = \frac{\partial L}{\partial z} \times \frac{\partial z}{\partial x}$$

UG. Local Gradient.



Backpropagation : Computational Graph

Interpreting each node :



Backpropagation : Computational Graph

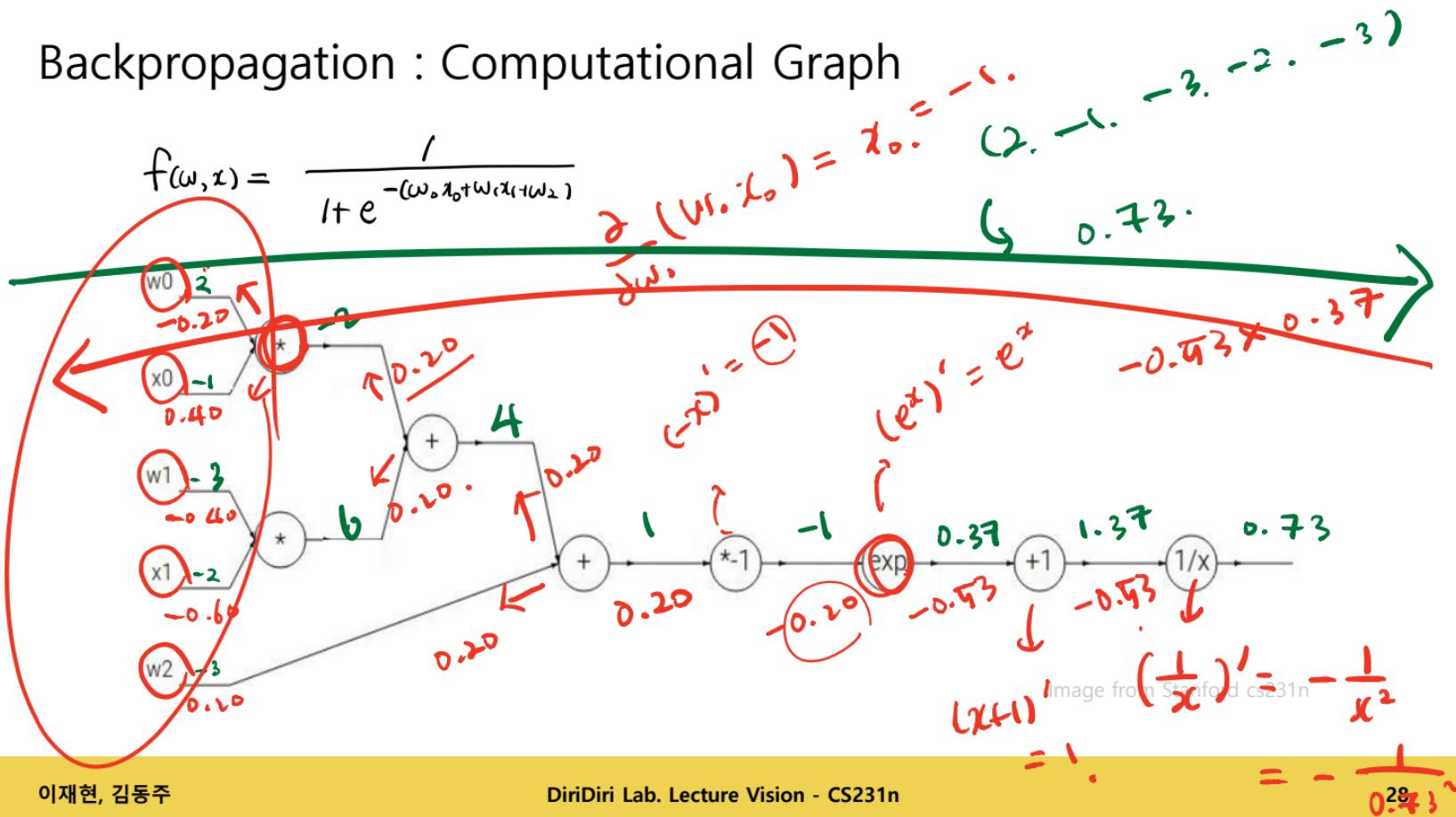
$$(x+y) \times z$$

$$\frac{\partial f}{\partial x} = z$$

Too simple...?

Try $f(\omega, x) = \frac{1}{1 + e^{-(\omega_0 x_0 + \omega_1 x_1 + \omega_2 x_2)}}$

Backpropagation : Computational Graph



Backpropagation : Computational Graph

Actually, inputs are not scalars

$$L = \frac{1}{N} \sum_{i=1}^N L_i(f(x_i, w; b), y_i) + \lambda R(w)$$

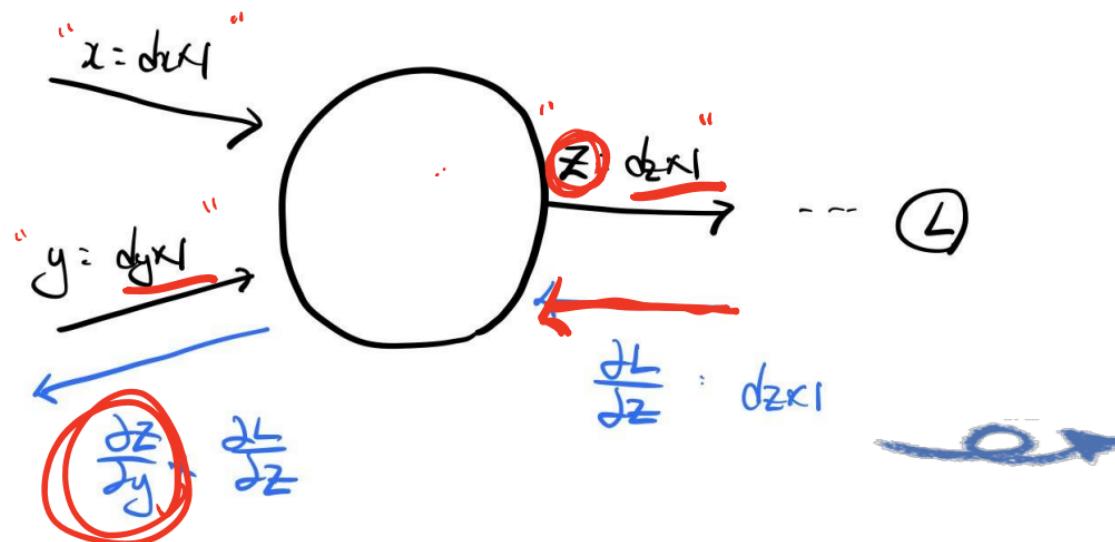
x : vector

W : matrix

b : vector

Backpropagation : Computational Graph

Vector in, Vector out

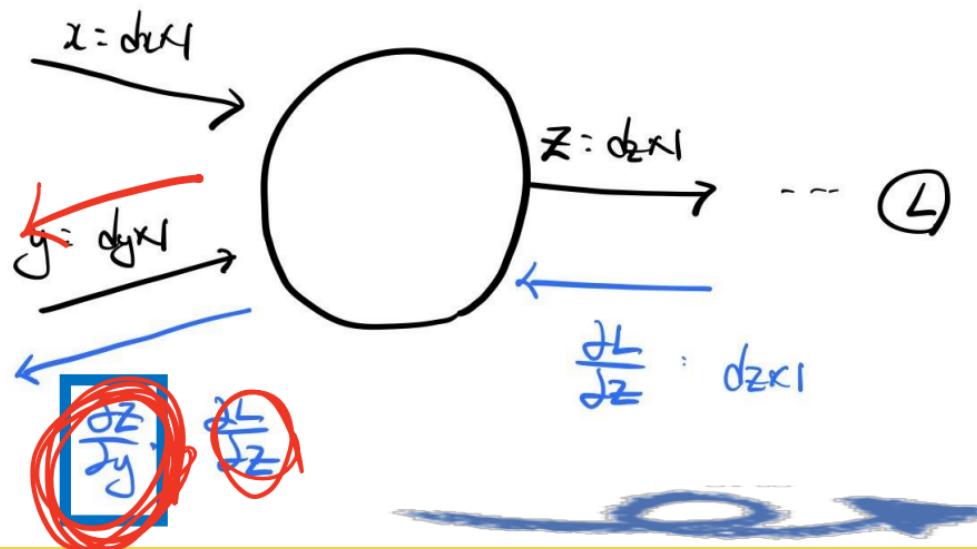


$$\frac{\partial L}{\partial z} = \left[\frac{\partial L}{\partial z_0}, \frac{\partial L}{\partial z_1}, \dots \right]^T$$

scalar L을 Z의 element들로
편미분한 vector

Backpropagation : Computational Graph

Vector in, Vector out



$$\begin{array}{c} \frac{\partial z}{\partial y} \\ \frac{\partial z}{\partial x} \\ \frac{\partial z}{\partial x} \end{array}$$

y

\vdots

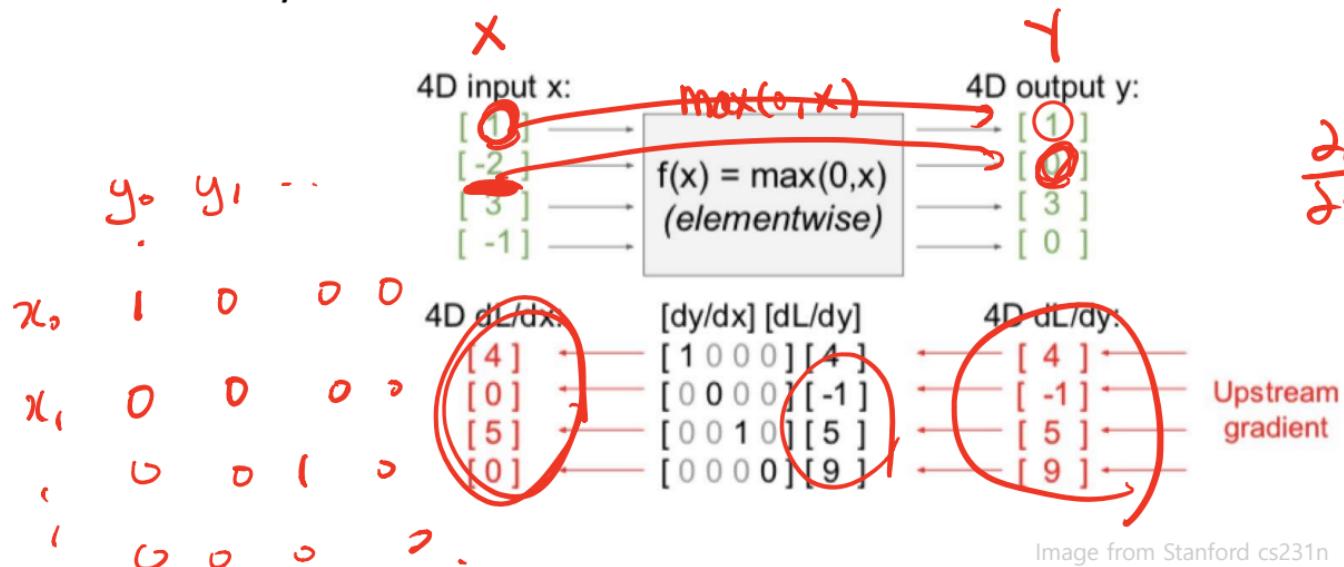
“Jacobian”

: element-wise derivative matrix

$$dy \times dz \left[\begin{array}{c} \frac{\partial z_0}{\partial y_0} \\ \frac{\partial z_0}{\partial y_1} \\ \vdots \\ \frac{\partial z_n}{\partial y_0} \\ \frac{\partial z_n}{\partial y_1} \\ \vdots \end{array} \right]$$

Backpropagation : Computational Graph

Vector in, Vector out

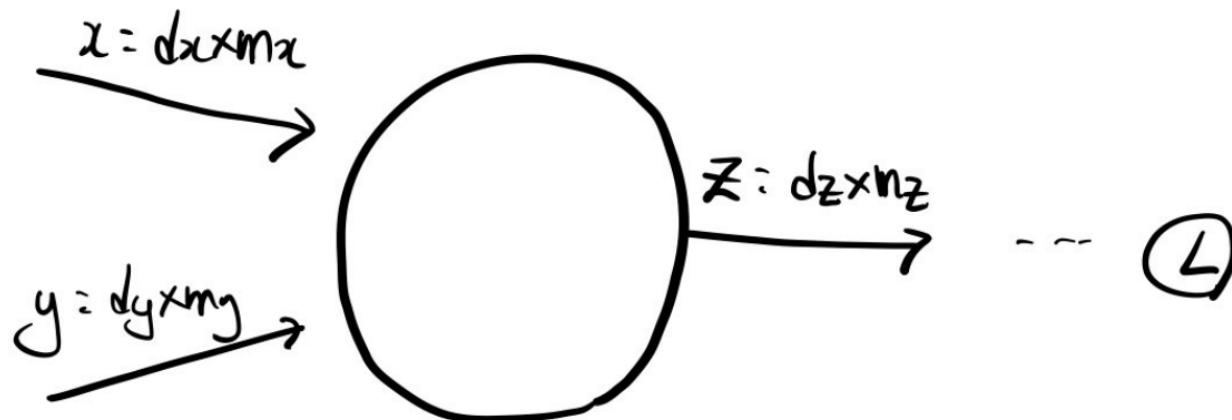


$$\frac{\partial L}{\partial y} = [4, -1, 5, e]^T$$

$$\frac{\partial L}{\partial x} = \underline{\frac{\partial y}{\partial x}} \times \underline{\frac{\partial L}{\partial y}}$$

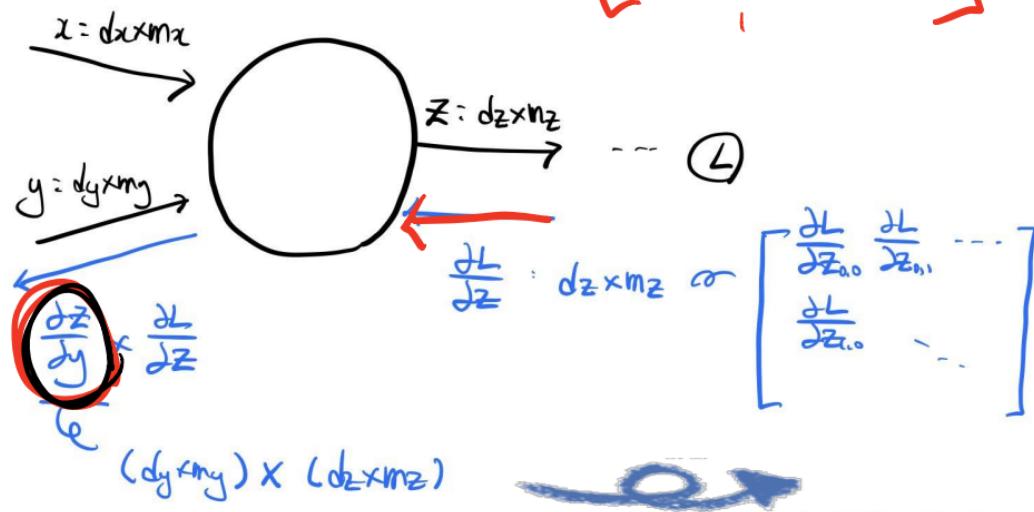
Backpropagation : Computational Graph

Matrix in, Matrix out

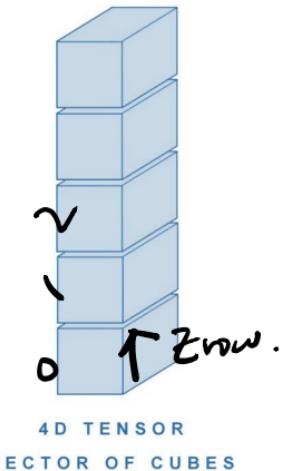


Backpropagation : Computational Graph

Matrix in, Matrix out $\frac{\partial z}{\partial y}$.

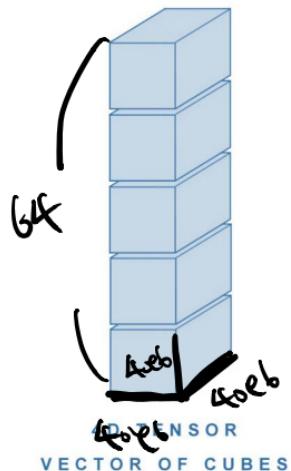


"General Jacobian은 4D Tensor"



Backpropagation : Computational Graph

Matrix in, Matrix out



if inputs are,

$$x : 64 * 4096$$

$$y : 4096 * 4096$$

the size of the 4D Jacobian Tensor is,

$$(64 * 4996) * (4096 * 4096) * 62bi = 256GB$$

Should find another way to derive the gradient

$$x \neq y$$

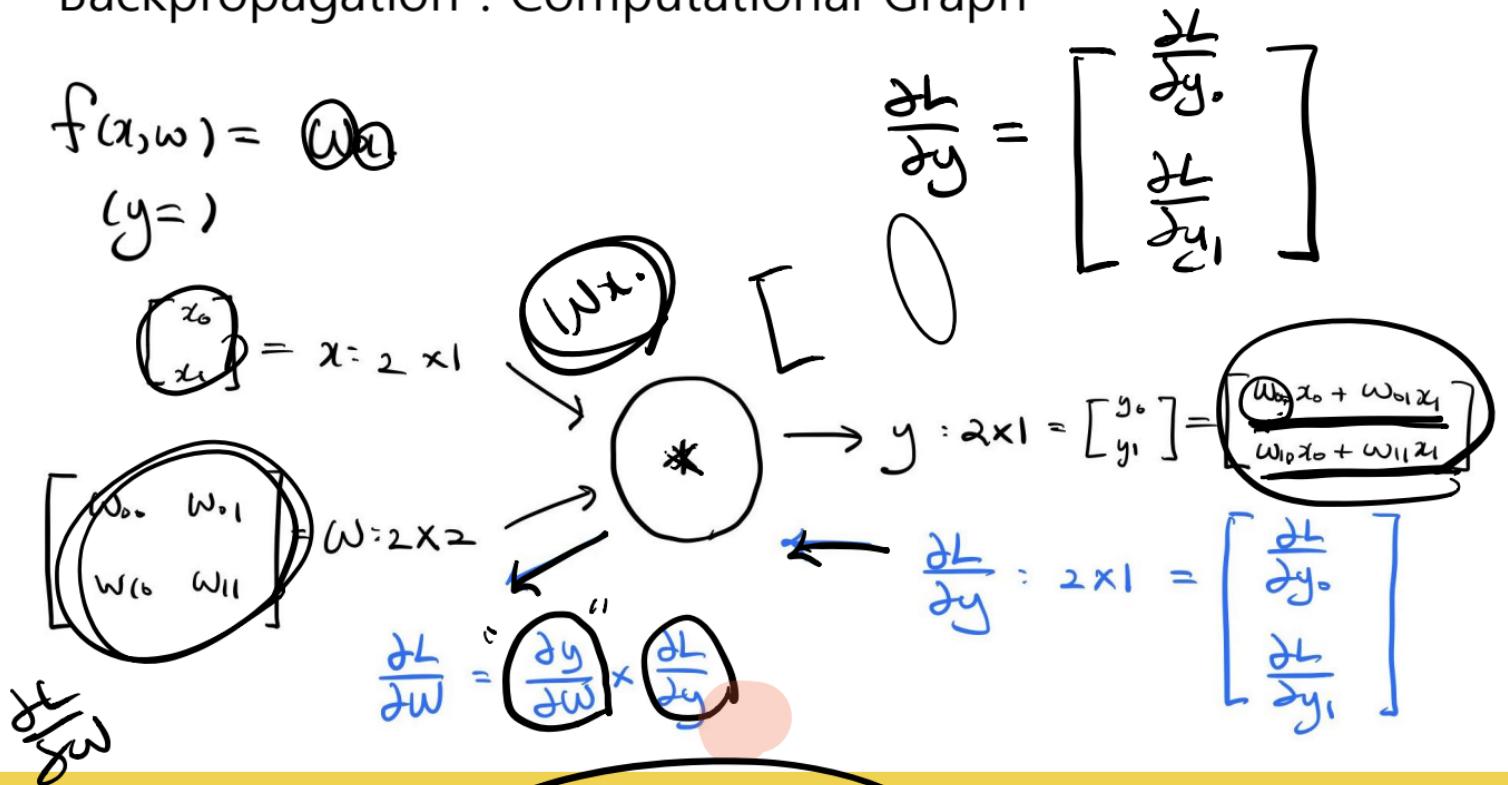
$$z = 64 \times 4096$$

$$\frac{\partial z}{\partial y}$$

$$\frac{\partial z}{\partial x}$$

$$\frac{\partial L}{\partial w} = \begin{bmatrix} 0 \\ 0 \end{bmatrix}$$

Backpropagation : Computational Graph



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

Backpropagation : Computational Graph "

$$\frac{\partial L}{\partial w}$$

instead of computing the General Jacobian, compute the gradient elementwise

$$\frac{\partial L}{\partial w} = \begin{bmatrix} \textcircled{0} \end{bmatrix}$$

$\frac{\partial L}{\partial w_{00}} = x_0 \frac{\partial L}{\partial y_0}, \quad \frac{\partial L}{\partial w_{01}} = x_1 \frac{\partial L}{\partial y_1}, \dots,$

$\frac{\partial L}{\partial w_{10}} = x_0 \frac{\partial L}{\partial y_1}, \quad \frac{\partial L}{\partial w_{11}} = x_1 \frac{\partial L}{\partial y_1}, \dots,$

$\frac{\partial L}{\partial w} = \frac{\partial L}{\partial y} x^T = \begin{bmatrix} \frac{\partial L}{\partial y_0} \\ \frac{\partial L}{\partial y_1} \end{bmatrix} [x_0 \ x_1]$

/ Backpropagation : Implementation $w \not\propto$

지금까지는, Computational Graph를 이용해,

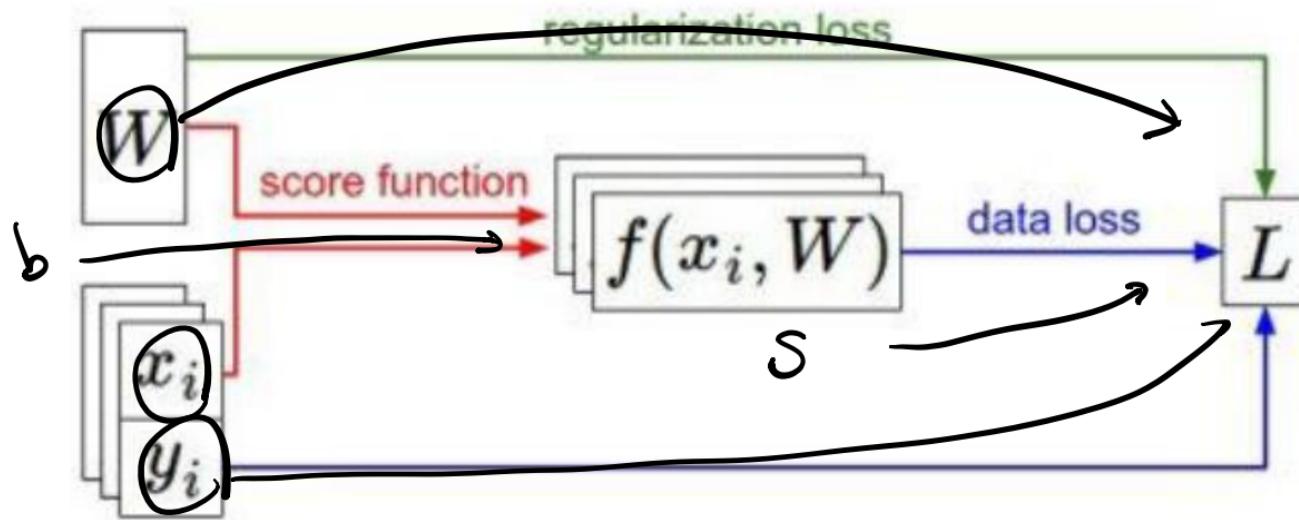
복잡한 함수의 gradient를 구하는 과정을 연습했다.

then, how can backpropagation be implemented in our linear classifier?

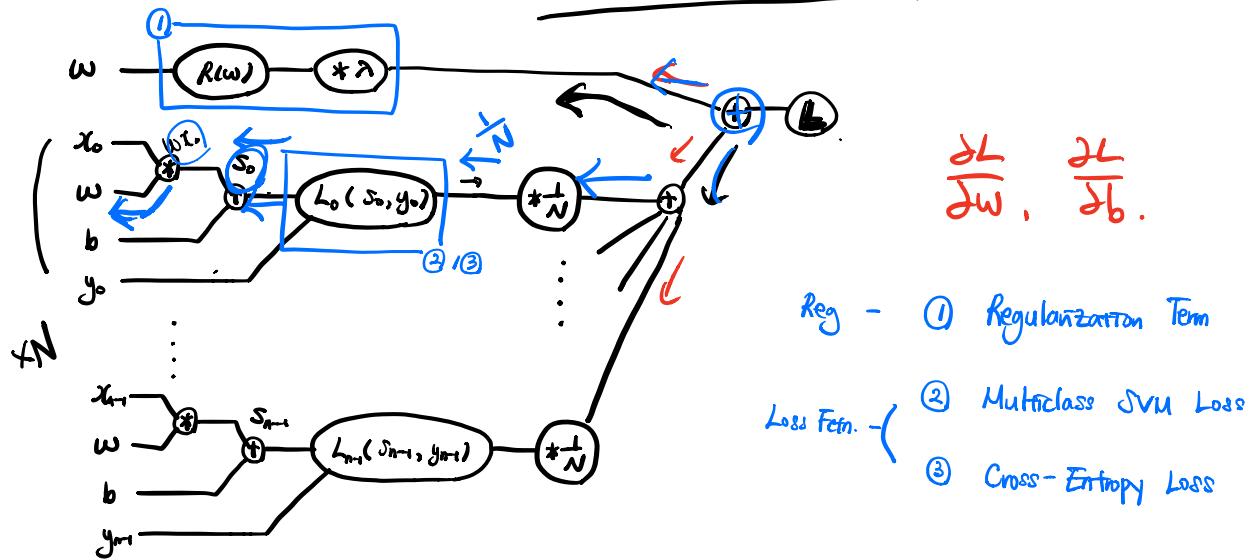


/

Backpropagation : Implementation



" Computational Graph of, $L = \frac{1}{N} \sum_{\bar{n}=1}^N L_{\bar{n}}(w x_{\bar{n}} + b, y_{\bar{n}}) + \lambda R(w)$ "



$$\frac{\partial R}{\partial w_{0,0}} = 2w_{0,0} \quad \dots \quad \frac{\partial R}{\partial w} = 2w$$

1) $\frac{dL}{dw} = ?$

① $R(w) = \sum_k \sum_a w_{ka}^2 \rightarrow \lambda \times R(w).$

$\frac{\partial}{\partial w} (\lambda R(w)) = \lambda \frac{\partial}{\partial w} (R(w)) = (2\lambda) \cdot w$

② M-SVM. $L_{\bar{n}} = \sum_{j \neq y_{\bar{n}}} \max(0, s_j - s_{y_{\bar{n}}} + \Delta)$

$s = w x_0 + b$

$$\frac{\partial L_{\bar{n}}}{\partial s} = \begin{bmatrix} \frac{\partial L}{\partial s_0} \\ \frac{\partial L}{\partial s_1} \\ \vdots \end{bmatrix}$$

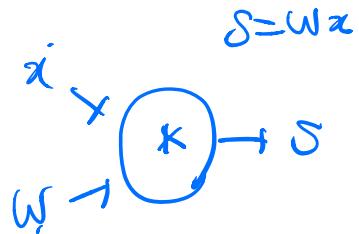
$\frac{\partial L_{\bar{n}}}{\partial w} = 1 \times ((s_k - s_{y_{\bar{n}}} + \Delta) \geq 0)$

$\frac{\partial L_{\bar{n}}}{\partial b} = \sum_{j \neq y_{\bar{n}}} (-1)(s_j - s_{y_{\bar{n}}} + \Delta) \geq 0$

$$\frac{\partial \bar{L}_i}{\partial s} = \begin{bmatrix} \vdots \\ \vdots \end{bmatrix}$$

$$\frac{\partial \bar{L}_i}{\partial s}$$

$$\begin{bmatrix} \frac{\partial \bar{L}_i}{\partial s_0} \\ \frac{\partial \bar{L}_i}{\partial s_1} \\ \vdots \end{bmatrix} \begin{bmatrix} x_0 & x_1 & \dots \end{bmatrix}$$



$$\frac{\partial s}{\partial w} = \boxed{\frac{\partial \bar{L}_i}{\partial s} x^T}$$

$$= \boxed{\frac{\partial \bar{L}_i}{\partial s_0} x_0, \frac{\partial \bar{L}_i}{\partial s_1} x_1, \dots}$$

$$\frac{\partial \bar{L}_i}{\partial w_k} = - \left(\sum_{j \neq y_i} 1 \cdot ((s_j - \delta_{y_i} + \Delta) \gamma_0) \right) \underline{x_i^T}$$

if $k = y_i$.

$\left(1 \cdot ((s_j - \delta_{y_i} + \Delta) \gamma_0) \right) x_i^T$ otherwise

③ Cross-Entropy

$$L_i = -\log \left(\frac{e^{s_{y_i}}}{\sum_j e^{s_j}} \right) \quad \frac{\partial L_i}{\partial s}$$

$$= -s_{y_i} + \log \left(\sum_j e^{s_j} \right) \quad \ln x \rightarrow \frac{1}{x}$$

$$\begin{aligned} \frac{\partial L_i}{\partial s_{y_i}} &= -1 + \frac{e^{s_{y_i}}}{\sum_j e^{s_j}} \\ \frac{\partial L_i}{\partial s_k} &= \frac{e^{s_k}}{\sum_j e^{s_j}} \end{aligned}$$

$$\frac{\partial L_i}{\partial s} \times x^T. = \frac{\partial L_i}{\partial w}$$

↓

$$\frac{\partial L}{\partial w} = (2\pi) \cdot W + \frac{1}{N} \sum_{i=1}^N \frac{\partial L_i}{\partial w} \quad \text{u.s.v.u.} \quad \text{ICE.}$$

!!

Backpropagation : Implementation

Implement the ~~AGD~~ (or just Gradient Descent Algorithm)

to your Linear Classifier!

“

“

Model

Model. fit()

Backpropagation : Implementation

앞으로 나올 MLP(Multilayer Perceptron), CNN(Convolutional Neural Network)는

훨씬 복잡한 Computational Graph 가짐.

but, 대부분의 function (or node)의 gradient는 검색을 통해 찾을 수 있으며,

대부분 라이브러리에 이미 구현되어 있다!

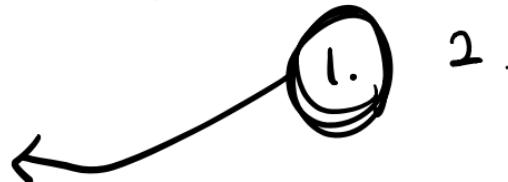
Review

1. "Optimization"

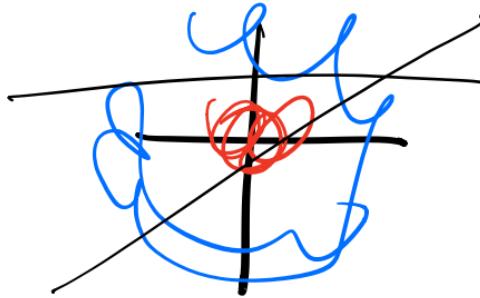
- SGD vs. Random Search
- Analytic Gradient vs. Numeric Gradient

2. "Backpropagation"

- Chain Rule
- Computational Graph
 - scalar .
 - vector .
 - matrix .
- Implementation •



Preview on Next Class

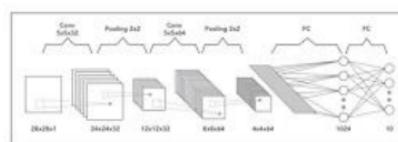
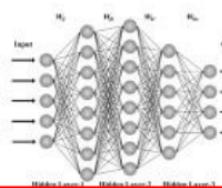


Optimization : Calculating Gradient

IDEA 2. Analytic

지금은 simple linear classifier만 다루고 있지만,

앞으로 다른 model은



then, how should we derive the gradient?

Multilayer Perceptron

Linear Classifier 단독으로는

풀 수 없는 문제들 존재...

여러 개를 중첩시켜 보자!