

Backprop and Neural Networks

Named Entity Recognition (NER)

↳ find and classify names in text

ex. Samuel Quinn was arrested in the Hilton Hotel in Paris in April 1989.

PER PER LOC LOC LOC DATE DATE

How? build up a context window of neighboring words

hand-labeled data \rightarrow logistic classifier 훈련

$$J_t(\theta) = \sigma(s) = \frac{1}{1 + e^{-s}}$$

predicted model
probability of class

Word vector 를 neural net에 넣고 행렬곱, 학습률 변환

Computing gradient by hand

Matrix calculus : fully vectorized gradients

Gradients

- one input & one output

$$f(x) = x^3$$

- It's gradient is its derivative

$$\frac{df}{dx} = 3x^2$$

- "input을 조금 바꾸면 output이 얼마나 바뀔까?"

- n inputs & one output

$$f(x) = f(x_1, x_2, \dots, x_n)$$

- Its gradient is a [vector of partial derivatives] with respect to each output

$$\frac{\partial f}{\partial x} = \left[\frac{\partial f}{\partial x_1}, \frac{\partial f}{\partial x_2}, \dots, \frac{\partial f}{\partial x_n} \right]$$

- n inputs & m outputs

$$f(x) = [f_1(x_1, x_2, \dots, x_n), \dots, f_m(x_1, x_2, \dots, x_n)]$$

- It's Jacobian is an $M \times n$ matrix of partial derivatives

$$\frac{\partial f}{\partial x} = \begin{bmatrix} \frac{\partial f_1}{\partial x_1} & \dots & \frac{\partial f_1}{\partial x_n} \\ \vdots & \ddots & \vdots \\ \frac{\partial f_m}{\partial x_1} & \dots & \frac{\partial f_m}{\partial x_n} \end{bmatrix} \quad \left(\frac{\partial f}{\partial x} \right)_{ij} = \frac{\partial f_i}{\partial x_j}$$

\Rightarrow Have every possible partial derivative of an output variable

Chain Rule

- For composition of one-variable functions: multiply derivatives

$$z = 3y$$

$$y = x^2$$

$$\frac{dz}{dx} = \frac{dz}{dy} \frac{dy}{dx} = 3 \cdot 2x = 6x$$

- For multiple variables at once: multiply Jacobians

$$h = f(z)$$

$$z = Wx + b$$

$$\frac{\partial h}{\partial x} = \frac{\partial h}{\partial z} \frac{\partial z}{\partial x} = \dots$$

$$(1) \frac{\partial}{\partial x} (Wx + b) = W, \quad (2) \frac{\partial}{\partial b} = I, \quad (3) \frac{\partial}{\partial u} (u^T h) = h^T$$

$\alpha(x)$.

$$\left. \begin{array}{l} S = u^T h \\ h = f(Wx + b) \\ x \end{array} \right\} \text{일 때, } \frac{\partial S}{\partial b} \text{ 찾기}$$

1. Break up equations into simple pieces

$$\begin{aligned} S &= u^T h \\ h &= f(Wx + b) \Rightarrow h = f(z) \\ x & \qquad \qquad z = Wx + b \end{aligned}$$

2. Apply the chain rule

$$\begin{aligned} \frac{\partial S}{\partial b} &= \cancel{\frac{\partial S}{\partial h} \frac{\partial h}{\partial z} \frac{\partial z}{\partial b}} \quad \text{Hadamard product} \\ &= \underbrace{u^T \cdot \text{diag}(f'(z)) \cdot I}_{\delta} = \underbrace{u^T \circ f'(z)}_{\delta} \end{aligned}$$

Re-using computation

$$\begin{aligned} \frac{\partial S}{\partial W} &= \cancel{\frac{\partial S}{\partial h} \frac{\partial h}{\partial z} \frac{\partial z}{\partial W}} \\ &= \underbrace{u^T \cdot \text{diag}(f'(z))}_{\delta} \cdot \frac{\partial z}{\partial W} \end{aligned}$$

$$= \underbrace{\delta^T}_{\text{local error signal}@z} \underbrace{x^T}_{\text{local input signal}}$$

* Why transpose?

$$\frac{\partial S}{\partial W} = \delta^T x^T$$

$[n \times m] \quad [n \times 1] \quad [1 \times m]$

this makes the dimensions work out!

Backpropagation 예제

re-use derivatives computed for higher layers in computing derivatives for lower layers to minimize computation

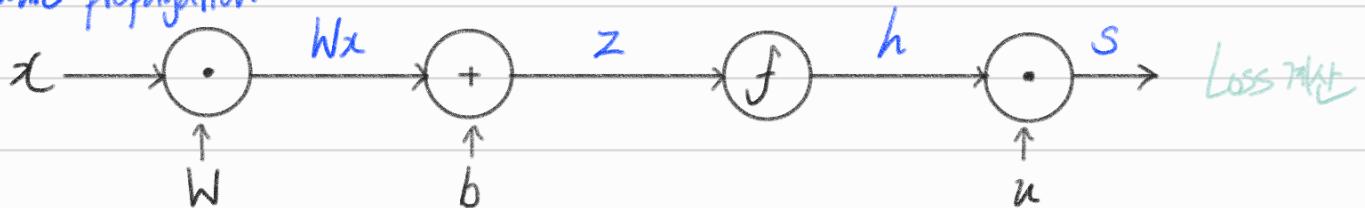
input: x

$$z = Wx + b$$

$$h = f(z) \quad (\text{Ex. sigmoid})$$

$$s = u^T h$$

Forward propagation



Backpropagation

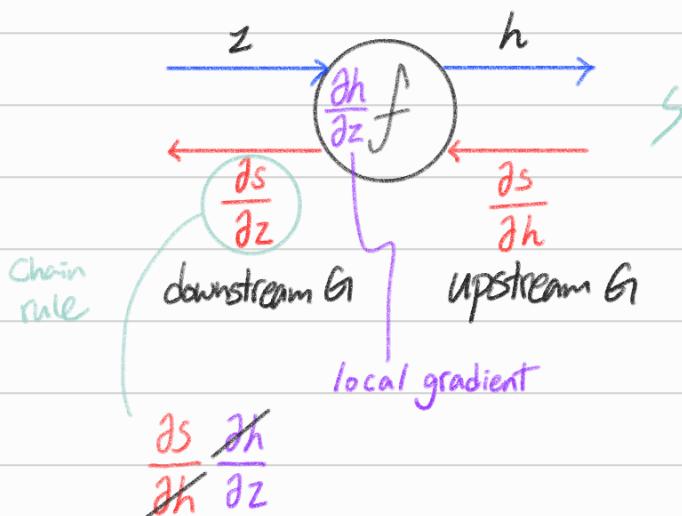
$$\frac{\partial s}{\partial b} \quad \frac{\partial h}{\partial z} \quad \frac{\partial s}{\partial h} \quad \frac{\partial s}{\partial s}$$

← ← ← ←

- pass along gradients

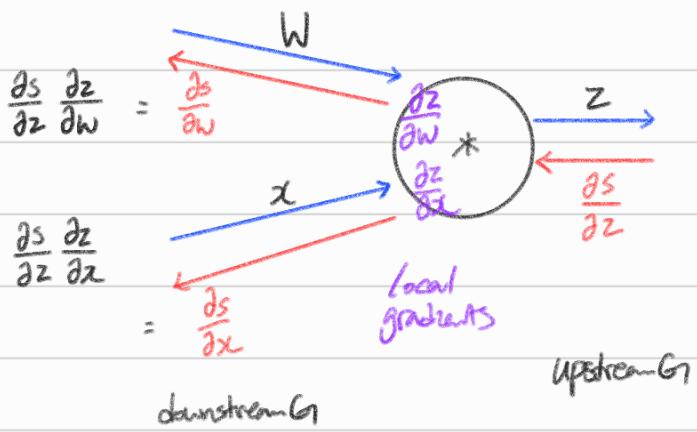
Single Node

- Node receives an "upstream gradient"
- Goal is to pass on the correct "downstream gradient"

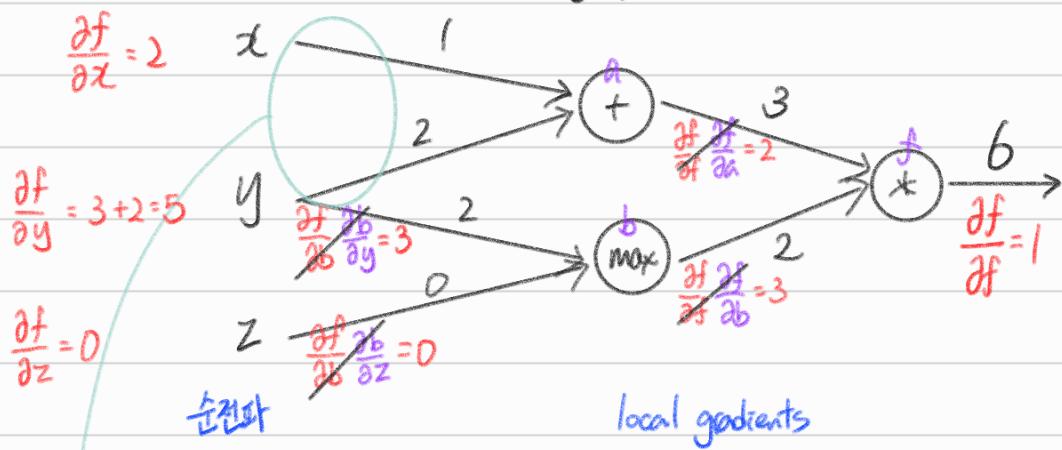


즉, downstream gradient = upstream gradient \times local gradient

Multiple inputs



$$\text{ex. } f(x, y, z) = (x+y) \max(y, z) \quad x=1, y=2, z=0$$



$$a = x+y$$

$$b = \max(y, z)$$

$$f = ab$$

$$\frac{\partial a}{\partial x} = 1, \frac{\partial a}{\partial y} = 1$$

$$\frac{\partial b}{\partial y} = 1 (y>z) = 1, \frac{\partial b}{\partial z} = 1 (z>y) = 0$$

$$\frac{\partial f}{\partial a} = b_2, \frac{\partial f}{\partial b} = a_3$$

$$\frac{\partial f}{\partial x} = \frac{\partial f}{\partial a} \frac{\partial a}{\partial x} = 2 \times 1 = 2$$

$$\frac{\partial f}{\partial y} = \frac{\partial f}{\partial a} \frac{\partial a}{\partial y} = 2 \times 1 = 2$$

$$\text{Downstream} = \text{Upstream} \times \text{Local}$$

$$D = U \cdot L$$

Node Intuitions

- + : distributes the upstream gradient
- max : rantes "
- * : switches "

Top3: Compute all the gradients at once

Analogous to using δ when we computed gradients by hand

F prop : visit nodes in topological sort order

B prop :

- initialize output gradient
- visit nodes in reverse order

big O() complexity of fprop & bprop is the same