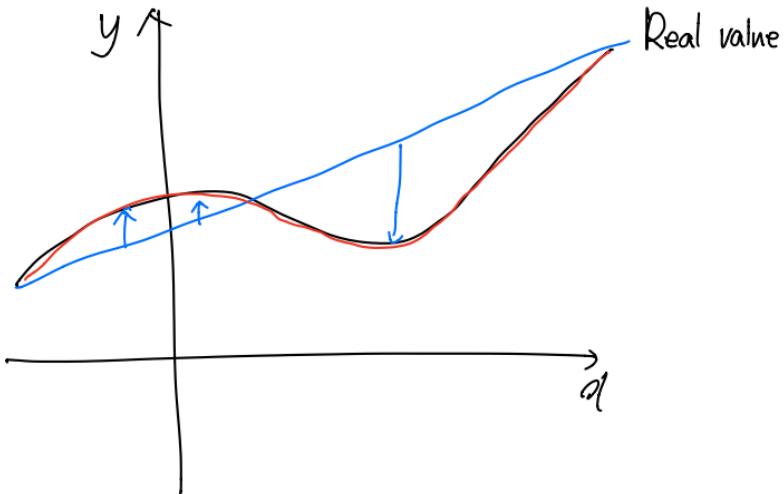


Problems of MLP

Overfitting

Model Capacity



Real value

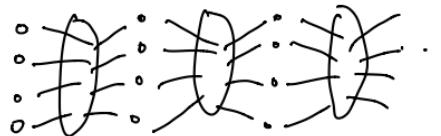
$$y = wx + b$$

$$y = w_1x^3 + w_2x^2 + w_3x + w_4$$

Parameter의 개수가 많음

Parameter의 개수가 많아지면 더 복잡한 모델을 예측 가능?

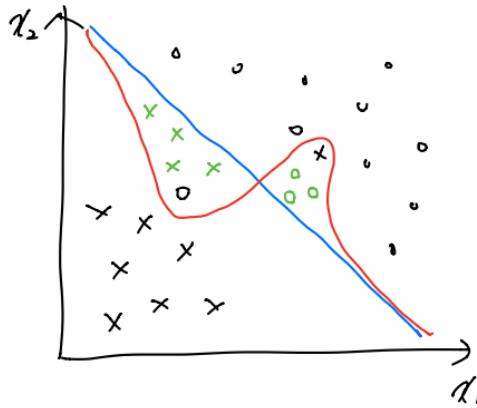
ex) MLP



Layer를 늘리면 parameter 개수 증가.

계속 늘려면 문제인가? → No

Overfitting

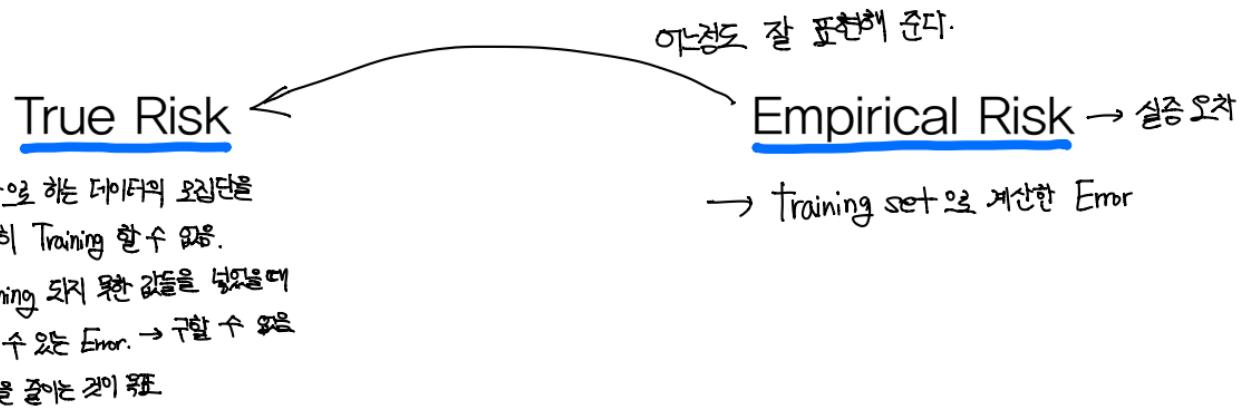


— : Simple한 예측

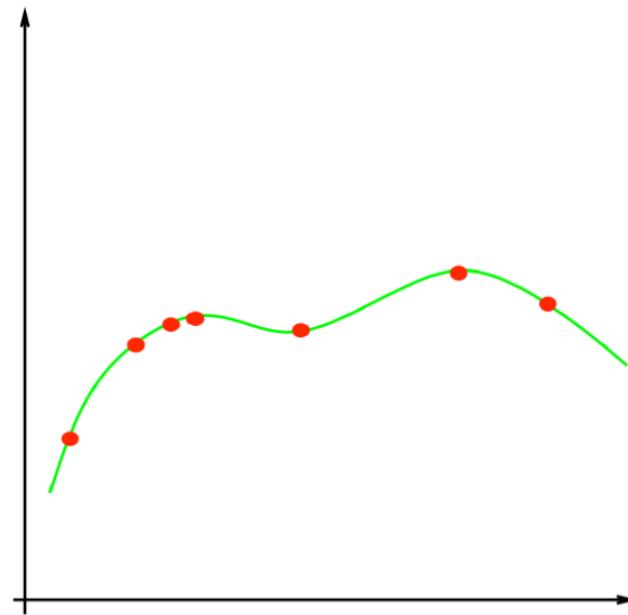
— : Layer를 늘려 세세한 부분까지 고려.

×, ○ : 실제 데이터들.
overfitting의 경우 이 칸들에 Error를 경험

Overfitting

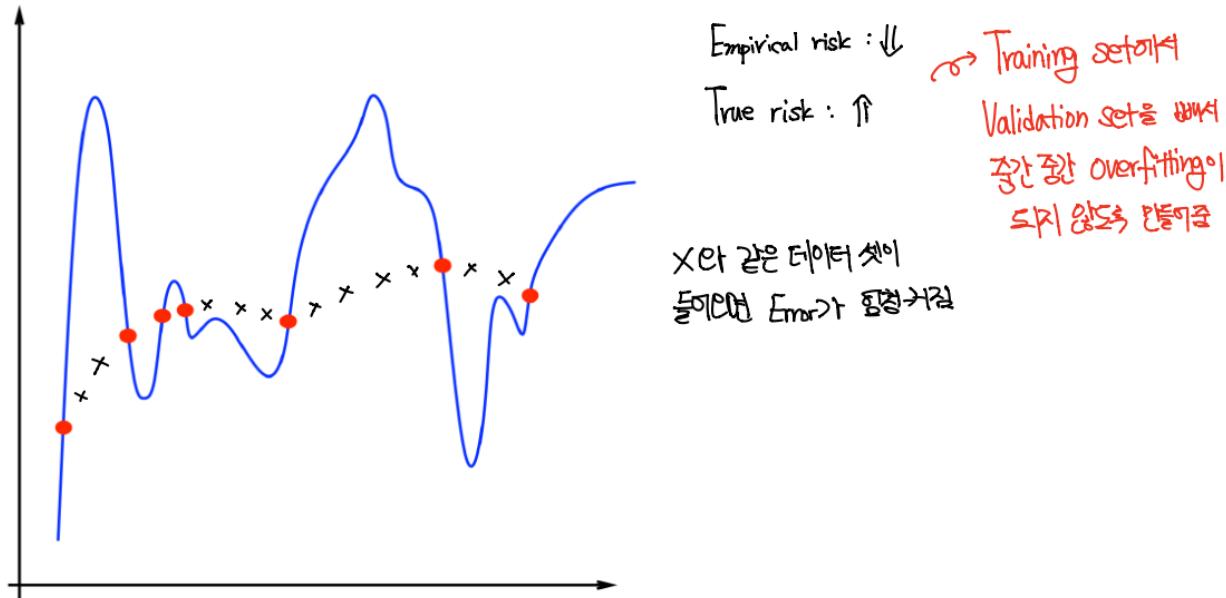


Overfitting



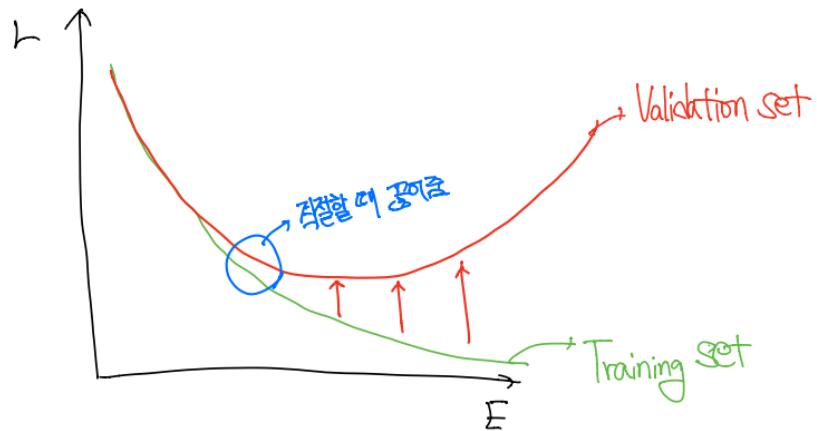
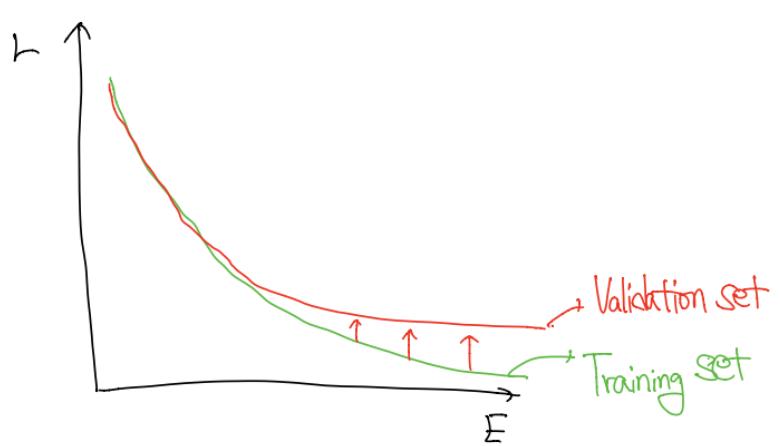
[T. Poggio]

Overfitting



[T. Poggio]

Overfitting



Summary

Overfitting occurs when model memorizes the train dataset thus cannot be applied to general dataset such as test set.

As model capacity increases, model become available to represent complicated systems but also be able to memorize the specific dataset.

We can know whether overfitting occurs or not by comparing training loss (empirical risk) and validation loss (approximated true risk)

Regularizations

L2 Regularization and Dropout

General data set에
적합하도록 만드는 것

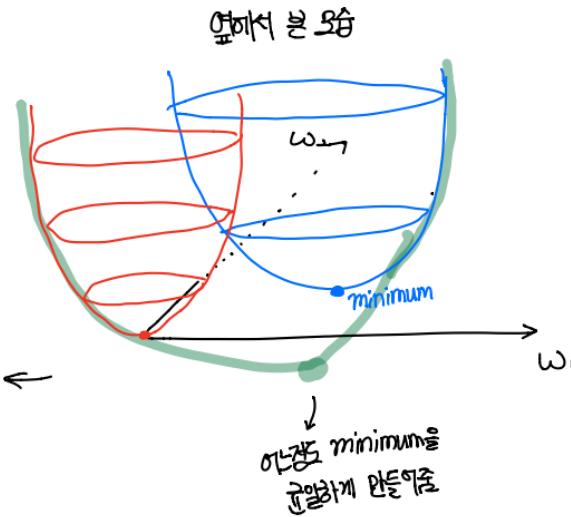
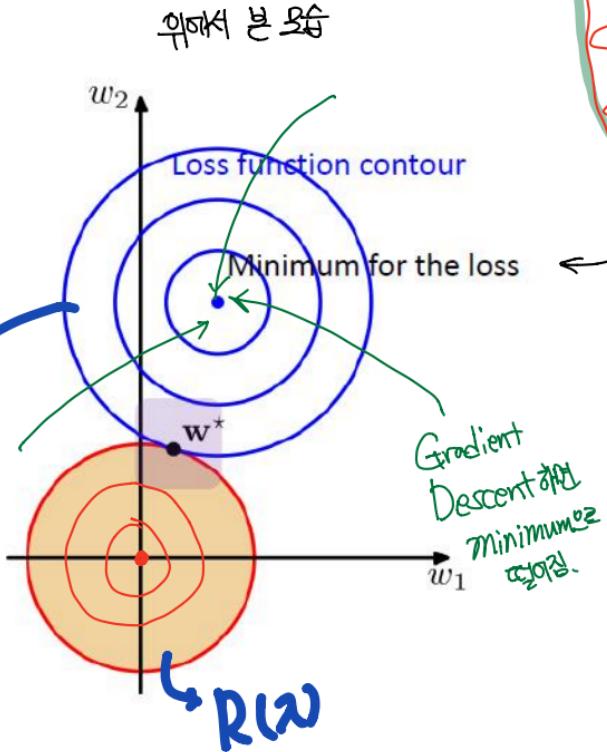
L2 Regularization

$$\text{Loss} = \text{MSE} + R(\lambda)$$

$\lambda \|w\|^2$
Constant

w (가중치)가 발산하지 않도록 힘을 주어
컨트롤 해줌.

MSE



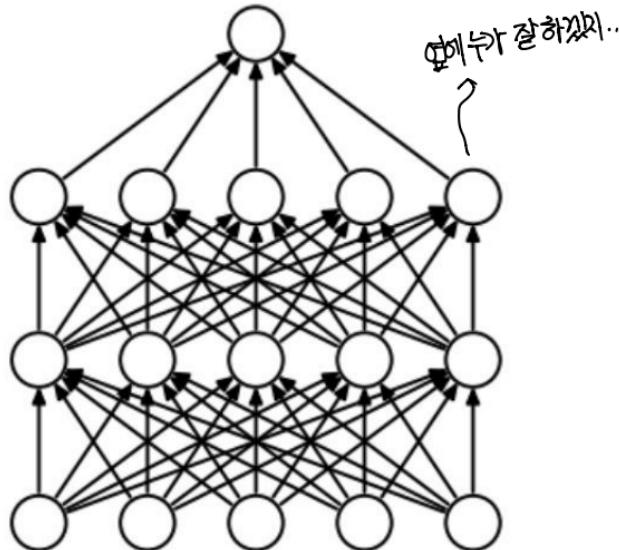
Dropout

↳ Training 할 때만 사용.

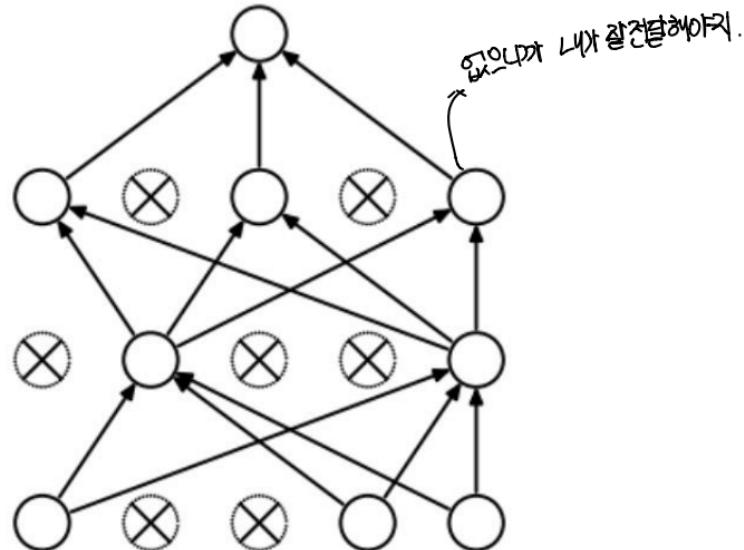
① 파라미터 줄이기.

② 간단한 여러 모델로 예측 후 합침.

⇒ Overfitting 줄이기.



(a) Standard Neural Net



(b) After applying dropout.

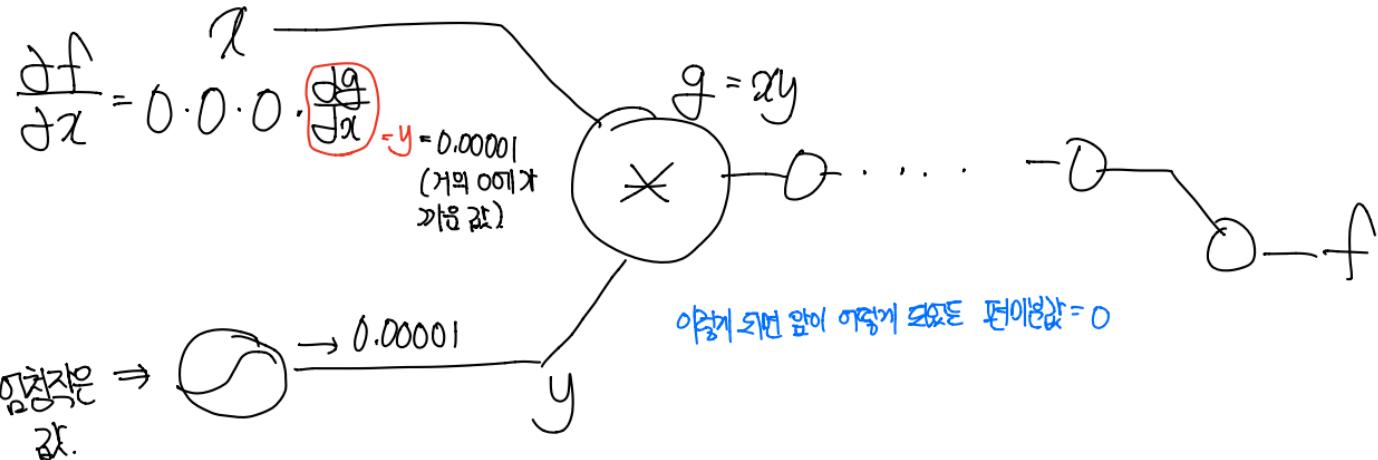
ex) 품질 관리에서
 → 내가 잘 했겠지~

 → 내가 잘 했지~

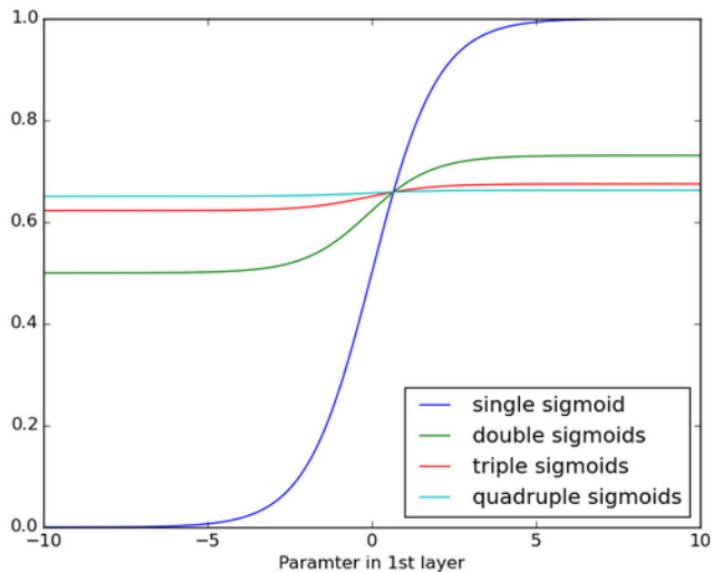
Problems of MLP

Gradient Vanishing

Gradient Vanishing

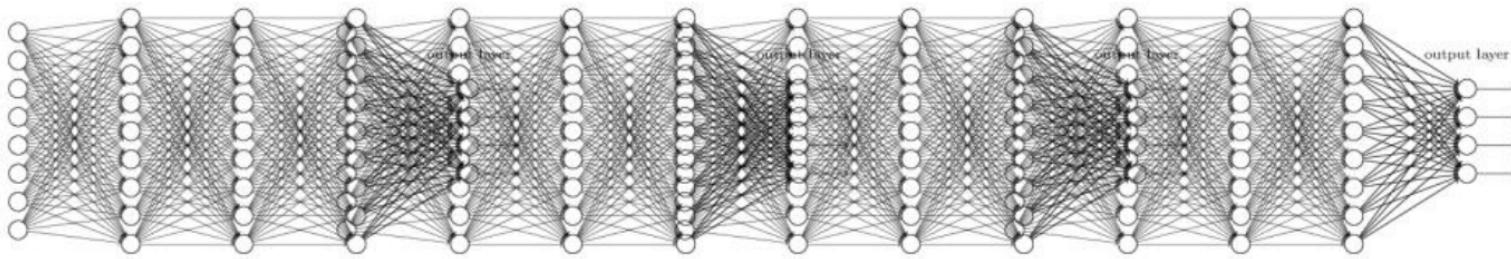


Gradient Vanishing

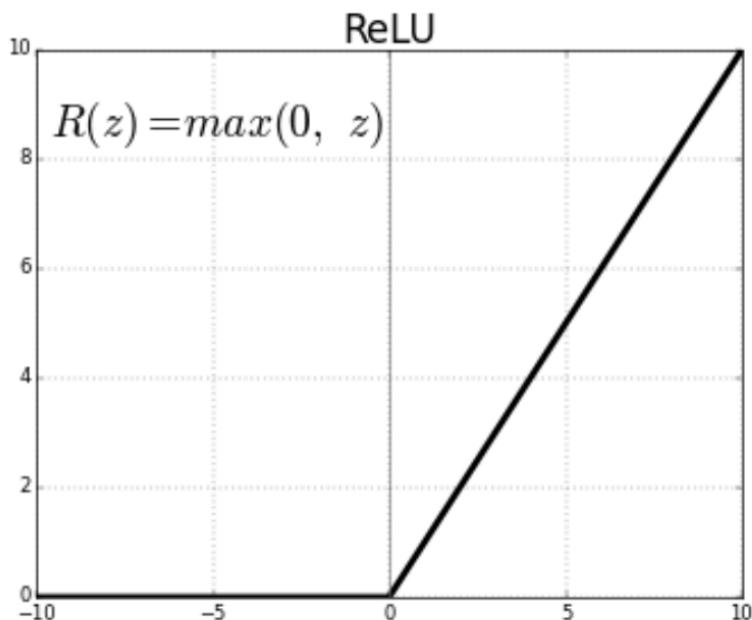


Sigmoid가 점점 절수록
‘0’에 가까운 값이 나옴

Gradient Vanishing



Rectified Linear Unit (ReLU) Activation

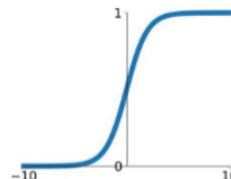


Review from Last Lecture

Activation Functions

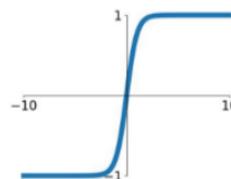
Sigmoid

$$\sigma(x) = \frac{1}{1+e^{-x}}$$



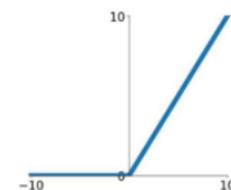
tanh

$$\tanh(x)$$



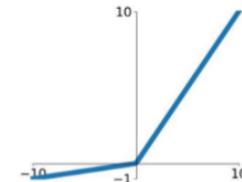
ReLU

$$\max(0, x)$$



Leaky ReLU

$$\max(0.1x, x)$$

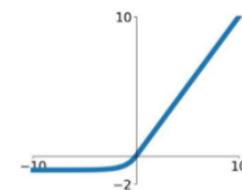


Maxout

$$\max(w_1^T x + b_1, w_2^T x + b_2)$$

ELU

$$\begin{cases} x & x \geq 0 \\ \alpha(e^x - 1) & x < 0 \end{cases}$$

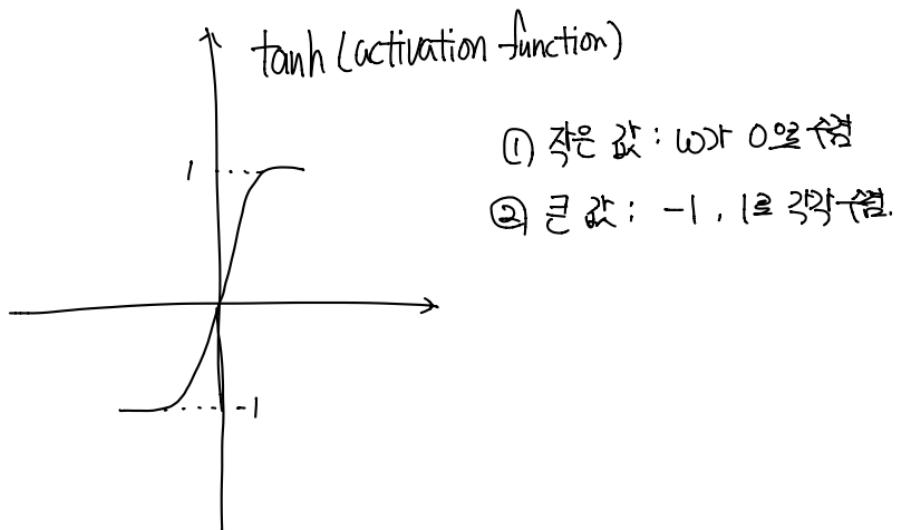


Other Techniques
Xavier Initialization and Batch Normalization

Xavier Initialization

만약 w 가 0으로 초기화 되었으면?

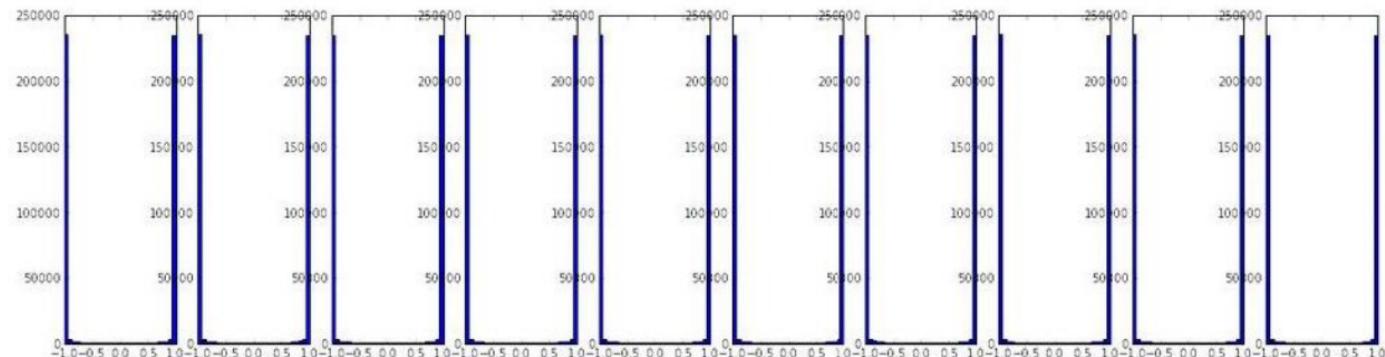
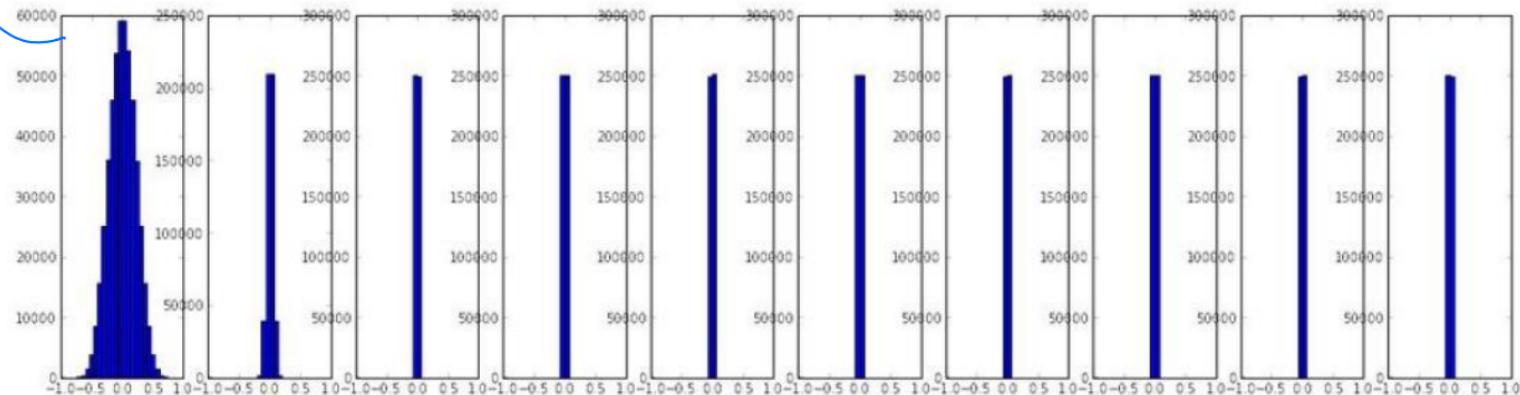
→ 2) Layer를 거쳐도 값이 0만 나와서 Loss 측정 불가능.



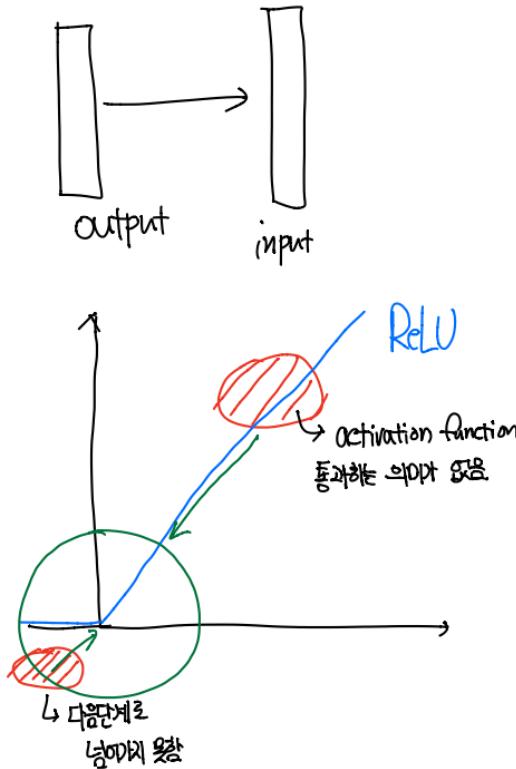
Xavier Initialization

tanh를 가장 흔적은 값이 되어 w_i 가 사라짐.(0으로 수렴함)

Standard gaussian $N(0, 1)$



Batch Normalization



Input: Values of x over a mini-batch: $\mathcal{B} = \{x_1 \dots m\}$;

Parameters to be learned: γ, β

Output: $\{y_i = \text{BN}_{\gamma, \beta}(x_i)\}$

$$\mu_{\mathcal{B}} \leftarrow \frac{1}{m} \sum_{i=1}^m x_i \quad // \text{mini-batch mean}$$

$$\sigma_{\mathcal{B}}^2 \leftarrow \frac{1}{m} \sum_{i=1}^m (x_i - \mu_{\mathcal{B}})^2 \quad // \text{mini-batch variance}$$

$$\hat{x}_i \leftarrow \frac{x_i - \mu_{\mathcal{B}}}{\sqrt{\sigma_{\mathcal{B}}^2 + \epsilon}} \quad // \text{normalize}$$

$$y_i \leftarrow \gamma \hat{x}_i + \beta \equiv \text{BN}_{\gamma, \beta}(x_i) \quad // \text{scale and shift}$$

Algorithm 1: Batch Normalizing Transform, applied to activation x over a mini-batch.