

Lecture 6 : Back Propagation

This will be your least favorite lecture, since it requires the most tedious derivations of the whole course. 하하하하

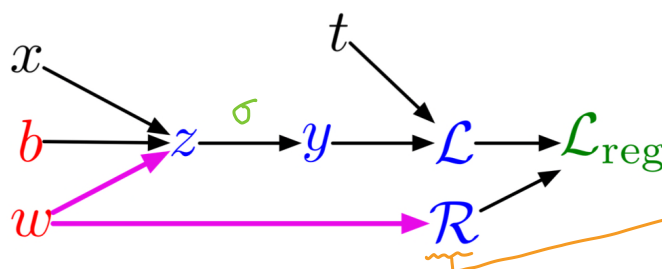
※ Lecture 6 , Lecture 8 둘다 BP 에 관한 챕터이며
6은 수학적 접근을 통한 자세한 이해를 위한 챕터 ,
8은 이를 굳이 매번 계산하지 않아도 되게끔 하는 automatic differentiation engine 구현에 대하여 설명한다.

1 Introduction

↳ BP란 신경망에서 손실함수에 대한 각 파라미터(w)의 편미분을 구하는 방식으로 이전에 Linear Regression, Logistic Regression 에서 진행한 것과 동일하게 이 편미분 값을 Gradient descent 에 활용한다.

2 The Chain Rule revisited

예시로 든 모델 설정 : 단일 변수 , 단일 출력 , 단일 샘플 (x, t)



$$z = wx + b \quad (1)$$

$$y = \sigma(z) \quad (2)$$

$$\mathcal{L} = \frac{1}{2}(y - t)^2 \quad (3)$$

$$\mathcal{R} = \frac{1}{2}w^2 \quad (4)$$

$$\mathcal{L}_{\text{reg}} = \mathcal{L} + \lambda \mathcal{R}. \quad (5)$$

Regularization term
이 강의에서는 크게 중요하지
않으니 이렇게 있구나 정도만.
(뒤에 다른 챕터에서 다룰)

이 상황에서 학습 = 성능개선 을 하고자 한다면

우리가 바꾸어야 하는 값 : w, b

⇒ 즉 우리는 \mathcal{L}_{reg} 에 대한 w 의 미분값 $\frac{\partial \mathcal{L}_{reg}}{\partial w}$, b 의 미분값 $\frac{\partial \mathcal{L}_{reg}}{\partial b}$ 가 필요

2.1 How you would have done it in calculus class

↳ 너네가 1~2학년때 미적분학에서 이것 어떻게 풀었냐면 ...

$$\begin{aligned}\mathcal{L}_{reg} &= \frac{1}{2}(\sigma(wx + b) - t)^2 + \frac{\lambda}{2}w^2 && \leftarrow \text{씩다 전개해서 풀었다 PTSD...} \\ \frac{\partial \mathcal{L}_{reg}}{\partial w} &= \frac{\partial}{\partial w} \left[\frac{1}{2}(\sigma(wx + b) - t)^2 + \frac{\lambda}{2}w^2 \right] \\ &= \frac{1}{2} \frac{\partial}{\partial w} (\sigma(wx + b) - t)^2 + \frac{\lambda}{2} \frac{\partial}{\partial w} w^2 \\ &= (\sigma(wx + b) - t) \frac{\partial}{\partial w} (\sigma(wx + b) - t) + \lambda w && \text{이제 Chain Rule!!} \\ &= (\sigma(wx + b) - t) \sigma'(wx + b) \frac{\partial}{\partial w} (wx + b) + \lambda w \\ &= (\sigma(wx + b) - t) \sigma'(wx + b) x + \lambda w \\ \frac{\partial \mathcal{L}_{reg}}{\partial b} &= \frac{\partial}{\partial b} \left[\frac{1}{2}(\sigma(wx + b) - t)^2 + \frac{\lambda}{2}w^2 \right] \\ &= \frac{1}{2} \frac{\partial}{\partial b} (\sigma(wx + b) - t)^2 + \frac{\lambda}{2} \frac{\partial}{\partial b} w^2 \\ &= (\sigma(wx + b) - t) \frac{\partial}{\partial b} (\sigma(wx + b) - t) + 0 \\ &= (\sigma(wx + b) - t) \sigma'(wx + b) \frac{\partial}{\partial b} (wx + b) \\ &= (\sigma(wx + b) - t) \sigma'(wx + b)\end{aligned}$$

★ 이러한 방식의 문제점

1. 계산식이 미친게 많다. + 실수를 하기 아주 쉽다 (위 식에서도 틀리게 있대요 찾아보아요)
2. 중복되는 term 이 너무 많다. $(wx + b)$ term은 값 도출까지 4번이나 적힌다.



BP 의 기본 아이디어는 이러한 중복되는 term을

하나의 모듈처럼 생각하여 공유하는 것이다 !

2.2 Multivariable chain rule: the easy case

↳ 다변수 일때는 Chain Rule 을 어떻게 적용할까?

review

단일 변수 일때 Chain Rule 은 다음과 같은 식이었다. $(f(g(x)))' = f(g(x))' \cdot g(x)' = \frac{df}{dg} \cdot \frac{dg}{dx}$

그럼 $(f(g(x), h(x)))' = ?$

Changing x slightly has two effects: it changes g slightly, and it changes h slightly. Each of these effects causes a slight change to f . For infinitesimal changes, these effects combine additively.

↳ 강의에서는 다음과 같이 설명함. x 값을 조정하면 g, h 값이 달라져 f 값이 바뀌는데, 만약 그 변화량이 아주 작다면 f 에 대한 미분 계산이 각 변수의 chain rule 결과를 합하는 것으로 계산한다.

$$\therefore (f(g(x), h(x)))' = \frac{\partial f}{\partial g} \cdot \frac{dg}{dx} + \frac{\partial f}{\partial h} \cdot \frac{dh}{dx}$$

└─ 이게 핵심

2.3 An alternative notation

↳ 위에서 한걸 좀 세련되고 눈에 들어오는 notation 으로 바꾸자!

(공식 표기법은 아니나, 이어지는 강의 notation 이해를 위해 필수)

함수 L 을 이루는 여러 변수 중 어떤 한 변수 v 에 대해

$$\overline{v} \triangleq \frac{\partial L}{\partial v}$$

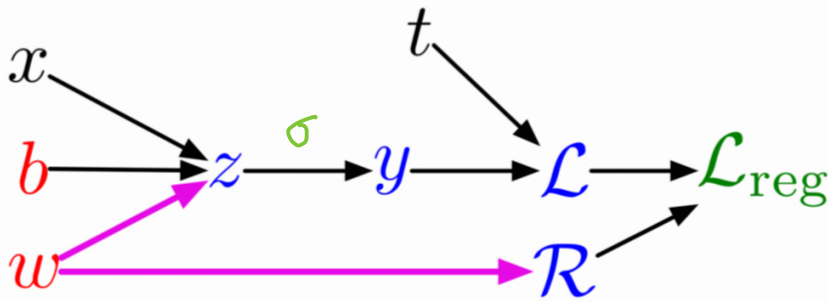
→ 신경망과 관련 지어 생각하자면,

- L = 편미분하고자 하는 가장 상위 함수 = 오차 함수
- v = 업데이트 하고자 하는 파라미터 = 특정 가중치

2.2 결과를 이 표기대로 바꾸면,

$$(f(g(x), h(x)))' = \frac{\partial f}{\partial g} \cdot \frac{dg}{dx} + \frac{\partial f}{\partial h} \cdot \frac{dh}{dx} = \overline{g} \cdot \frac{dg}{dx} + \overline{h} \cdot \frac{dh}{dx}$$

2.4 Using the computation graph



BP in pseudo code

위 연산 흐름 그래프에서 각 지점을 V , 해당 위치에서 편미분 값을 \bar{V} 라고 할때,

For $i = 1, \dots, N$

\Rightarrow Forward Pass

Compute v_i as a function of $\text{Pa}(v_i)$

$$v_N = 1$$

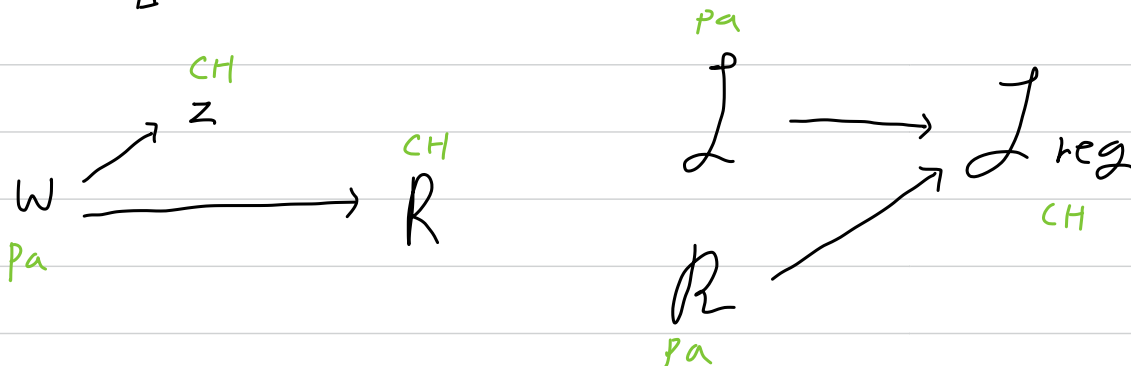
For $i = N - 1, \dots, 1$

\Rightarrow Backward Pass

$$\bar{v}_i = \sum_{j \in \text{Ch}(v_i)} \bar{v}_j \frac{\partial v_j}{\partial v_i}$$

$\text{Pa}() = \text{parent}$, $\text{Ch}() = \text{children}$ 을 의미하며

연산 흐름 그래프에서 부모 자식 관계는 다음과 같다.



수도 코드 흐름대로 $\frac{\partial \mathcal{L}_{\text{reg}}}{\partial w}(\bar{w})$, $\frac{\partial \mathcal{L}_{\text{reg}}}{\partial b}(\bar{b})$ 를 계산해보자.

↳ 2.1과 비교해보자.

$$\overline{\mathcal{L}_{\text{reg}}} = 1$$

$$\overline{\mathcal{R}} = \overline{\mathcal{L}_{\text{reg}}} \frac{d\mathcal{L}_{\text{reg}}}{d\mathcal{R}}$$

$$= \overline{\mathcal{L}_{\text{reg}}} \lambda$$

$$\overline{\mathcal{L}} = \overline{\mathcal{L}_{\text{reg}}} \frac{d\mathcal{L}_{\text{reg}}}{d\mathcal{L}}$$

$$= \overline{\mathcal{L}_{\text{reg}}}$$

$$\overline{y} = \overline{\mathcal{L}} \frac{d\mathcal{L}}{dy}$$

$$= \overline{\mathcal{L}} (y - t)$$

$$\overline{z} = \overline{y} \frac{dy}{dz}$$

$$= \overline{y} \sigma'(z)$$

$$\overline{w} = \overline{z} \frac{\partial z}{\partial w} + \overline{\mathcal{R}} \frac{d\mathcal{R}}{dw}$$

$$= \overline{z} x + \overline{\mathcal{R}} w$$

$$\overline{b} = \overline{z} \frac{\partial z}{\partial b}$$

$$= \overline{z}$$

최종정리

$$\overline{\mathcal{L}_{\text{reg}}} = 1$$

$$\overline{\mathcal{R}} = \overline{\mathcal{L}_{\text{reg}}} \lambda$$

$$\overline{\mathcal{L}} = \overline{\mathcal{L}_{\text{reg}}}$$

$$\overline{y} = \overline{\mathcal{L}} (y - t)$$

$$\overline{z} = \overline{y} \sigma'(z)$$

$$\overline{w} = \overline{z} x + \overline{\mathcal{R}} w$$

$$\overline{b} = \overline{z}$$

☆

• 중복 연산이 없다.

• 각 \bar{v} 가 마치 하나의 모듈처럼 처리되고 공유된다.

↳ ex) 만약 다른 Loss Function 을 쓰려면 전체를 다 계산할 필요없이 \bar{y} 만 바꾸면 된다.

3 Backprop on a multilayer net

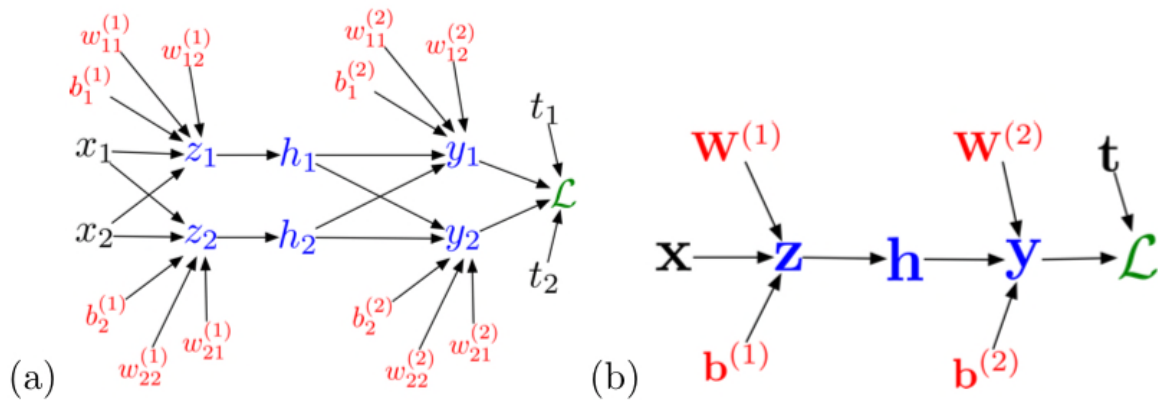


Figure 2: (a) Full computation graph for the loss computation in a multilayer neural net. (b) Vectorized form of the computation graph.

↳ 강의 자료 밑: "식이 조금 지저분해질 뿐, 똑같다",
" \bar{h} 만 조금 달라지는데 multivariate chain rule 이 적용된다. "

Forward - Pass

$$z_i = \sum_j w_{ij}^{(1)} x_j + b_i^{(1)}$$

$$h_i = \sigma(z_i)$$

$$y_k = \sum_i w_{ki}^{(2)} h_i + b_k^{(2)}$$

$$\mathcal{L} = \frac{1}{2} \sum_k (y_k - t_k)^2$$

Backward Pass

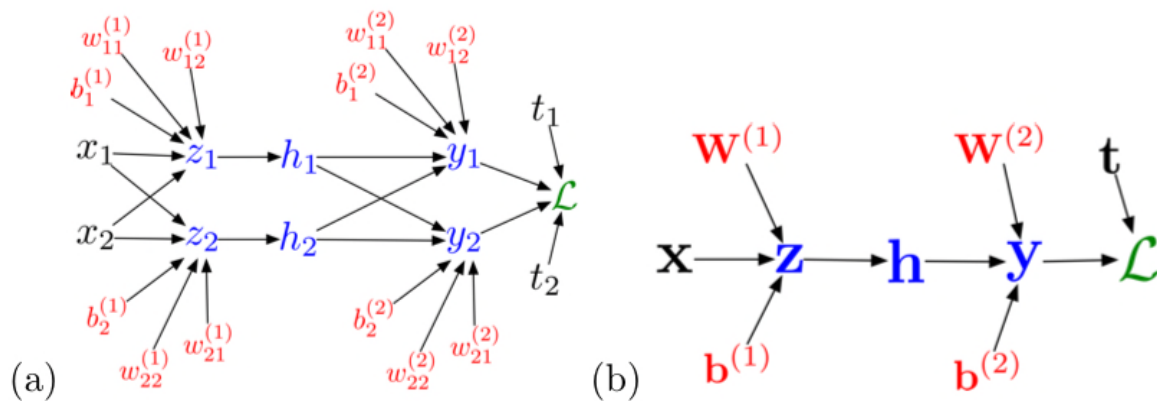


Figure 2: (a) Full computation graph for the loss computation in a multi-layer neural net. (b) Vectorized form of the computation graph.

$$\overline{\mathcal{L}} = 1$$

$$\overline{y_k} = \overline{\mathcal{L}} (y_k - t_k)$$

$$\overline{w_{ki}^{(2)}} = \overline{y_k} h_i$$

$$\overline{b_k^{(2)}} = \overline{y_k}$$

$$\overline{h_i} = \sum_k \overline{y_k} w_{ki}^{(2)} \quad \star \Rightarrow \text{화살표가 두우우우개}$$

$$\overline{z_i} = \overline{h_i} \sigma'(z_i)$$

$$\overline{w_{ij}^{(1)}} = \overline{z_i} x_j$$

$$\overline{b_i^{(1)}} = \overline{z_i}$$

식이 혹시 헷갈리다면

$$\overline{V_1} = \overline{V_2} \cdot \square \text{ 꼴의}$$

식에서 $V_1 \leftrightarrow V_2$ 사이
관계식을 Forward pass
에서 찾아서 V_1 으로
미분 해보자.

$$\bar{\mathcal{L}} = 1$$

$$\bar{y} = \bar{\mathcal{L}}(y - t)$$

$$\bar{W}^{(2)} = \bar{y} h^T$$

$$\bar{b}^{(2)} = \bar{y}$$

$$\bar{h} = W^{(2)T} \bar{y}$$

$$\bar{z} = \bar{h} \cdot \sigma'(z)$$

$$\bar{W}^{(1)} = \bar{z} x^T$$

$$\bar{b}^{(1)} = \bar{z}$$

⇒ Vectorized Notation

탐원 반응보고 skip 결정.

Note: $\bar{W}^{(2)} = 2 \times 2 \text{ matrix} = \begin{pmatrix} \bar{w}_{11} & \bar{w}_{12} \\ \bar{w}_{21} & \bar{w}_{22} \end{pmatrix}$
 $\bar{y} = 2 \times 1 \text{ matrix} = \begin{pmatrix} \bar{y}_1 \\ \bar{y}_2 \end{pmatrix}$
 $h^T = 1 \times 2 \text{ matrix} = (h_1, h_2)$

4 Appendix: why the weird notation?

↳ not important, why \bar{V} ?

+ Lecture 8 이 뭐하는지 그냥 그림만 한번 훑기