


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

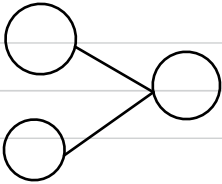
손실함수: 학습할 수 있도록 해주는 지표
(loss function)

Optimizer

How. 경사하강법. 등등.

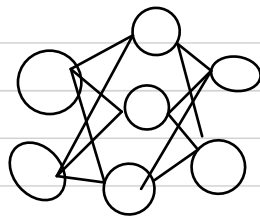
<이 손실함수의 결과 값을 가장 작게 만드는 가중치 매개변수를 찾는 것이 학습의 목표>

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'

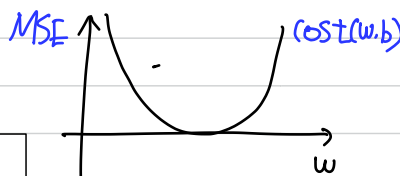
= '데이터에서 목표한 결과를 사람의 개입없이'

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_i 의 활성화 함수가
소프트맥스 시그모이드가 아닐때.
 $y_i = h(w_i x_i + b)$ (지수함수 형태)
 $= h(w_1 x_1 + w_2 x_2 + \dots + b)$
 $= h(\theta_0 + \theta_1 x_1 + \theta_2 x_2 + \dots + \theta_n x_n)$
 $= h(\theta^T x)$

θ_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.
(∂_x)

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

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

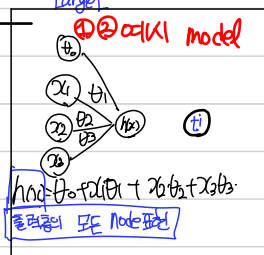
$$= \frac{1}{2m} \sum_{i=1}^m \left(\theta_0 + \theta_1 x_{i1} + \theta_2 x_{i2} + \theta_3 x_{i3} - t_i \right)^2$$

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

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

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

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



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

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

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

$$\frac{\partial \text{cost}(w, b)}{\partial w} = \frac{2}{n} \sum_{i=1}^m ((w x_i + b) - t_i) x_i$$

$\frac{\partial \text{Loss}}{\partial w_j} = \frac{2}{n} \sum_{i=1}^m (y_i - t_i) x_{ij}$

이것은 정답에 영향을 미칠뿐
반응은 없음 x

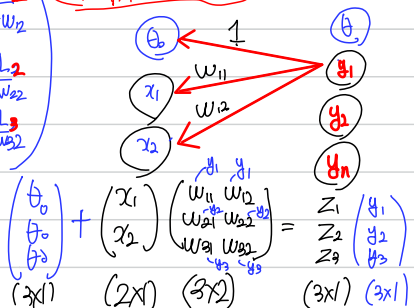
$w_{\text{new}} = w - \eta \frac{\partial \text{Loss}}{\partial w} \rightarrow \begin{pmatrix} w_{11} & w_{12} \\ w_{21} & w_{22} \end{pmatrix} = \begin{pmatrix} w_{11} & w_{12} \\ w_{21} & w_{22} \end{pmatrix} - \eta \begin{pmatrix} \frac{\partial \text{Loss}}{\partial w_{11}} & \frac{\partial \text{Loss}}{\partial w_{12}} \\ \frac{\partial \text{Loss}}{\partial w_{21}} & \frac{\partial \text{Loss}}{\partial w_{22}} \end{pmatrix}$

③ $\frac{\partial L_1}{\partial w_1} \frac{\partial L_1}{\partial w_2}$

$\frac{\partial L_2}{\partial w_{21}} \frac{\partial L_2}{\partial w_{22}}$

$\frac{\partial L_3}{\partial w_{31}} \frac{\partial L_3}{\partial w_{32}}$

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



$$b + wX = z, y$$

편향, 가중치

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) $h(x) = \text{sigmoid}(x)$

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

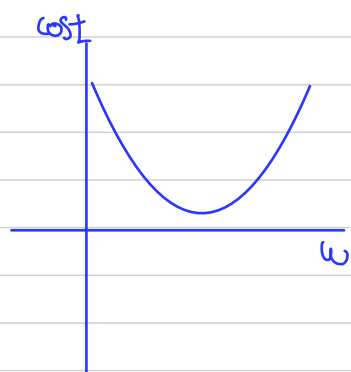
$$\text{cost} = \frac{1}{n} \sum_i (y_i - t_i)^2$$

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

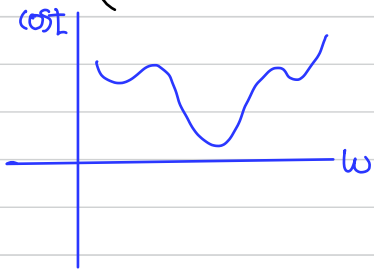
$$\text{cost} = \frac{1}{n} \sum_i (y_i - t_i)^2$$

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

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



○ w 에 대한 2차 함수

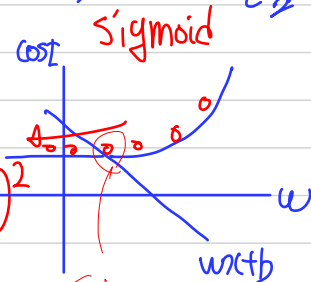


○ 심한 비볼록 형태
 ○ 극소지역에 갇힐 위험

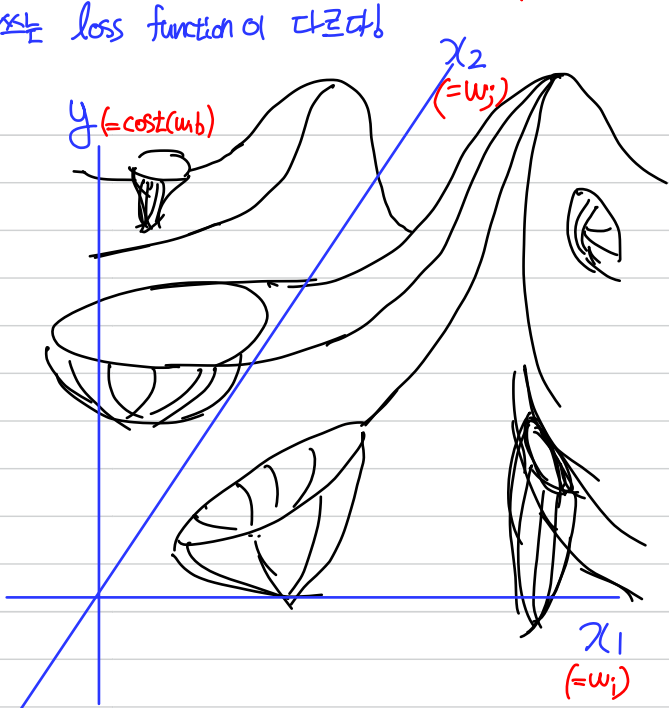
relu

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

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

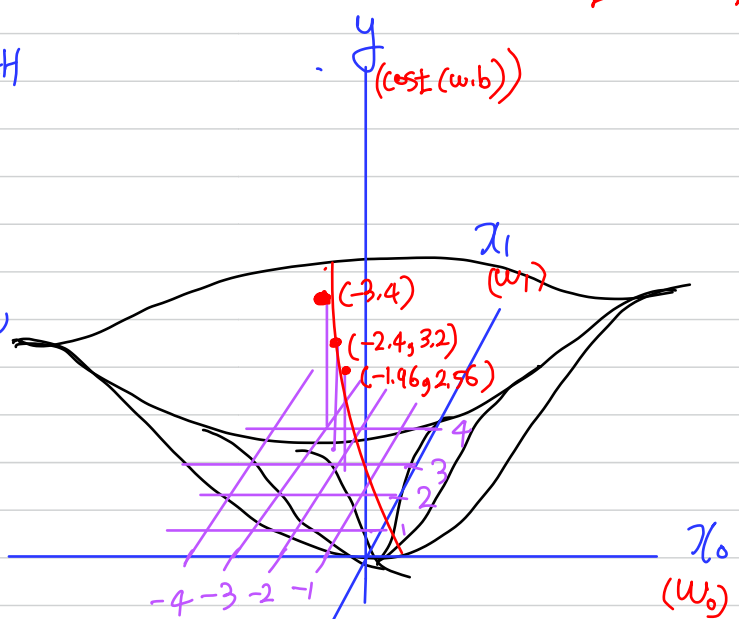


$$\text{cost} = \begin{cases} \frac{1}{n} \sum_i (wx_i + b - t_i)^2 & (wx + b \geq 0) \\ \frac{1}{n} \sum_i (0 - t_i)^2 & (wx + b < 0) \end{cases}$$



○ 여러가지 극소지역에 갇힐 상황이 많음

How solve



(gradient descent)
 경사 하강법
 (optimizer의 1종)

6장에서..

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

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

이런면 좌표가 오른쪽으로 옮겨져서 Loss가 더 되려할 것 같은 쪽으로 좌표가 이동하겠지!
 (완전, 왼쪽)

(기울기가 양의 기울기면)
 (더러지 말고 빼줘라!)

$(-3, 4)$
 \downarrow
 $(-2.4, 3.2)$
 \downarrow
 $(-1.96, 2.56)$
 \vdots
 $(\div 0, \div 0)$

$$y = x_0^2 + x_1^2$$

$$\frac{dy}{dx_0} = 2x_0 \quad \frac{dy}{dx_1} = 2x_1$$

ex) $(-3, 4)$ 에서의 기울기, $\alpha = 0.1$

sol) $\frac{dy}{dx_0} \big|_{x_0=-3} = -6 \quad \frac{dy}{dx_1} \big|_{x_1=4} = 8$

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

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

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

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

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

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

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

Loss function을 미분할때 activation function이 영향을 주는 상황

9 $h(x) = \text{Sigmoid}$ 일때

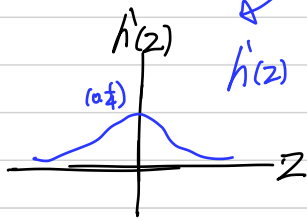
$$\frac{d \text{cost}}{dw} \cdot \frac{d \text{cost}}{dx} \text{의 차이}$$

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

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

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

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



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

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
강신화기에 하극장인 걸리게 된다.

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

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

$$\begin{aligned} h &= \frac{1}{1 + e^{-(wx+b)}} \\ \frac{d h(wx+b)}{dx} &= h \cdot (1-h) \cdot w \\ h'(wx+b) &= h \cdot (1-h) \\ h'(wx+b) &= \left(\frac{1}{1 + e^{-(wx+b)}} \right) \cdot \left(1 - \frac{1}{1 + e^{-(wx+b)}} \right) \\ &= \frac{1}{2} \cdot \left(1 - \frac{1}{2} \right) \cdot 3 \\ &= \frac{3}{4} \quad (x=0) \end{aligned}$$

$\begin{pmatrix} x_1 \\ x_2 \\ x_3 \\ \vdots \\ x_n \end{pmatrix} \rightarrow$ 다 열거가
어려움!

⇒ 따라서 Sigmoid는 Cross Entropy를 사용하게 된다

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

나 다른 optimizer를 사용하면 계가 안떨까? MSE의 한계가 아니까?

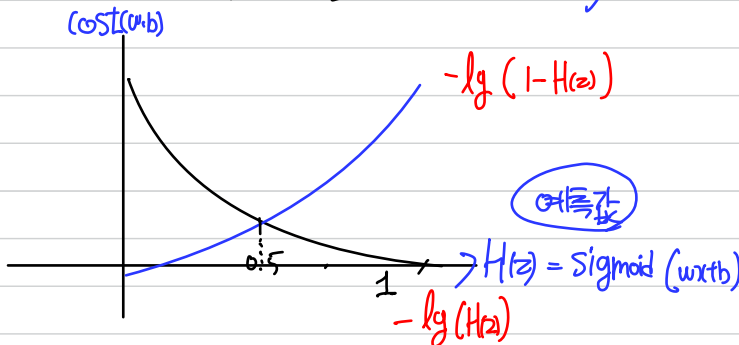
(20)

Cross Entropy

→ Sigmoid와 같은 activation function 에 쓰는 Loss function.

binary cross Entropy.

- 클럭값 0~1 \longleftrightarrow 실제값 0 & 1 $\xrightarrow{F, T}$
- 차이가 작은건 작게. 큰건는 크게 $\Rightarrow \log$ 함수가 적절



- 실제값이 1일때 예측값 $H(z)=1$ 이면 오차는 0 $-lg(H(z))$
- 실제값이 0일때 예측값 $H(z)=1$ 이면 오차는 ∞ $-lg(1-H(z))$

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

$$\text{cost}(H(z), y) = -(\underbrace{y_1}_{p_1} \cdot \lg(\underbrace{H(z)_1}_{p_1})) - (\underbrace{(1-y_1)}_{y_2} \cdot \lg(\underbrace{1-H(z)_1}_{p_2}))$$

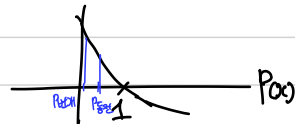
$$= -(y_1 \cdot \lg(p_1) + y_2 \cdot \lg(p_2))$$

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

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

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 잘

= Q 에 대한 정보량을 P 에 대해 평균

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

$$\begin{aligned} H_{P,Q}(Y|X) &= -\sum_i \sum_{y \in \Omega(Y)} P(y_i | x_i) \cdot \lg Q(y_i | x_i) \\ &= -\sum_i \left[\underbrace{P(y_1 | x_i)}_{P(y_1)} \cdot \lg Q(y_1 | x_i) + P(y_2 | x_i) \cdot \lg Q(y_2 | x_i) \right] \end{aligned}$$

$$= -\sum_i P(y_1 | x_i) \cdot \lg Q(y_1 | x_i) - (1 - P(y_1 | x_i)) \lg Q(y_2 | x_i)$$

$$= -\sum_i P(y_1) \cdot \lg Q(y_1) + \{1 - P(y_1)\} \lg Q(y_2) \quad (1-y_1)$$

Cross entropy

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

