

---

---

---

---

---



# Weight Decay (init-weight)

가중치 초기값은 큰 값으로 설정되면 안된다.  $\rightarrow \frac{dE}{dw_i}$  에서  $w_{t+1} = w_t - \eta \cdot \frac{dE}{dw_t}$  로

abs(w)가 큰 값을 가지면 sigmoid에선 0 or 1 모든 가중치의 값이 똑같이 갱신되기 때문에 reduction 폭을 좁게 줌  

$$\begin{pmatrix} w_{11} & w_{12} \\ w_{21} & w_{22} \end{pmatrix} = \begin{pmatrix} 2 & 2 \\ 2 & 2 \end{pmatrix} \rightarrow \text{가중치를 여러개 갖는 의미가 사라짐}$$

그럼 weight 들의 초기값이 학습에 영향을 줄까?

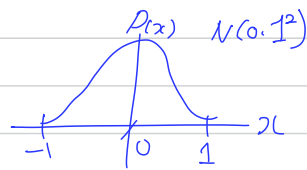
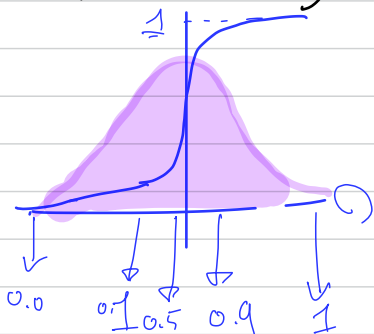
$\Rightarrow$  activation function에 따라 원하는

분포가 다름.

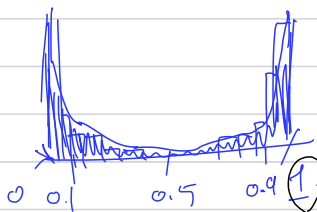
if) activation = Sigmoid(x)

$$w = N(0, 1)$$

$$\begin{pmatrix} w_{11} & w_{12} \\ w_{21} & w_{22} \end{pmatrix} = \begin{pmatrix} N(0, 1) & N(0, 1) \\ N(0, 1) & N(0, 1) \end{pmatrix}$$

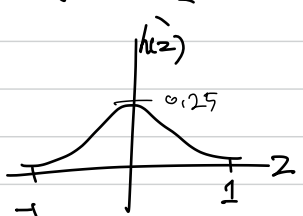


<활성화 분포>



histogram

모든  $h(x; w_i)$ 이 0 or 1을 갖는다.  
 출력값이자 입력값의  $y_i = 0$ 이면  $w_i \cdot y_i$  term 이 0이되기에 가중치가 갱신되지 않고  
 $y_i = 1$ 이면 back Propagation



$$\frac{dE}{dw} = \frac{dE}{dO} \frac{dO}{dz} \frac{dz}{dw} = \left(\frac{1}{2} - O\right) (h(z)) (1-h(z)) w$$

$\rightarrow$  Vanishing gradient Problem

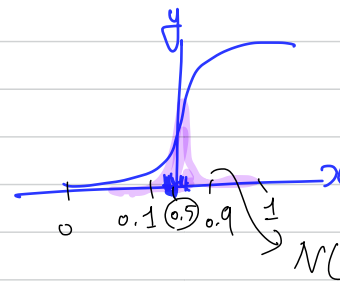
0이 되게 함

$\rightarrow$  갈수록  $\left(\frac{dE}{dw}\right)$  값, 미분값이 작아져 학습이 안되는 현상.

$$w_{t+1} = w_t + \eta \cdot \frac{dE}{dw}$$

$$w = N(0, 0.01)$$

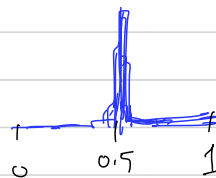
$$\begin{pmatrix} w_{11} & w_{12} \\ w_{21} & w_{22} \end{pmatrix} = \begin{pmatrix} N(0, 0.01) & N(0, 0.01) \\ N(0, 0.01) & N(0, 0.01) \end{pmatrix}$$



$$h\left(\begin{pmatrix} w_{11} & w_{12} \\ w_{21} & w_{22} \end{pmatrix}\right) = \begin{pmatrix} 0.5 & 0.5 \\ 0.5 & 0.5 \end{pmatrix}$$

0과 1과 관계  $0 \div N(0, 0.01)$

활성화 분포



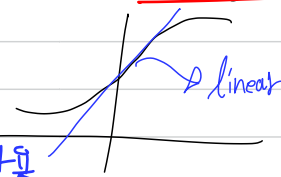
histogram

모든  $h(x; w_i)$  값들이 0.5로 고정된다.  
 $\Rightarrow$  가중치, layer를 쌓는 이유가 없다.  
 $\Rightarrow$  학습이 잘된다.  
 $\Rightarrow$  각각의 가중치가 asymmetry, 비대칭.  
 다양하게 가져야 좋은 학습이 된다.

# Xavier with Sigmoid

○ Xavier Glorot  
○ Yoshua Bengio  
○ 2010

- activation function 이 linear 할때!  $\rightarrow$  Sigmoid는 좌우대칭, 가운데 뾰족 선형.
- tanh 일때 가장 바람직!  
(S 자형)



선형적.

ex)  $w_1x + w_2x^2 + w_3x^3 + \dots$   
 $\Rightarrow$  선형적

- 앞 계층의 node가  $n$ 개 일때  $\sigma = \sqrt{\frac{1}{fan_{in}}}$ 의 분포를 사용

○ LeCun Normal initialization

- 뒤 층까지 고려한 분포  $\sigma = \sqrt{\frac{2}{fan_{in} + fan_{out}}}$

$$\sigma = \sqrt{\frac{6}{fan_{in} + fan_{out}}}$$

○ Xavier Normal initialization

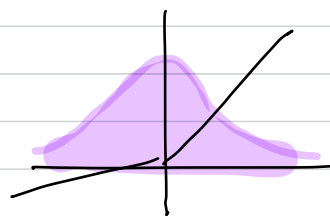
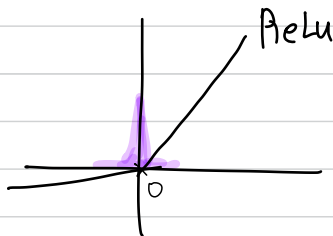
○ Xavier Uniform distribution

- input의 크기와 output의 크기가 커질수록 분산은 작게

# He with PReLU

○ 2015

- activation function 이 non-linear 할때
- PReLU 계열에 어울림



1 layer ...

5 layer

○ Normal distribution 과 Uniform distribution이 있는데 명확한 선택기준 없지만 Normal을 쓰는 드래곤과 he 논문에서 말함.

$N(0, 0.01^2)$

Xavier initialize

- 0에 몰려 weight가 갱신 불가

- layer에 얼마 못가

$N(0, 0.01^2)$  때와 마찬가지로 0에 몰려 weight update가 거의 되지 않음

- He initialization

앞 계층 node가  $n_{in}$ 개 일때

$$\sigma = \sqrt{\frac{2}{n_{in}}}$$

$$\sigma = \sqrt{\frac{6}{fan_{in}}}$$

$\rightarrow$  직관적으론 음수대 값들이 다 0이 됐기에 그만큼 다양성을 위해 그 값을 해준다.

- He normal distribution

- He uniform distribute

2 conv

$$\text{kernel shape} = \text{kernel size} + (\text{input channel, filters})$$

output R.G.B → feature map

$$(5 \times 5) \quad (3, 6) \Rightarrow \text{conv2D}(5, 5, 6)$$

$$= (5, 5, 3, 6)$$

input channel → output channel = filters.

1 conv

$$\text{kernel shape} = 4 + (1, 20)$$

$$= (4, 1, 20)$$

$$\text{receptive field} = \prod_{i=0}^{k-2} \text{kernel shape}[i] \quad \left( k = \text{kernel shape의 차원수} \right)$$

$$= 4 \quad (\text{conv1D}) \quad k=4$$

$$= 25 \quad (\text{conv2D}) \quad k=3$$

fan in = receptive field × input channel

(input tensor의 차원수)

$$= 25 \times 3 \quad (\text{kernel shape}[k-2]) \quad 2D$$

input channel

$$= 4 \times 1 \quad 1D$$

fan out = receptive field × output channel

(output tensor의 차원수)

$$= 25 \times 6 \quad (\text{kernel shape}[k-1]) \quad 2D$$

output channels

$$= 4 \times 20 \quad (\text{kernel shape}[k-1]) \quad 1D$$

# Batch Normalization

batch-norm 각 층의 활성화(들)를 적당히 퍼서  
가중치의 크기를 일관성 있게.

○ 2015 Sergey Ioffe, Christian Szegedy

○ 단순 whitening이 아닌 신경망안에 포함되어  
 $\mu, \sigma$ 를 조정

- scale과 shift 연산을 위해  $\gamma, \beta$ 가 들어가게 되어
  - 정칙화 부분을 원래대로 되돌리는 identity mapping 가능
  - 학습을 통해  $\gamma, \beta$ 를 정렬 수 있음.

- standard = 표준화  $\rightarrow N(0,1)$   $\frac{x-\mu}{\sigma}$
- min-max = 정칙화  $\rightarrow (0,1)$  범위  $\frac{x - \min X}{\max X - \min X}$   
(normalization)  $\left\{ \begin{array}{l} \text{정규화} \\ \text{정렬} \end{array} \right.$

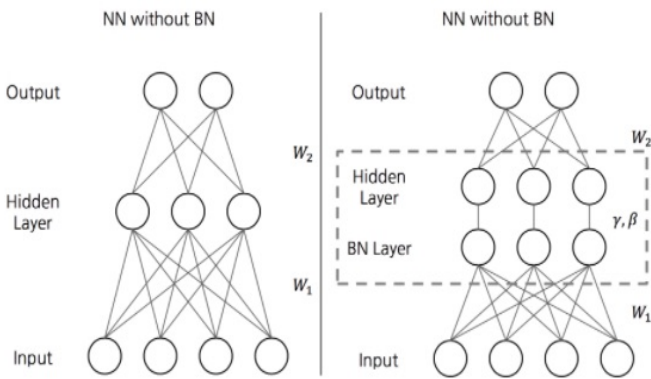
input  $\vec{x} = \{x_1, x_2, x_3, \dots, x_n\}$   
output  $BN_{\gamma, \beta}(\vec{x})$

mean  $\mu_{\vec{x}} = \frac{1}{n} \sum_{i=1}^n x_i$

Variance  $\sigma_{\vec{x}}^2 = \frac{1}{n} \sum_{i=1}^n (x_i - \mu_{\vec{x}})^2$

normalize  $\hat{x}_i = \frac{x_i - \mu_{\vec{x}}}{\sqrt{\sigma_{\vec{x}}^2 + \epsilon}}$   $\epsilon = 10^{-8}$

shift scale  $y_i = \gamma \cdot \hat{x}_i + \beta = BN_{\gamma, \beta}(\vec{x})$   
Scale shift output



○ 더 큰 learning rate 사용 가능.

internal covariate shift 감소, Parameter scale 영향(x)

더 큰 weight가 더 작은 gradient를 유도하게

Parameter growth가 안정화 되는 효과 (반복의 missing)

○ mini-batch를 어떻게 선택하느냐에 따라 data sample에 대한 다른 결과가 나와

$\Rightarrow$  drop out, 앙상블과 같은 효과, 의존도 떨어짐

$\Rightarrow$  더 general한 model을 learning 하는 효과가 생김.

## chain rule

i)  $\frac{dl}{dx_i} = \frac{dl}{dy_i} \cdot \frac{dy_i}{dx_i} = \frac{dl}{dy_i} \cdot \gamma$

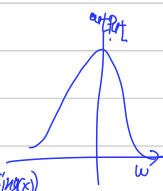
ii)  $\frac{dl}{d\sigma_{\vec{x}}^2} = \frac{dl}{d\hat{x}_i} \cdot \frac{d\hat{x}_i}{d\sigma_{\vec{x}}^2} = \frac{dl}{d\hat{x}_i} (x_i - \mu_{\vec{x}}) \cdot \frac{1}{\sqrt{\sigma_{\vec{x}}^2 + \epsilon}} \cdot \frac{-1}{2(\sigma_{\vec{x}}^2 + \epsilon)^{3/2}}$   
 $= \frac{1}{\sqrt{\sigma_{\vec{x}}^2 + \epsilon}} (x_i - \mu_{\vec{x}}) \cdot \frac{-1}{2(\sigma_{\vec{x}}^2 + \epsilon)^{3/2}}$   
 $= \sum_{i=1}^m \frac{dl}{d\hat{x}_i} (x_i - \mu_{\vec{x}}) \cdot \frac{-1}{2(\sigma_{\vec{x}}^2 + \epsilon)^{3/2}}$

iii)  $\frac{dl}{d\mu_{\vec{x}}} = \frac{dl}{d\sigma_{\vec{x}}^2} \cdot \frac{d\sigma_{\vec{x}}^2}{d\mu_{\vec{x}}} + \frac{dl}{d\hat{x}_i} \cdot \frac{d\hat{x}_i}{d\mu_{\vec{x}}}$   
 $= \sum_{i=1}^m \frac{dl}{d\sigma_{\vec{x}}^2} \cdot \left( 2 \cdot \frac{1}{n} \sum_{i=1}^n (x_i - \mu_{\vec{x}}) \right) + \frac{dl}{d\hat{x}_i} \cdot \left( \frac{-1}{\sqrt{\sigma_{\vec{x}}^2 + \epsilon}} \right) \cdot \sum_{i=1}^m \frac{1}{n}$

iv)  $\frac{dl}{dx_i} = \frac{dl}{d\hat{x}_i} \cdot \frac{d\hat{x}_i}{dx_i} + \frac{dl}{d\sigma_{\vec{x}}^2} \cdot \frac{d\sigma_{\vec{x}}^2}{dx_i} + \frac{dl}{d\mu_{\vec{x}}} \cdot \frac{d\mu_{\vec{x}}}{dx_i}$   
 $= \frac{dl}{d\hat{x}_i} \left( \frac{1}{\sqrt{\sigma_{\vec{x}}^2 + \epsilon}} \right) + \frac{dl}{d\sigma_{\vec{x}}^2} \cdot \frac{2}{n} (x_i - \mu_{\vec{x}}) + \frac{dl}{d\mu_{\vec{x}}} \cdot \frac{1}{n}$

v)  $\frac{dl}{d\gamma} = \frac{dl}{dy_i} \cdot \frac{dy_i}{d\gamma} = \sum_{i=1}^m \frac{dl}{dy_i} \cdot \hat{x}_i$

vi)  $\frac{dl}{d\beta} = \frac{dl}{dy_i} \cdot \frac{dy_i}{d\beta} = \sum_{i=1}^m \frac{dl}{dy_i} \cdot 1$



# Drop out.

`np.random.randn(10, 100) * 0.01`

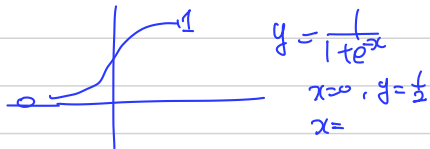
표준정규분포  $N(0, 0.01^2)$  = 표준편차가 0.01인  
정규분포

`6 * np.random.randn(...) + 11`

표정. 정규분포  $N(\mu, \sigma^2)$

$N(0, 1)$  를 init 으로 사용하는 weight로

5 layer sigmoid를 넘길때 활성화 함수의 출력값의 분포가 안정적



0과 1의 쪽에 치우치게 됨.  
 $\Rightarrow$  역전파의 기울기가 작아짐.

