



# CS231n Lecture 6

- ☑ BOAZ 10기 박성현
- ☑ BOAZ 11기 김태희
- ☑ BOAZ 11기 홍지민

## 1. One time setup

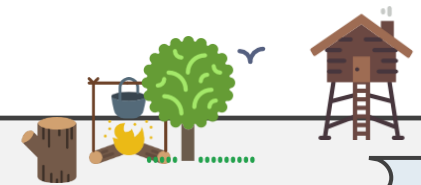
*activation functions, preprocessing, weight initialization, regularization, gradient checking*

## 2. Training dynamics

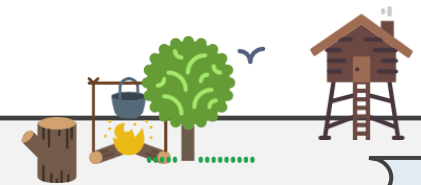
*babysitting the learning process, parameter updates, hyperparameter optimization*

## 3. Evaluation

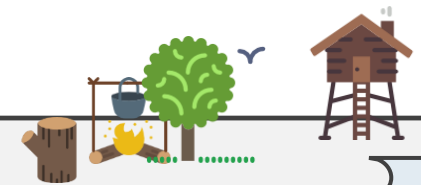
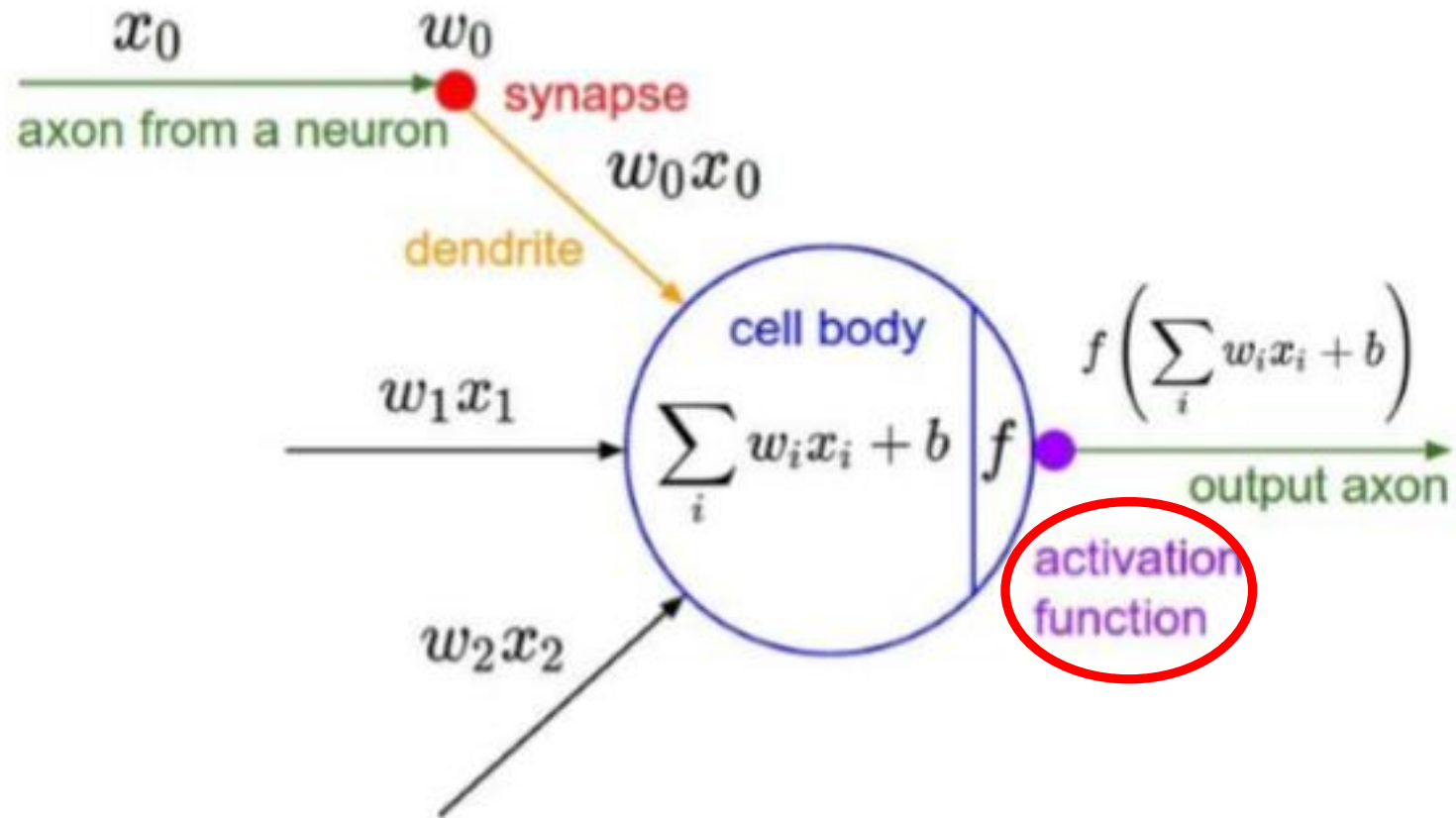
*model ensembles*



- Activation Functions
- Data Preprocessing
- Weight Initialization
- Batch Normalization
- Babysitting the Learning Process
- Hyperparameter Optimization

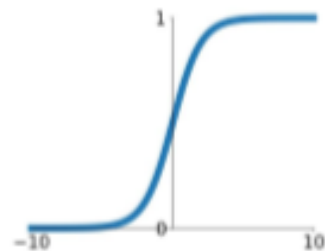


# Activation Functions



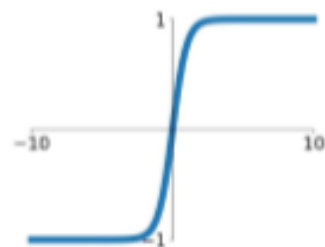
## Sigmoid

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



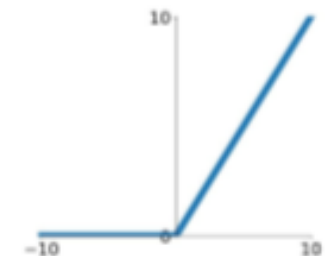
## tanh

$$\tanh(x)$$



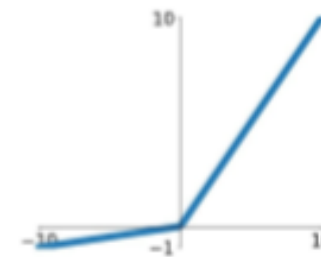
## ReLU

$$\max(0, x)$$



## Leaky ReLU

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

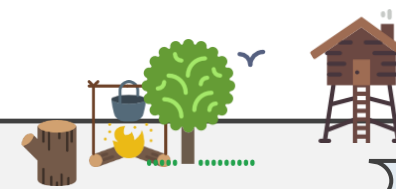
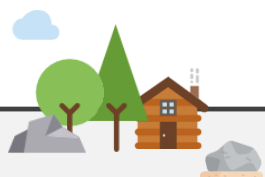
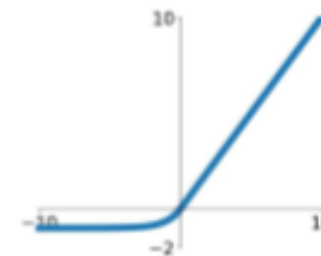


## Maxout

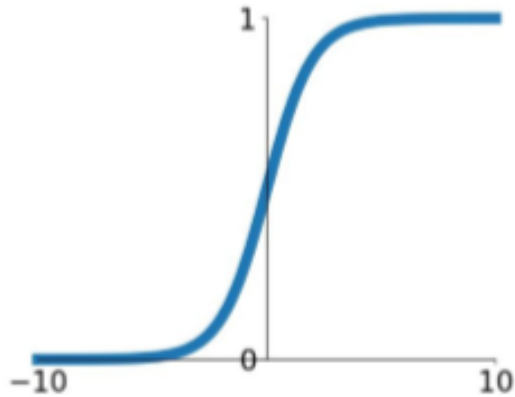
$$\max(w_1^T x + b_1, w_2^T x + b_2)$$

## ELU

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



# Activation Functions (Sigmoid)



**Sigmoid**

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

- Squashes numbers to range [0,1]
- Historically popular since they have nice interpretation as a saturating “firing rate” of a neuron

3 problems:

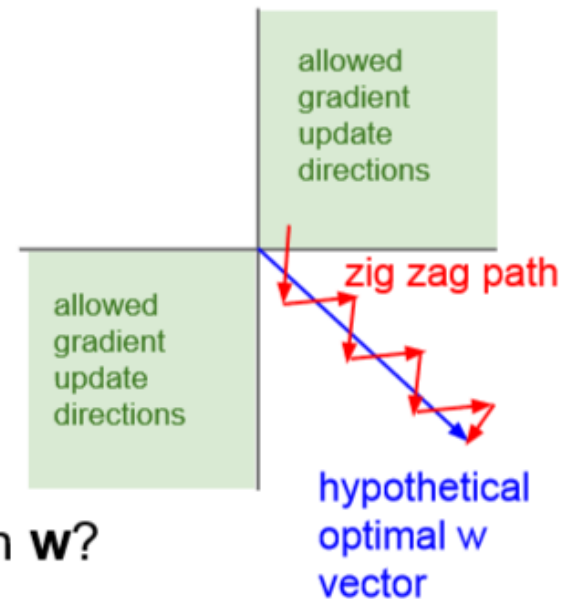
1. Saturated neurons “kill” the gradients
2. Sigmoid outputs are not zero-centered
3.  $\exp()$  is a bit compute expensive



# Activation Functions (Sigmoid)

Consider what happens when the input to a neuron is always positive...

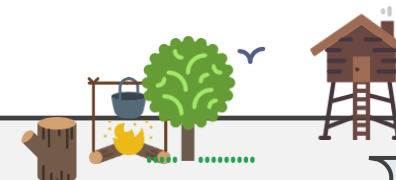
$$f\left(\sum_i w_i x_i + b\right)$$



What can we say about the gradients on **w**?

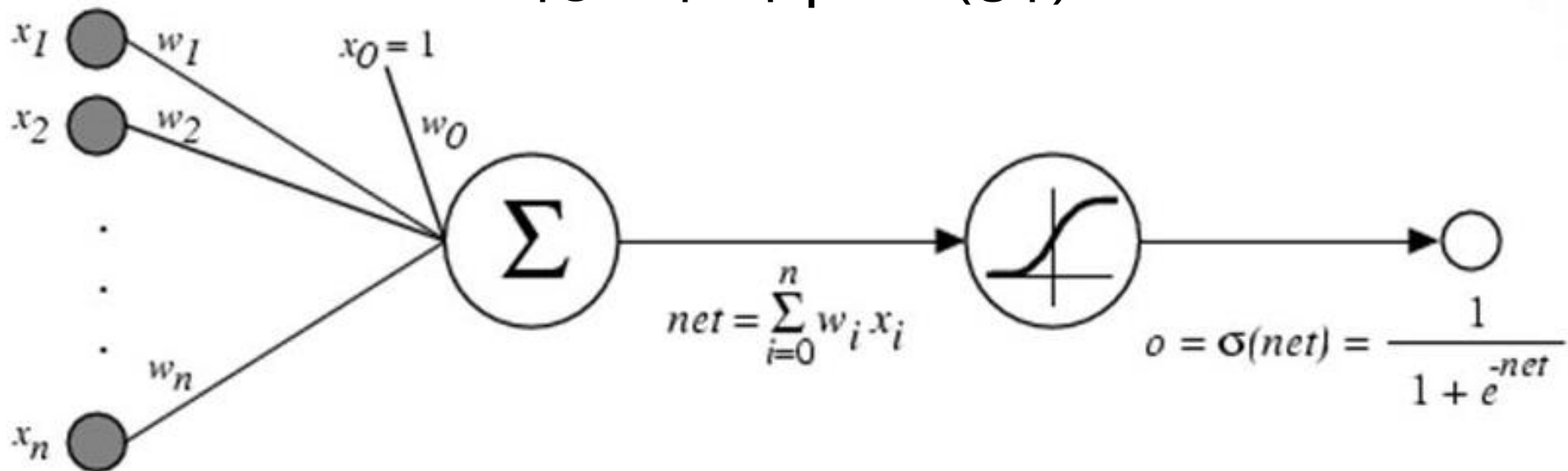
Always all positive or all negative :(  
(this is also why you want zero-mean data!)

Input이 모두 positive이고 sigmoid를 사용할 경우,  
W의 update가 모두 positive 또는 모두 negative 방향으로만 가능  
→ 해결방법 : X를 **zero-mean**으로 바꿔준다!



# Activation Functions (Sigmoid)

가정 :  $x$ 가 모두 positive(양수)



[W gradient]

Gradient > 0

Gradient < 0

[net gradient]

Gradient > 0

Gradient < 0

[upstream gradient]

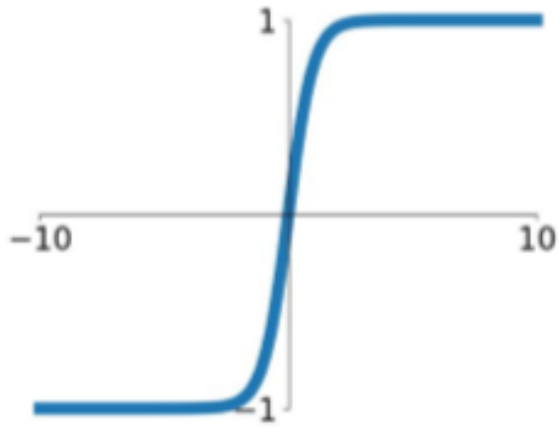
Gradient > 0

Gradient < 0





# Activation Functions (Tanh)

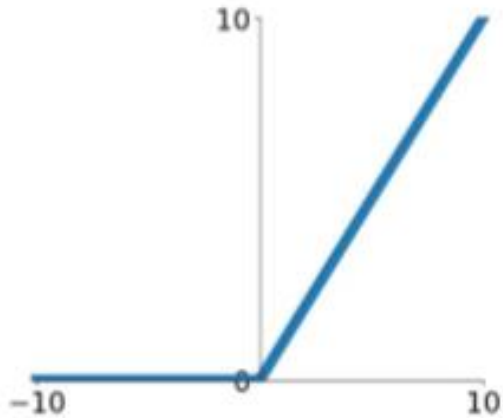


$\tanh(x)$

- Squashes numbers to range  $[-1,1]$
- zero centered (nice)
- still kills gradients when saturated :(

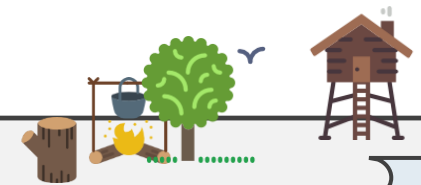


# Activation Functions (ReLU)

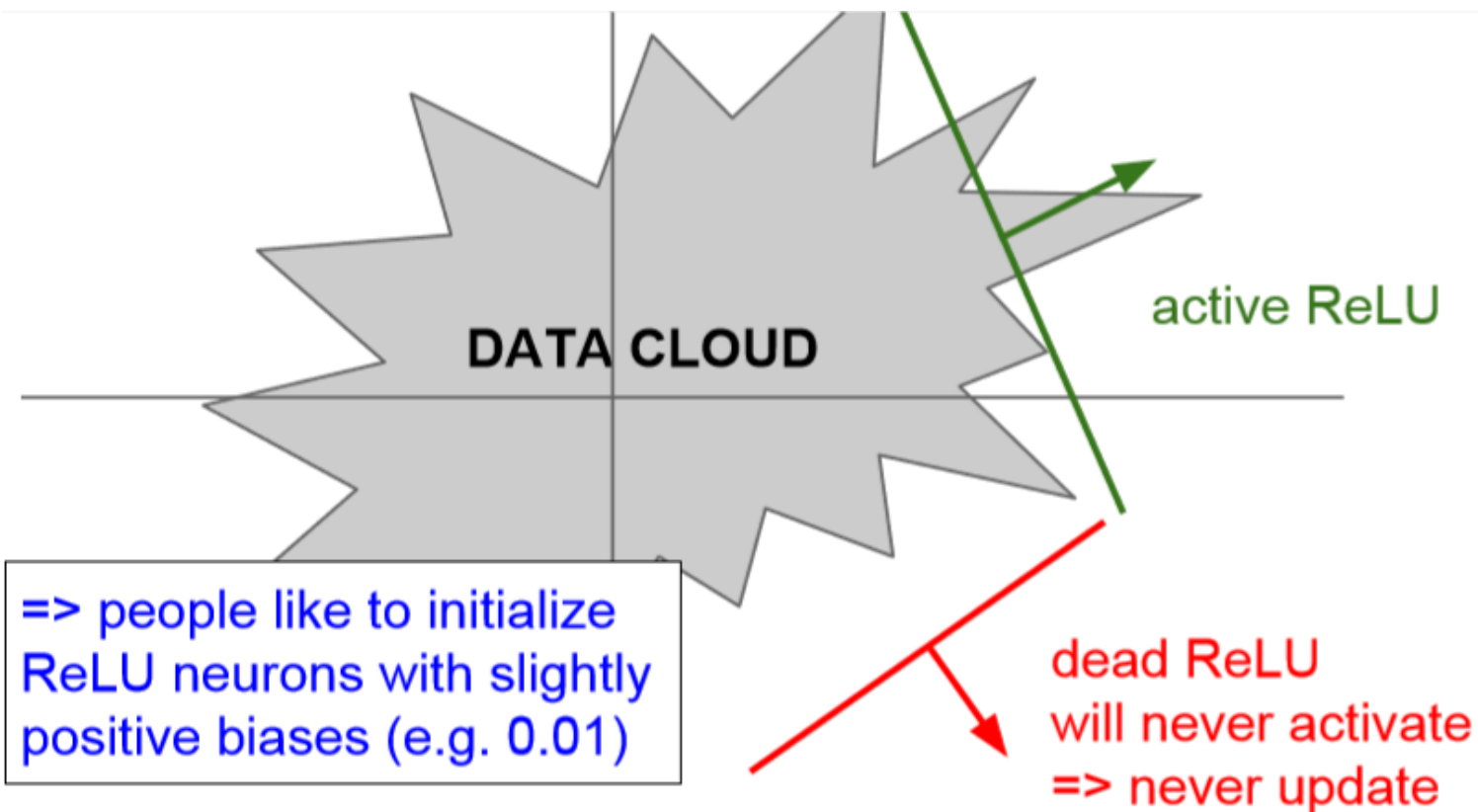


**ReLU**  
(Rectified Linear Unit)

- Computes  $f(x) = \max(0, x)$
- Does not saturate (in +region)
- Very computationally efficient
- Converges much faster than sigmoid/tanh in practice (e.g. 6x)
- Actually more biologically plausible than sigmoid
- Not zero-centered output
- An annoyance:



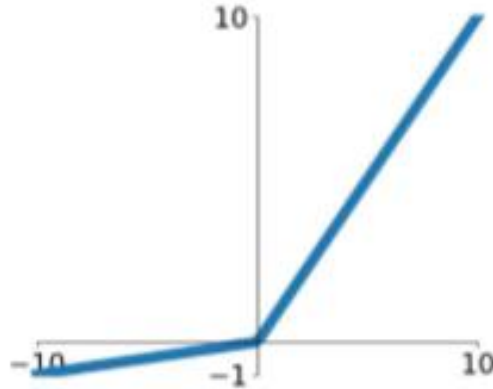
# Activation Functions (ReLU)



Input이 모두 negative가 될 경우, dead ReLU 발생  
dead ReLU가 발생하면, activate도 안되고, update도 되지 않음 (학습 X)



# Activation Functions (Leaky ReLU & PReLU)



**Leaky ReLU**

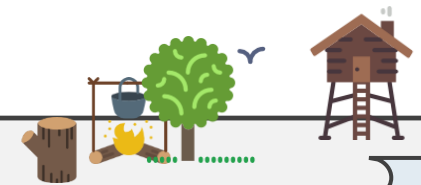
$$f(x) = \max(0.01x, x)$$

- Does not saturate
- Computationally efficient
- Converges much faster than sigmoid/tanh in practice! (e.g. 6x)
- **will not “die”.**

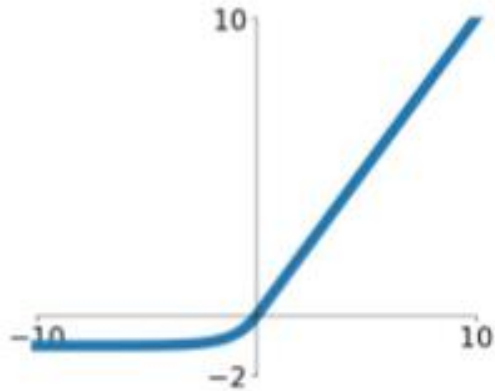
**Parametric Rectifier (PReLU)**

$$f(x) = \max(\alpha x, x)$$

backprop into  $\alpha$   
(parameter)



# Activation Functions (ELU)



$$f(x) = \begin{cases} x & \text{if } x > 0 \\ \alpha (\exp(x) - 1) & \text{if } x \leq 0 \end{cases}$$

- All benefits of ReLU
- Closer to zero mean outputs
- Negative saturation regime compared with Leaky ReLU adds some robustness to noise

- Computation requires exp()

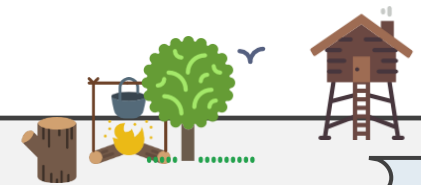


# Activation Functions (Maxout)

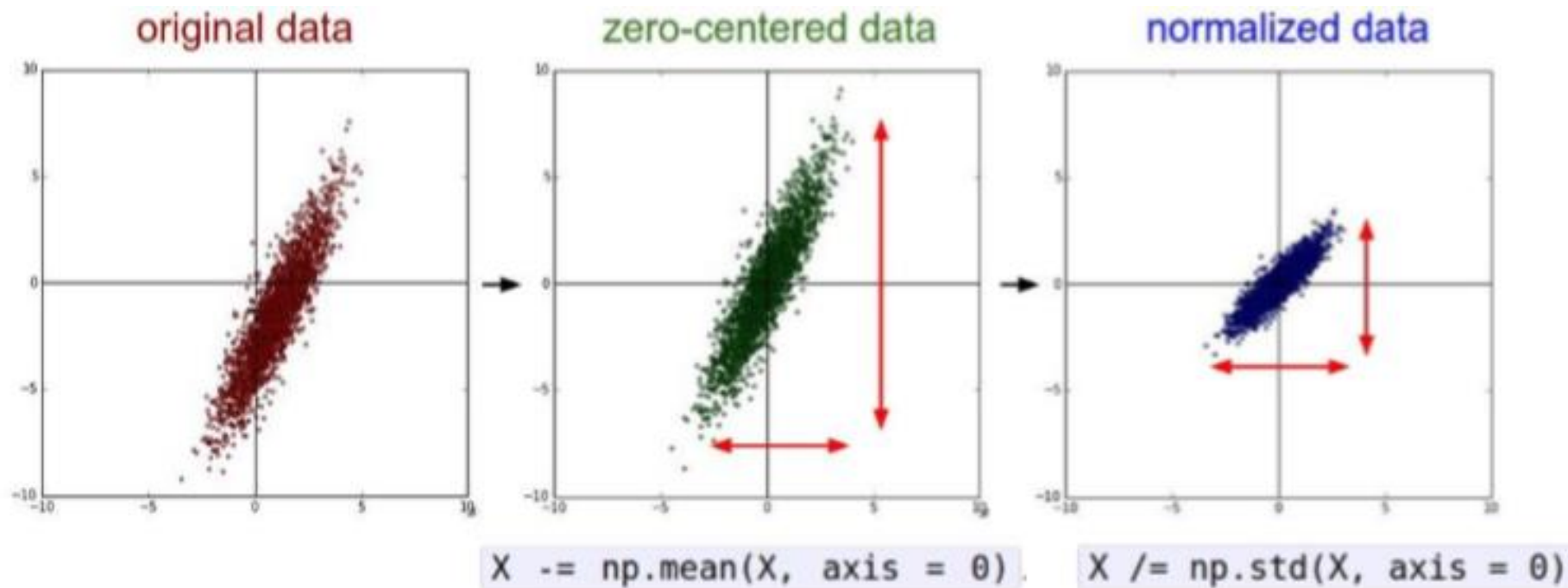
- Does not have the basic form of dot product -> nonlinearity
- Generalizes ReLU and Leaky ReLU
- Linear Regime! Does not saturate! Does not die!

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

Problem: doubles the number of parameters/neuron :(







PCA, Whitening 등의 방법을 사용 or Image의 mean을 빼주는 방법 사용

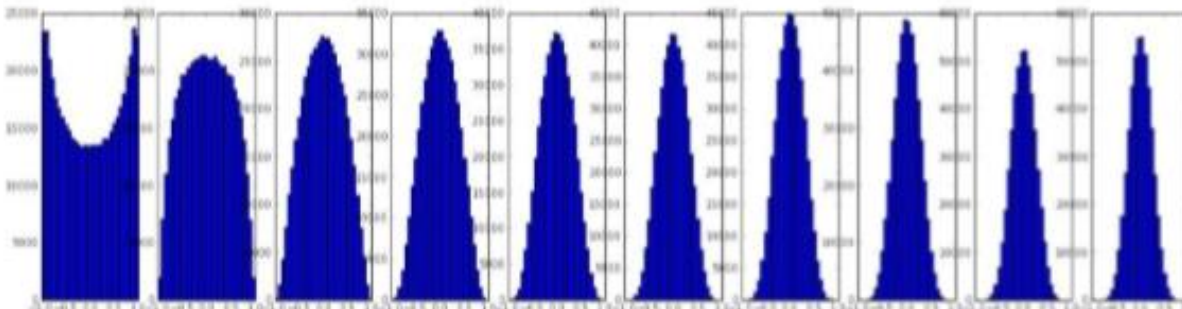
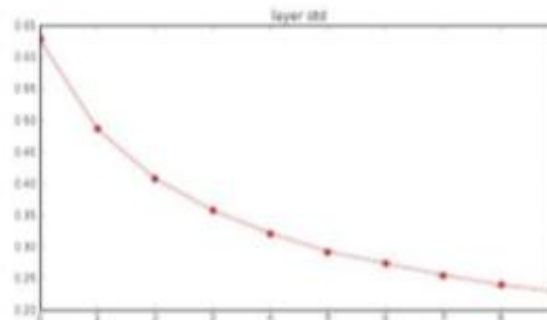
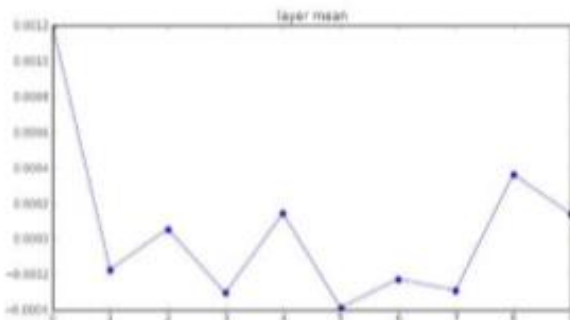


# Weight Initialization

input layer had mean 0.001880 and std 1.001311  
 hidden layer 1 had mean 0.001198 and std 0.627953  
 hidden layer 2 had mean -0.000175 and std 0.486051  
 hidden layer 3 had mean 0.000055 and std 0.407723  
 hidden layer 4 had mean -0.000306 and std 0.357108  
 hidden layer 5 had mean 0.000142 and std 0.320917  
 hidden layer 6 had mean -0.000389 and std 0.292116  
 hidden layer 7 had mean -0.000228 and std 0.273387  
 hidden layer 8 had mean -0.000291 and std 0.254935  
 hidden layer 9 had mean 0.000361 and std 0.239266  
 hidden layer 10 had mean 0.000139 and std 0.228008

```
W = np.random.randn(fan_in, fan_out) / np.sqrt(fan_in) # layer initialization
```

“Xavier initialization”  
 [Glorot et al., 2010]



**Reasonable initialization.**  
 (Mathematical derivation  
 assumes linear activations)

- fan\_in : input 노드의 개수
- fan\_out : output 노드의 개수



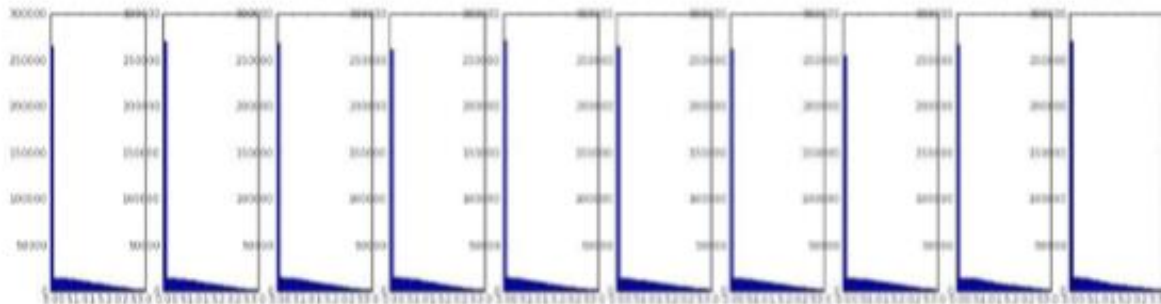
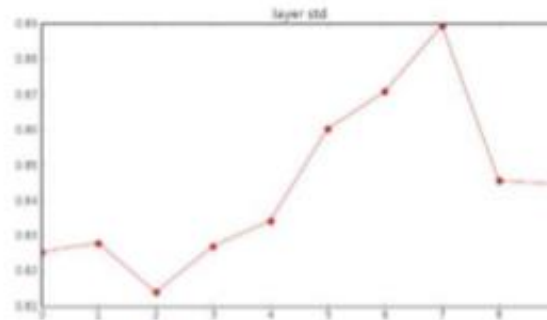
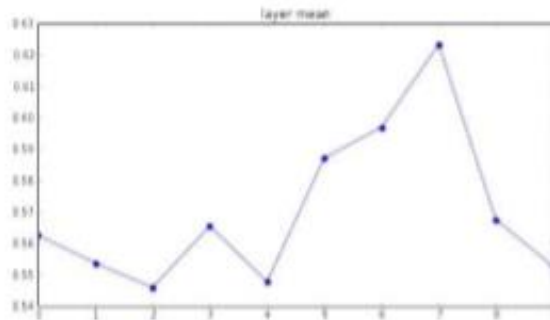


# Weight Initialization

input layer had mean 0.000501 and std 0.999444  
 hidden layer 1 had mean 0.562488 and std 0.825232  
 hidden layer 2 had mean 0.553614 and std 0.827835  
 hidden layer 3 had mean 0.545867 and std 0.813855  
 hidden layer 4 had mean 0.565396 and std 0.826982  
 hidden layer 5 had mean 0.547678 and std 0.834692  
 hidden layer 6 had mean 0.587183 and std 0.860635  
 hidden layer 7 had mean 0.596867 and std 0.870610  
 hidden layer 8 had mean 0.623214 and std 0.889348  
 hidden layer 9 had mean 0.567498 and std 0.845357  
 hidden layer 10 had mean 0.552531 and std 0.844523

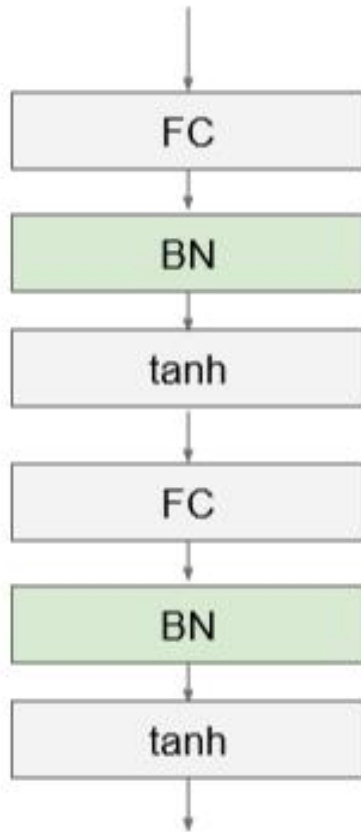
```
W = np.random.randn(fan_in, fan_out) / np.sqrt(fan_in/2) # layer initialization
```

He et al., 2015  
 (note additional /2)



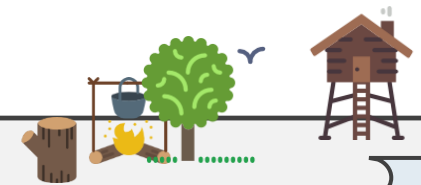
- fan\_in : input 노드의 개수
- fan\_out : output 노드의 개수





Usually inserted after Fully Connected or Convolutional layers, and before nonlinearity.

$$\hat{x}^{(k)} = \frac{x^{(k)} - E[x^{(k)}]}{\sqrt{\text{Var}[x^{(k)}]}}$$



# Batch Normalization

Normalize:

$$\hat{x}^{(k)} = \frac{x^{(k)} - \mathbb{E}[x^{(k)}]}{\sqrt{\text{Var}[x^{(k)}]}}$$

And then allow the network to squash the range if it wants to:

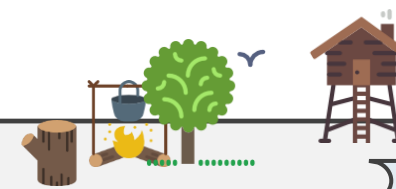
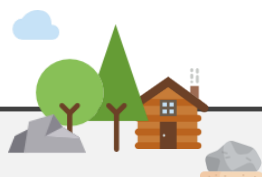
$$y^{(k)} = \gamma^{(k)} \hat{x}^{(k)} + \beta^{(k)}$$

Note, the network can learn:

$$\gamma^{(k)} = \sqrt{\text{Var}[x^{(k)}]}$$

$$\beta^{(k)} = \mathbb{E}[x^{(k)}]$$

to recover the identity mapping.



# Batch Normalization

**Input:** Values of  $x$  over a mini-batch:  $\mathcal{B} = \{x_{1\dots m}\}$ ;

Parameters to be learned:  $\gamma, \beta$

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

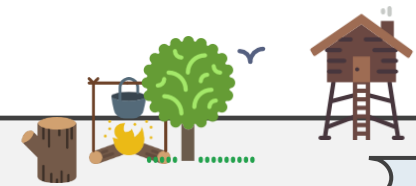
$$\mu_{\mathcal{B}} \leftarrow \frac{1}{m} \sum_{i=1}^m x_i \quad // \text{ mini-batch mean}$$

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

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

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

- Improves gradient flow through the network
- Allows higher learning rates
- Reduces the strong dependence on initialization
- Acts as a form of regularization in a funny way, and slightly reduces the need for dropout, maybe

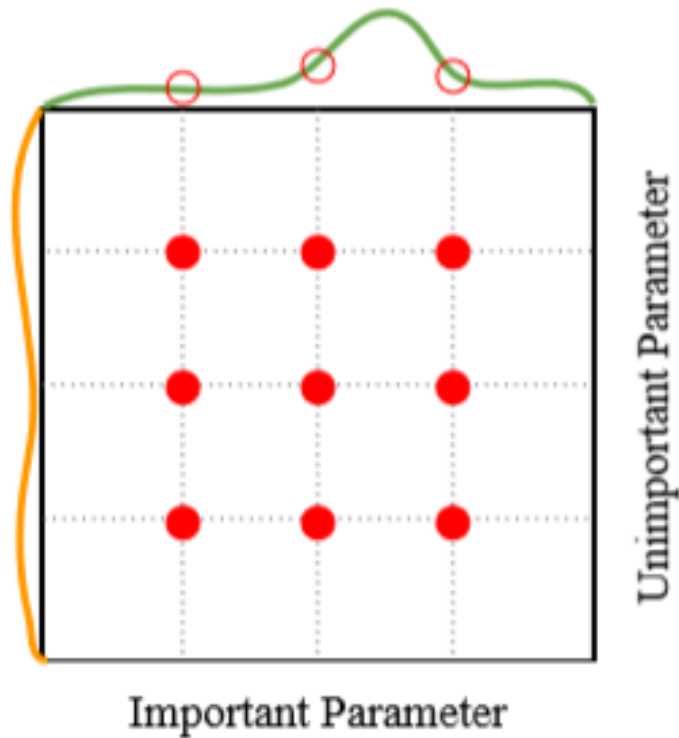


1. Preprocess the data
2. Choose the architecture
3. Overfitting이 되지 않도록, Validation set에 대해서도 Loss, Acc 확인
4. Learning rate를 적절하게 선택
  - Loss not going down : Learning rate too low
  - Loss exploding : Learning rate too high

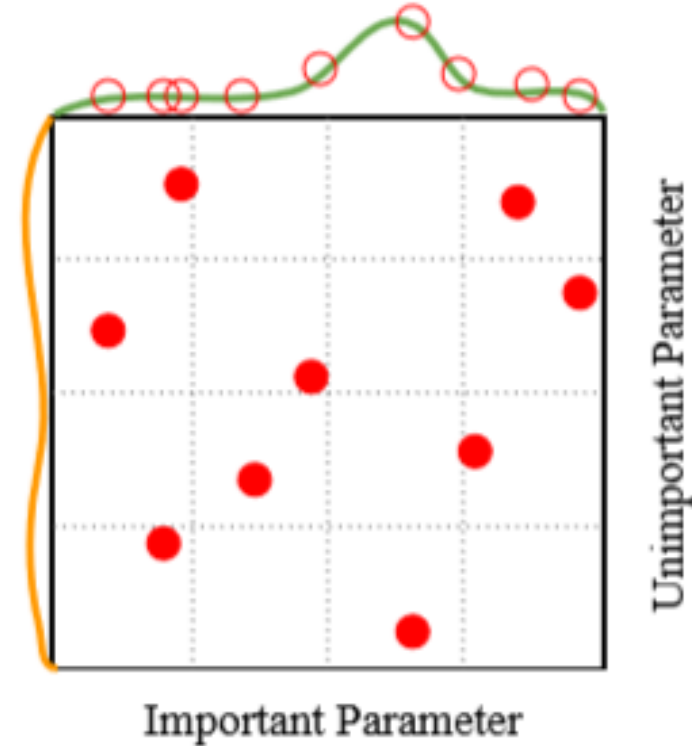


# Hyperparameter Optimization

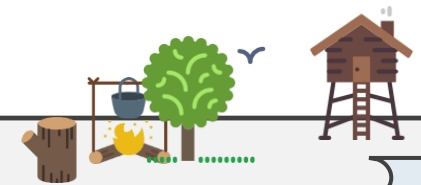
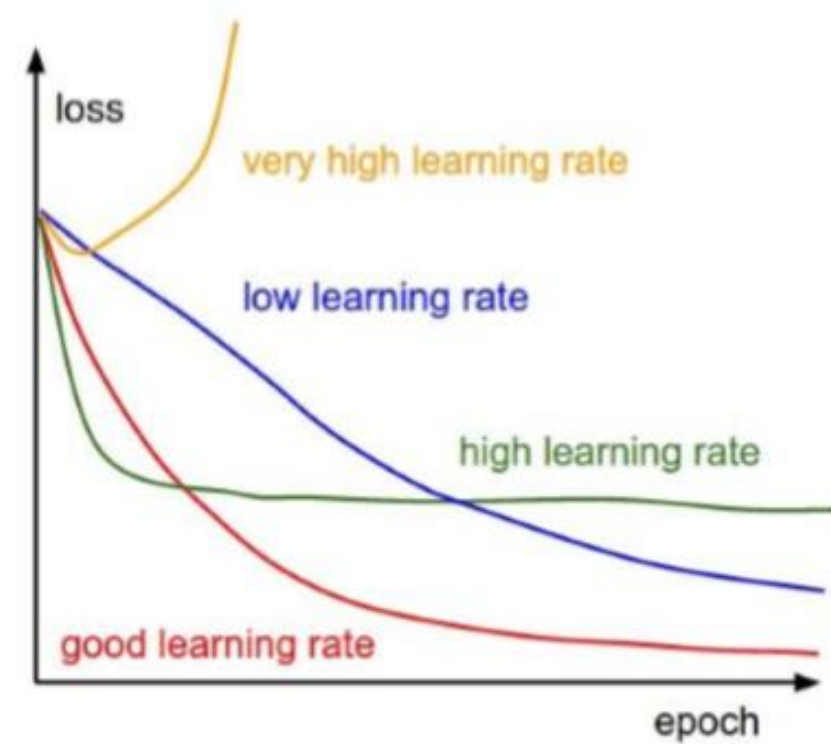
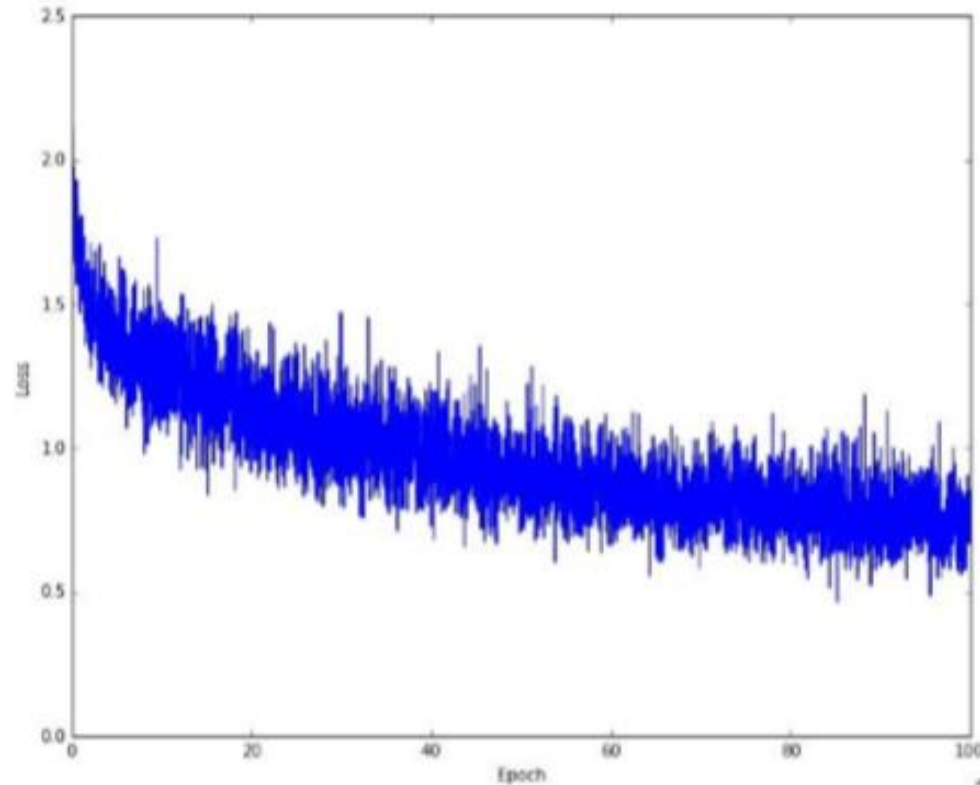
Grid Layout



Random Layout

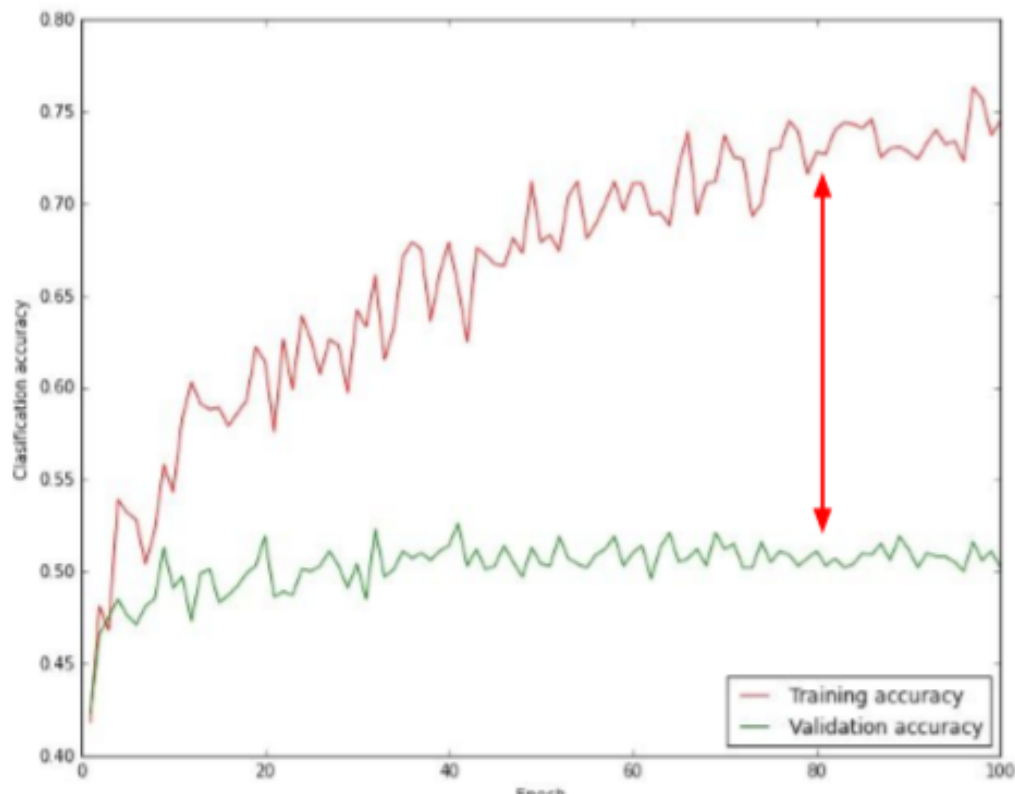


Monitor and visualize the loss curve





Monitor and visualize the accuracy:



**big gap = overfitting**

**=> increase regularization strength?**

**no gap**

**=> increase model capacity?**





CS231n : <http://cs231n.stanford.edu/syllabus.html>

