

---

---

---

---

---



**학습**: 훈련 데이터로부터 가중치 매개변수의 최적값을 자동으로 획득.  
(w, b)

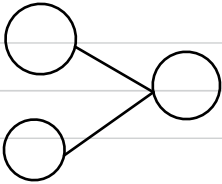
**손실함수**: 학습할 수 있도록 해주는 지표  
(loss function)

optimizer

How. 경사하강법. 등등.

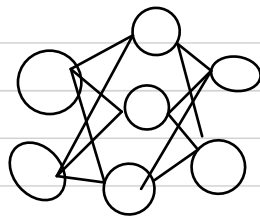
<이 손실함수의 결과 값을 가장 작게 만드는 가중치 매개변수를 찾는 것이 학습의 목표>

Perception



- 단층. 선형
- 직접 w, b를 선정

neural network



- 다층. 비선형
- 자동으로 w, b 선정

<기본 코딩>

$$X + \text{알고리즘} \Rightarrow Y$$

<deep learning> = end to end machine learning

$$X + Y \Rightarrow \text{알고리즘} (=가중치)$$

= '처음부터 끝까지'

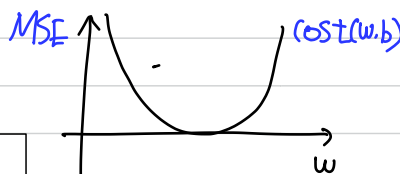
= '데이터에서 목표한 결과를 사람의 개입없이'

어떻게 이 지표를 가지고 자동으로 w, b를 가질수 있는 거지?

loss function

① MSE

$$\text{cost}(w, b) = \frac{1}{n} \sum_{i=1}^k (y_i - t_i)^2$$



○ y'의 활성화 함수가  
소프트맥스 시그모이드가 아닐때.  
(지수함수 활용)

$$\begin{aligned} y_i &= h(w_i x_i + b) \\ &= h(w_1 x_1 + w_2 x_2 + \dots + b) \\ &= h(\theta_0 + \theta_1 x_1 + \theta_2 x_2 + \dots + \theta_n x_n) \\ &= h(\theta^T x) \end{aligned}$$

$\theta_0$  = bias

$\theta_i (i=1, 2, 3, \dots) = \text{weight}$

- ① ② 출력층 node가 1개일때
- ③ 출력층 node가 n개일때

○ 가질게는 가중치 매개변수에 대한  
손실 함수의 기울기.  
( $\frac{\partial}{\partial x}$ )

①

$$\text{cost}(w, b) = J(\theta) = \frac{1}{n} \sum_{i=1}^m (y_i - t_i)^2$$

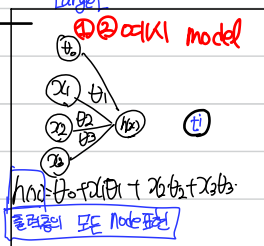
$$\begin{aligned} J(\theta) &= \frac{1}{2m} \sum_{i=1}^m (h(x_i) - t_i)^2 \\ &= \frac{1}{2m} \sum_{i=1}^m \left( \theta_0 + \theta_1 x_{i1} + \theta_2 x_{i2} + \theta_3 x_{i3} - t_i \right)^2 \end{aligned}$$

$$= \frac{1}{2m} \sum_{i=1}^m (h(x_i) - t_i) x_{i1}$$

ex)  $\frac{\partial J(\theta)}{\partial \theta_0} = \frac{1}{2m} \sum_{i=1}^m (y_i - t_i) \theta_0$  ( $\theta = x_0$ )

$$\therefore \frac{\partial J(\theta)}{\partial w_j} = \frac{1}{2m} \sum_{i=1}^m (y_i - t_i) x_{ij}$$

$$\therefore \frac{\partial J(\theta)}{\partial w} = \frac{1}{2m} X^T (y - t)$$



②

$$\begin{aligned} \text{cost}(w, b) &= \frac{1}{n} \sum_{i=1}^m (y_i - t_i)^2 \\ \frac{\partial \text{cost}(w, b)}{\partial w} &= \frac{1}{n} \sum_{i=1}^m \left( (w x_i + b) - t_i \right) \frac{\partial (w x_i + b)}{\partial w} \\ \frac{\partial \text{cost}(w, b)}{\partial w} &= \frac{1}{n} \sum_{i=1}^m (y_i - t_i) x_i \end{aligned}$$

③

$$\begin{aligned} \frac{\partial \text{cost}(w, b)}{\partial w} &= \frac{2}{n} \sum_{i=1}^m (y_i - t_i) x_i \\ &\rightarrow \text{정답 node} \\ &\rightarrow \text{출력층 모든 node의 편} \\ &\rightarrow \text{이것은 정답에 영향을 미칠뿐} \\ &\rightarrow \text{반응은 없음} \end{aligned}$$

③

$$\begin{aligned} \frac{\partial L_1}{\partial w_1} &= \frac{\partial L_1}{\partial w_2} \\ \frac{\partial L_2}{\partial w_1} &= \frac{\partial L_2}{\partial w_2} \\ \frac{\partial L_3}{\partial w_1} &= \frac{\partial L_3}{\partial w_2} \end{aligned}$$

$$\frac{1}{n} \sum_{i=1}^m (y_i - t_i)^2$$

$$\begin{pmatrix} \theta_0 \\ \theta_1 \\ \theta_2 \end{pmatrix} + \begin{pmatrix} x_1 \\ x_2 \end{pmatrix} \begin{pmatrix} w_{11} & w_{12} \\ w_{21} & w_{22} \\ w_{31} & w_{32} \end{pmatrix} = \begin{pmatrix} z_1 \\ z_2 \\ z_3 \end{pmatrix} \begin{pmatrix} y_1 \\ y_2 \\ y_3 \end{pmatrix}$$

$$b + wX = z, y$$

편향, 편향

무슨 문제를 풀고 싶냐, 즉 마지막 활성화 함수가 무엇에냐에 따라 쓰는 loss function이 다르다!

○  $y_i$ 가 선형으로 유지될 때 (ex  $h(x)=x$ )  

$$J = Wx + b$$

$y_i$ 가 비선형으로 있을 때  
 ex)  $h(x) = \text{sigmoid}(x)$   

$$y = \text{Sigmoid}(wx + b)$$

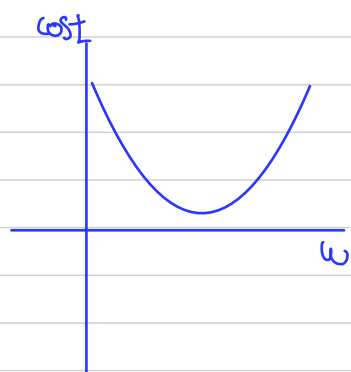
$$\text{cost} = \frac{1}{n} \sum_i (y_i - t_i)^2$$

$$= \frac{1}{n} \sum_i (wx + b - t_i)^2$$

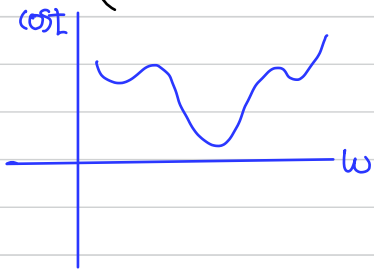
$$\text{cost} = \frac{1}{n} \sum_i (y_i - t_i)^2$$

$$= \frac{1}{n} \sum_i (\text{Sigmoid}(wx + b) - t_i)^2$$

$$= \frac{1}{n} \sum_i \left( \frac{1}{1 + e^{-(wx + b)}} - t_i \right)^2$$



○  $w$ 에 대한 2차 함수

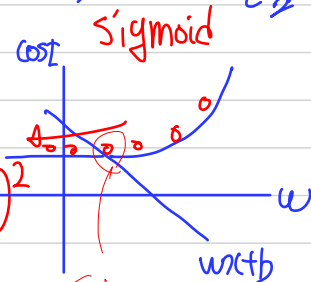


○ 심한 비볼록 형태  
 ○ 극소지역에 갇힐 위험

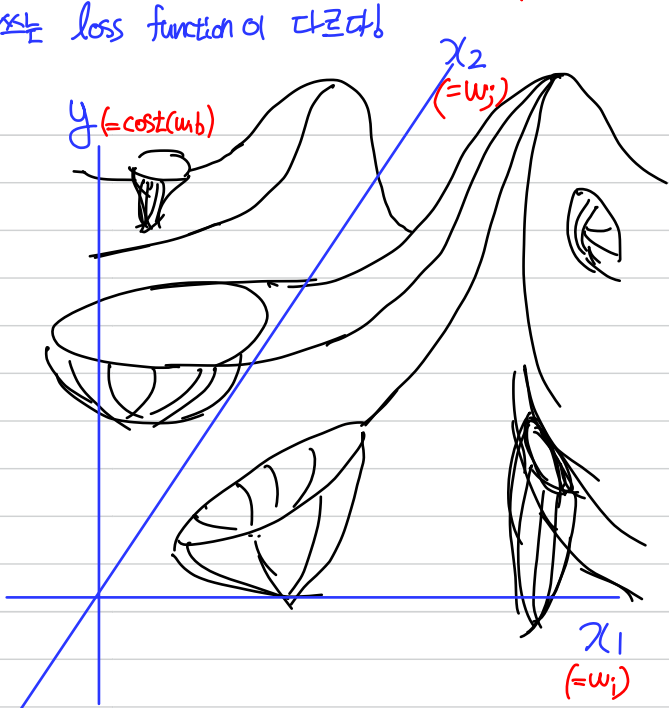
relu

$$\text{cost} = \frac{1}{n} \sum_i (\text{relu}(wx + b) - t_i)^2$$

$$= \frac{1}{n} \sum_i (\max(0, wx + b) - t_i)^2$$

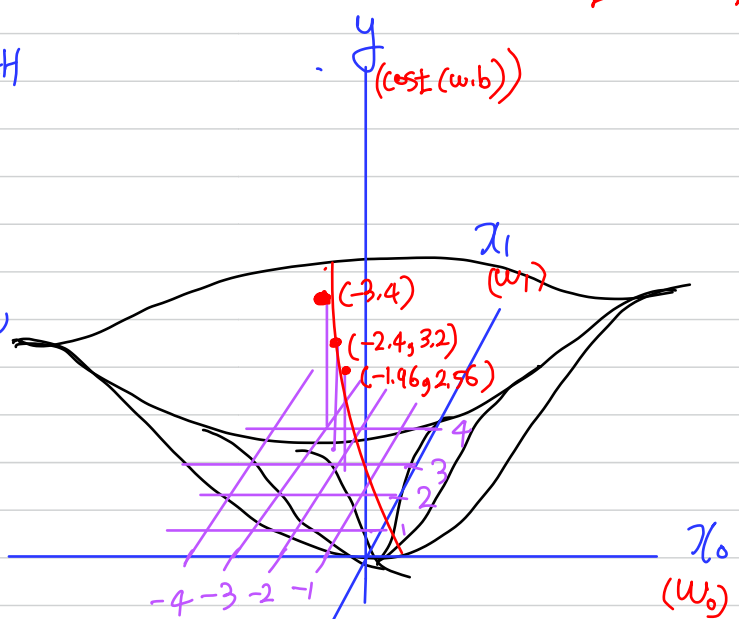


$$\text{cost} = \begin{cases} \frac{1}{n} \sum_i (wx + b - t_i)^2 & (wx + b \geq 0) \\ \frac{1}{n} \sum_i (0 - t_i)^2 & (wx + b < 0) \end{cases}$$



○ 여러가지 극소지역에 갇힐 위험이 있음

How solve



(gradient descent)  
 강하 하강법  
 (optimizer의 1종)

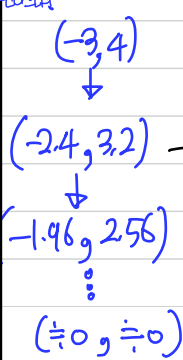
6장에서..

$$W = W - \eta \frac{dL}{dW}$$

$\frac{dL}{dW}$  기울기. 손실함수가 최소값이 되는 방향을 가리킴  
 평면에 땀줄수도 있어.  
 학습률: 너무 작으면 최소값에 도달하기도 전에, 너무 크게 하면 발산해버려. (전달 타고 넘어)  
 Step: 몇번 반복했나?  
 이상적으로 무한 반복하면 최소값은 같은 좌표  $(w_1, w_2)$ 로 옮겨져!

이제부터 좌표가 오른쪽으로 옮겨져서 Loss가 더 되려할 것 같은 쪽으로 좌표가 이동해!  $(w_1, w_2)$

(기울기가 양의 기울기면) 더하지 말고 빼줘라!



$y = x_0^2 + x_1^2$   
 $\frac{dy}{dx_0} = 2x_0$      $\frac{dy}{dx_1} = 2x_1$   
 ex)  $(-3, 4)$ 에서의 기울기,  $\alpha = 0.1$   
 solve)  $\frac{dy}{dx_0} \big|_{x_0=-3} = -6$      $\frac{dy}{dx_1} \big|_{x_1=4} = 8$   

$$x_0 = x_0 - \alpha \cdot \frac{dy}{dx_0}$$

$$-2.4 = -3 - (0.1)(-6)$$

$$x_1 = x_1 - \alpha \cdot \frac{dy}{dx_1}$$

$$3.2 = 4 - (0.1)(8)$$

$$-1.96 = -2.4 - (0.1)(-4.8)$$

$$2.56 = 3.2 - (0.1)(6.4)$$

○  $w$ 에 대한 이차함수로 표현이 가능  
 ○  $\frac{dL}{dW}$ : 기울기가 가려지는 방향은 함수의 등적값을 가장 크게 줄이는 방향  
 ○ activation function. 즉. 회귀함수 일때는 위와같은 MSE가 좋음  
 ○ 특성들 2개

# Loss function을 미분할때 activation function이 영향을 주는 상황

9  $h(x) = \text{Sigmoid}$  일때

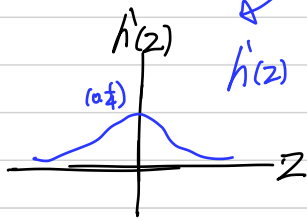
$$\frac{d \text{cost}}{dw} \cdot \frac{d \text{cost}}{dx} \text{의 차이}$$

$$\text{cost}(w, b) = \frac{1}{n} \sum_i (y_i - t_i)^2 = \frac{1}{n} \sum_i (h(wx+b) - t_i)^2 \quad \text{where } h(x) = \frac{1}{1 + e^{-x}}$$

$$\frac{d \text{cost}(w, b)}{dw} = \frac{1}{2n} \sum_i (y_i - t_i) \cdot (y_i)$$

$$= \frac{1}{2n} \sum_i (y_i - t_i) \cdot (h(wx+b))'$$

$$= \frac{1}{2n} \sum_i (y_i - t_i) \cdot (y_i)' \cdot x$$



$$h'(z) \quad (z = wx+b)$$

작은 미분계수 때문에 weight, bias가 제대로 업데이트 되지 않는다.

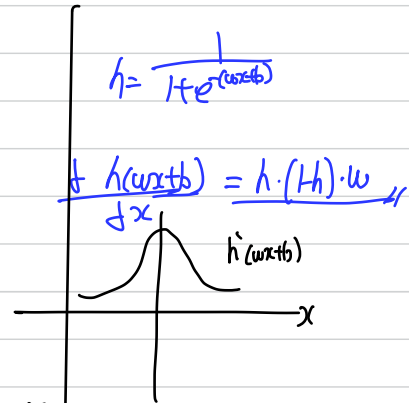
$z$ 가 너무 크거나 너무 작으면  $h'(z)$ 의 기울기 값이 0에 가까워진다.

⇒ 작은 미분계수 때문에 weight, bias가 제대로 업데이트 되지 않는다.

⇒ 비선형은 미분계수 변화율이 너무 적으니 갱신하기가 하극장인 걸리게 된다.

$$\frac{d \text{cost}(w, b; x_{ij})}{dw_{ij}} = \frac{1}{n} \sum_i (h(x_i) - t_i) h'(z) \cdot x_{ij}$$

$$\frac{d \text{cost}(w, b; x_{ij})}{db_{ij}} = \frac{1}{n} \sum_i (h(x_i) - t_i) h'(z)$$



$$\left( \frac{1}{1 + e^{-x}} \right) \left( 1 - \frac{1}{1 + e^{-x}} \right) \cdot 3$$

$$= \frac{1}{2} \cdot \left( 1 - \frac{1}{2} \right) \cdot 3$$

$$= \frac{3}{4} \quad (x=0) \checkmark$$

$$\begin{pmatrix} x_1 \\ x_2 \\ x_3 \\ \vdots \\ x_n \end{pmatrix} \rightarrow \text{다 열벡터가 아니다!}$$

⇒ 따라서 Sigmoid는 Cross Entropy를 사용하게 된다

global minimum이 있는데 local minimum에서 멈출수있는 MSE의 한계를 위해 Cross entropy를 사용!

나 다른 optimizer를 사용하면 계가 안날까? MSE의 한계가 아니까?

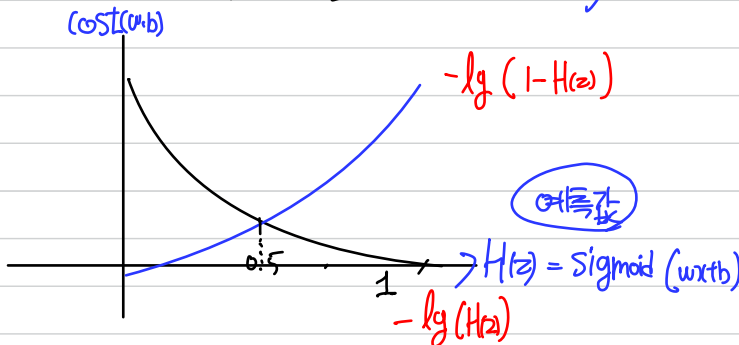
(20)

# Cross Entropy

→ Sigmoid와 같은 activation function 에 쓰는 Loss function.

## binary cross Entropy.

- 클럭값 0~1 ↔ 실제값 0 & 1  $\left. \begin{matrix} F, T \end{matrix} \right\} \Rightarrow \log$  함수가 적절
- 차이가 작은건 작게. 큰건는 크게



- 실제값이 1일때 예측값  $H(z)=1$  이면 우측은 0  $-\lg(H(z))$
- 실제값이 0일때 예측값  $H(z)=1$  이면 우측은  $\infty$   $-\lg(1-H(z))$

해고 통합하면

$$\begin{aligned} \text{cost}(H(z), y) &= -y \cdot \lg(H(z)) - (1-y) \cdot \lg(1-H(z)) \\ &= -\frac{1}{n} \sum_i (y_i \cdot \lg(H(z_i)) + (1-y_i) \cdot \lg(1-H(z_i))) \end{aligned}$$

$$\begin{aligned} \text{if) 실제값이 } 1 \Rightarrow y=1 \quad \text{cost}(H(z), y) &= -\lg(H(z)) \\ \text{" } 0 \Rightarrow y=0 \quad \text{cost}(H(z), y) &= -\lg(1-H(z)) \end{aligned}$$

즉 실제값  $y$ , 예측값  $H(z)$ 의 차이가 커지면 cost가 커지고.  
" " 차이가 작아지면 cost가 낮아진다.

일반화를 한다.

## Categorical Cross entropy

$$\text{cost}(H(z), y) = -(\underbrace{y_1}_{p_1} \cdot \lg(\underbrace{H(z)}_{p_2})) - (\underbrace{(1-y_1)}_{y_2} \cdot \lg(\underbrace{1-H(z)}_{p_2}))$$

$(= 0 \sim 1 \text{ 사이의 예측값})$

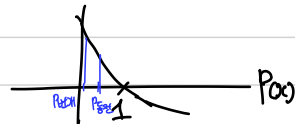
$$= -(y_1 \cdot \lg(p_1) + y_2 \cdot \lg(p_2))$$

$$\xrightarrow{\text{softmax classifier loss}} = -\sum_i y_i \cdot \lg(p_i) \quad \begin{matrix} \rightarrow \text{예측값} \\ \rightarrow \text{실제값} \end{matrix}$$

$$\therefore \text{cost}(H(z), y) = -\frac{1}{n} \sum_i \sum_j y_i^{(j)} \cdot \lg(p_j^{(i)})$$

entropy = 정보량의 기댓값

$$\lg\left(\frac{1}{p(x)}\right) = -\lg(p(x))$$



확률이 낮을수록 높은 정보량을 갖

$$\text{ex) } P(\text{동전}) = \frac{1}{2} \text{ VS } P(\text{한자}) = \frac{1}{256}$$

다시, 정보의 기댓값(entropy)는 어떤 랜덤변수  $X$ 에  $i=1, 2, \dots, n$  의 요소가 들어있을 때

$$-\sum_i A_i \cdot \lg(P(A_i))$$

그럼 cross entropy는

다른 사건 확률을 곱해서 entropy 계산  
= 두 개 확률분포  $P$ 와  $Q$ 에 대한  
하나의 사건  $X$ 가 갖는 정보량으로 장  
=  $Q$ 에 대한 정보량을  $P$ 에 대해 평균

$$H_{P,Q} = -\sum_i P(x_i) \cdot \lg Q(x_i)$$

$$\begin{aligned} H_{P,Q}(Y|X) &= -\sum_i \sum_{y \in \Omega(Y)} P(y_i | x_i) \cdot \lg Q(y_i | x_i) \\ &= -\sum_i \left[ \underbrace{P(y_i | x_i)}_{P(y_i)} \cdot \lg Q(y_i | x_i) + \underbrace{P(x_i | y_i)}_{P(x_i)} \cdot \lg Q(y_i | x_i) \right] \\ &= -\sum_i P(y_i | x_i) \cdot \lg Q(y_i | x_i) - (1 - P(y_i | x_i)) \lg(1 - Q(y_i | x_i)) \\ &= -\sum_i P(y_i) \cdot \lg Q(y_i) + \{1 - P(y_i)\} \lg(1 - Q(y_i)) \end{aligned}$$

Cross entropy

$$= -Q(x=0) \cdot \lg P(x=0) - Q(x=1) \cdot \lg P(x=1)$$

