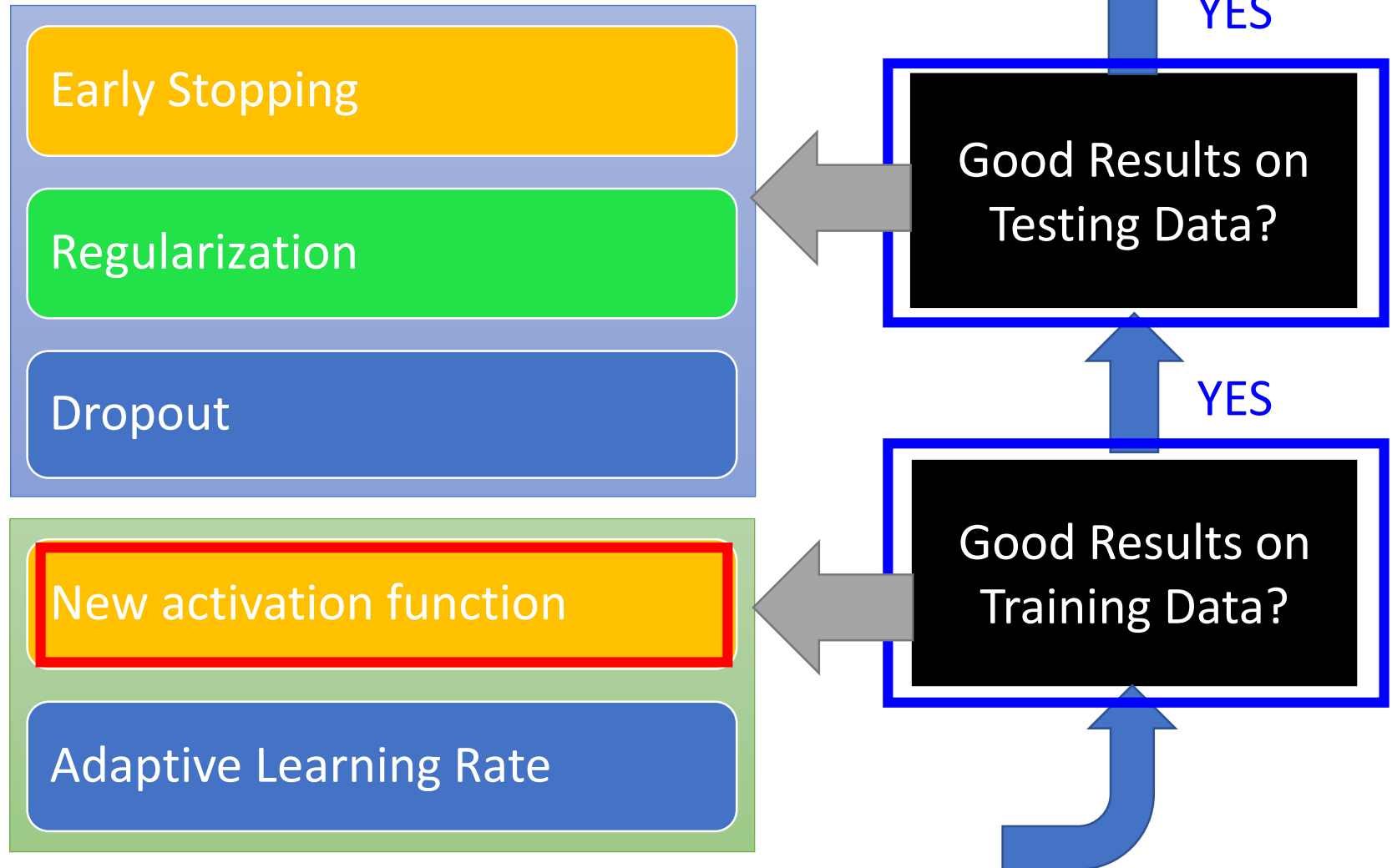
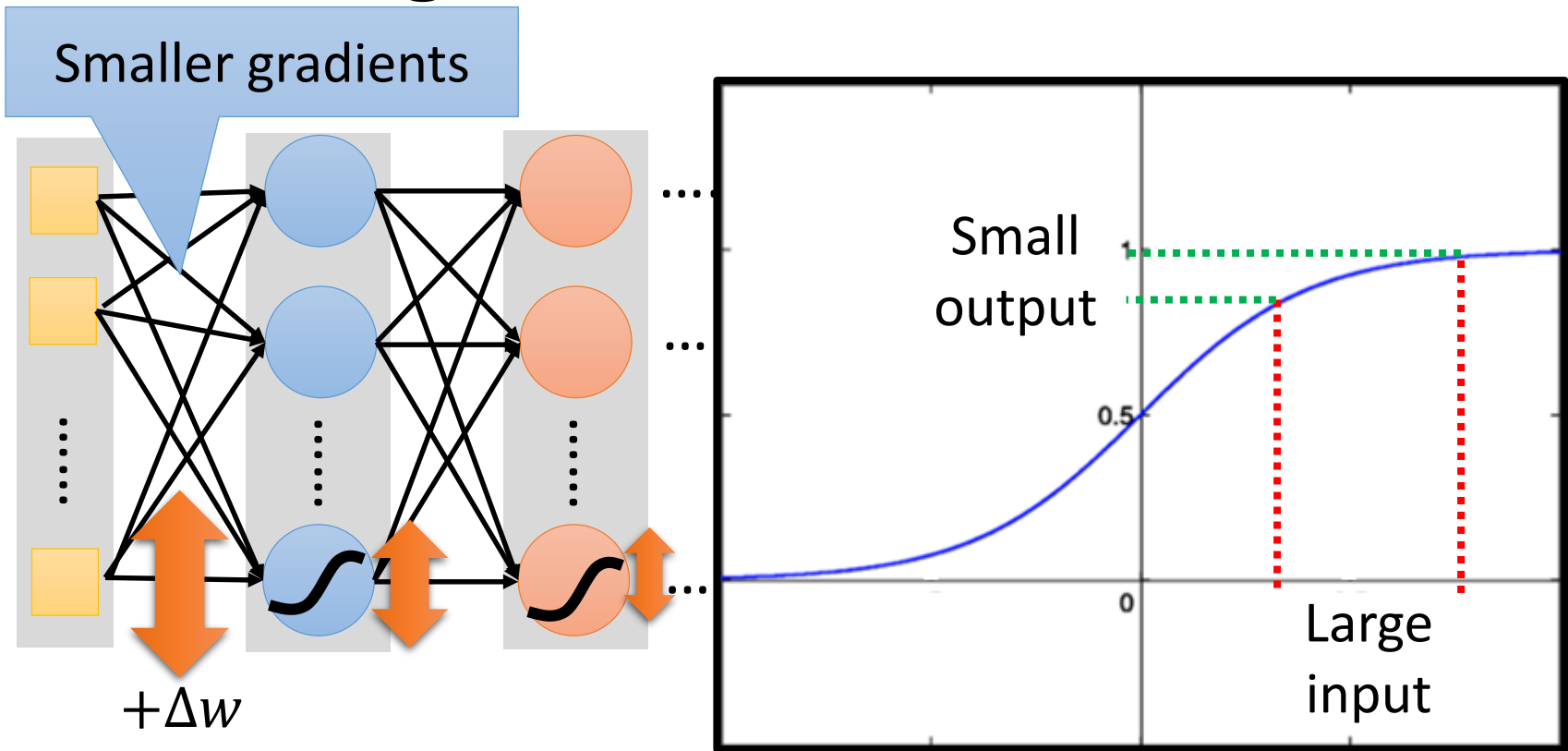


# Tips for Deep Learning

# Recipe of Deep Learning



# Vanishing Gradient Problem

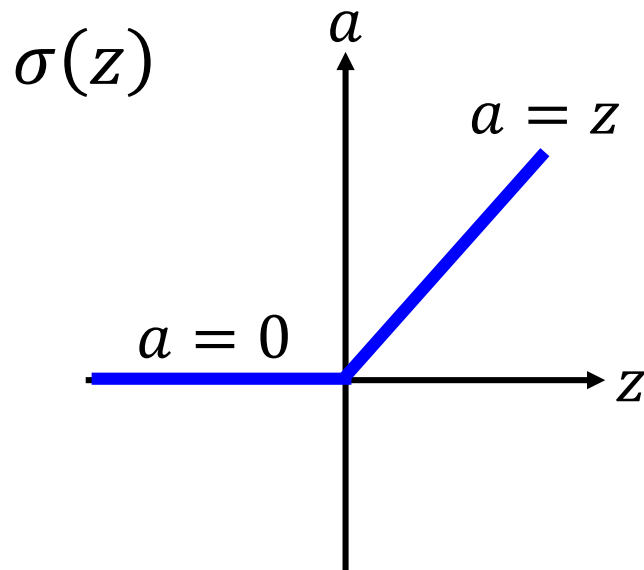


Intuitive way to compute the derivatives ...

$$\frac{\partial l}{\partial w} = ? \quad \frac{\Delta l}{\Delta w}$$

# ReLU

- Rectified Linear Unit (ReLU)



[Xavier Glorot, AISTATS'11]  
[Andrew L. Maas, ICML'13]  
[Kaiming He, arXiv'15]

## Reason:

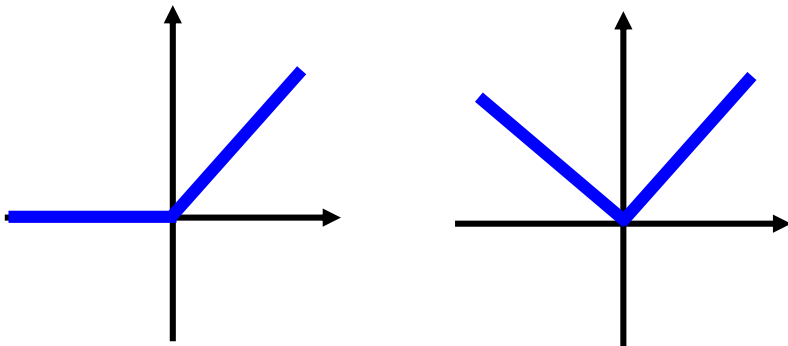
1. Fast to compute
2. Biological reason
3. Infinite sigmoid with different biases
4. Vanishing gradient problem

# Maxout

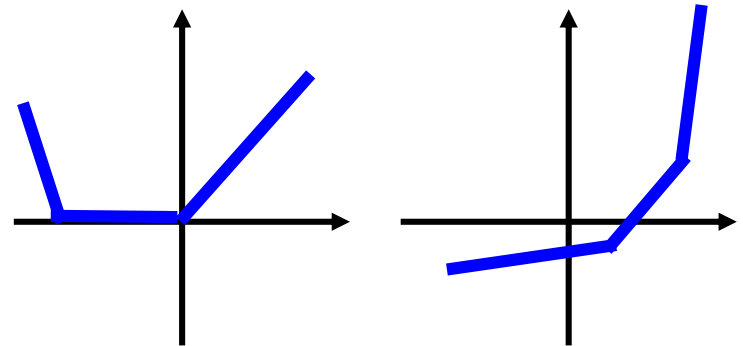
ReLU is a special cases of Maxout

- Learnable activation function [\[Ian J. Goodfellow, ICML'13\]](#)
  - Activation function in maxout network can be any piecewise linear convex function
  - How many pieces depending on how many elements in a group

2 elements in a group

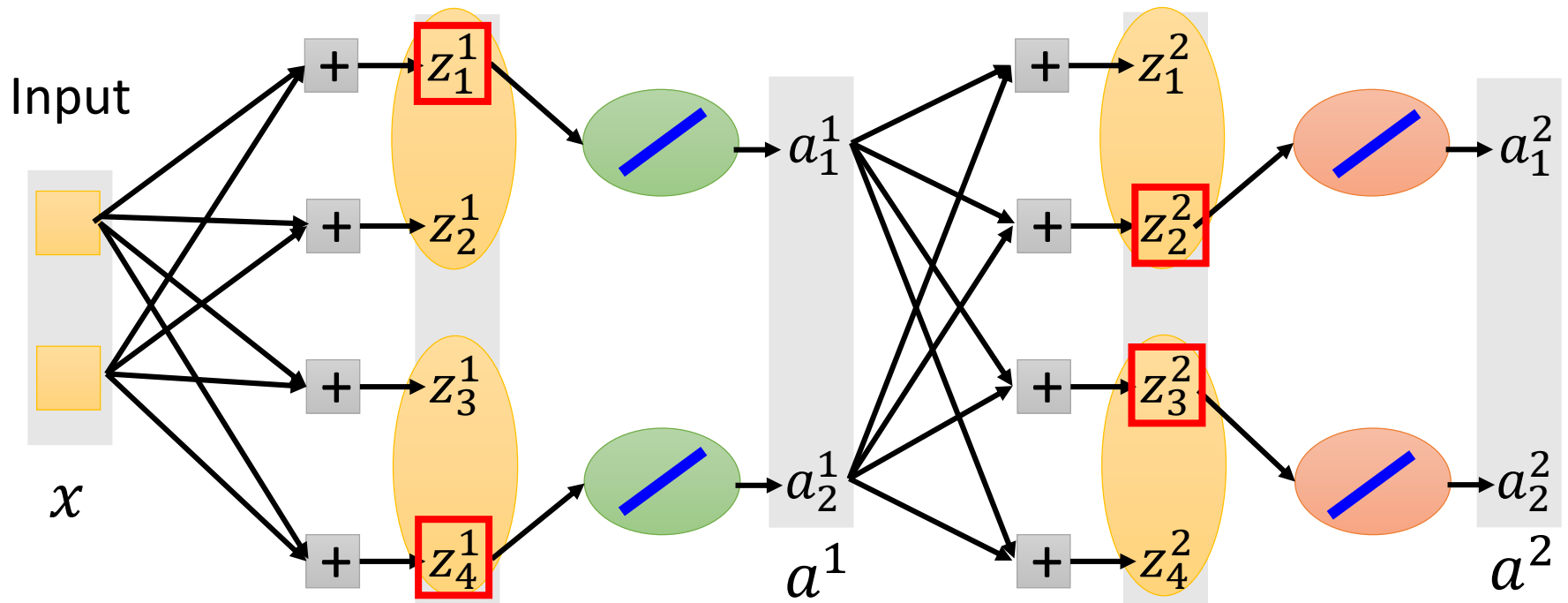


3 elements in a group



# Maxout - Training

- Given a training data  $x$ , we know which  $z$  would be the max



- Train this thin and linear network

Different thin and linear network for different examples

# Adam

## RMSProp + Momentum

### RMSProp

Error Surface can be very complex when training NN.

$$w^1 \leftarrow w^0 - \frac{\eta}{\sigma^0} g^0 \quad \sigma^0 = g^0$$

$$w^2 \leftarrow w^1 - \frac{\eta}{\sigma^1} g^1 \quad \sigma^1 = \sqrt{\alpha(\sigma^0)^2 + (1 - \alpha)(g^1)^2}$$

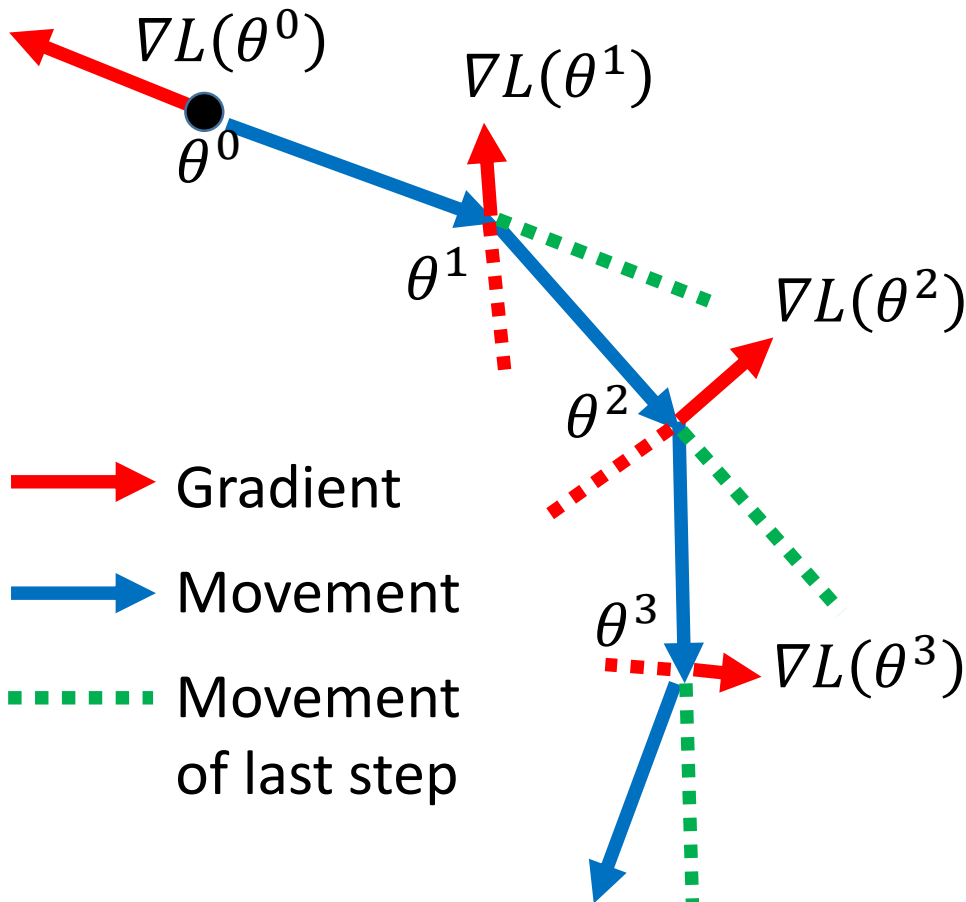
$$w^3 \leftarrow w^2 - \frac{\eta}{\sigma^2} g^2 \quad \sigma^2 = \sqrt{\alpha(\sigma^1)^2 + (1 - \alpha)(g^2)^2}$$

⋮

$$w^{t+1} \leftarrow w^t - \frac{\eta}{\sigma^t} g^t \quad \sigma^t = \sqrt{\alpha(\sigma^{t-1})^2 + (1 - \alpha)(g^t)^2}$$

# Momentum

Movement: movement of last step minus gradient at present



Start at point  $\theta^0$

Movement  $v^0=0$

Compute gradient at  $\theta^0$

Movement  $v^1 = \lambda v^0 - \eta \nabla L(\theta^0)$

Move to  $\theta^1 = \theta^0 + v^1$

Compute gradient at  $\theta^1$

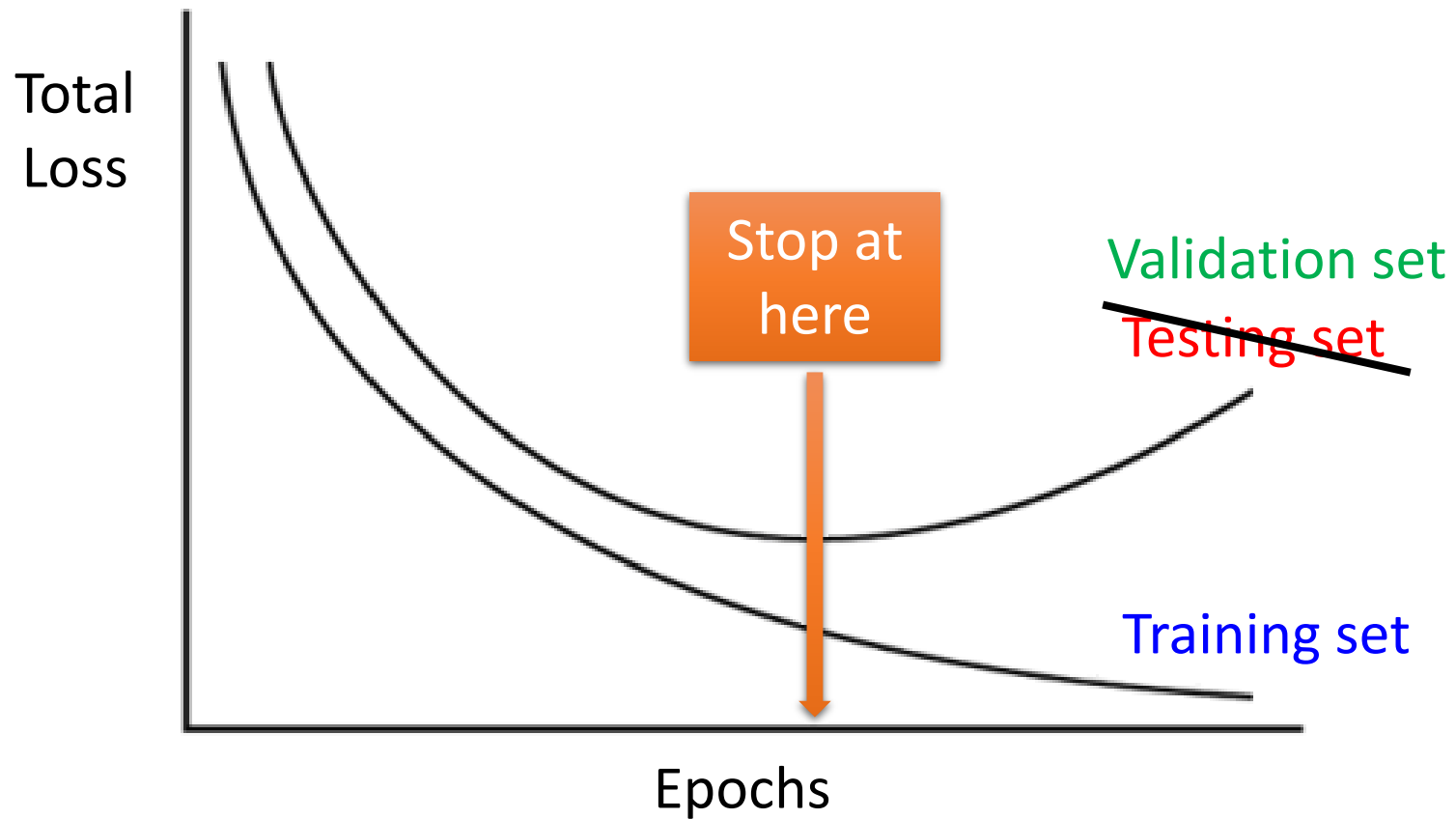
Movement  $v^2 = \lambda v^1 - \eta \nabla L(\theta^1)$

Move to  $\theta^2 = \theta^1 + v^2$

Movement not just based on gradient, but previous movement.



# Early Stopping



Keras: <http://keras.io/getting-started/faq/#how-can-i-interrupt-training-when-the-validation-loss-isnt-decreasing-anymore>

# Regularization

- New loss function to be minimized
  - Find a set of weight not only minimizing original cost but also close to zero

$$L'(\theta) = \underbrace{L(\theta)} + \lambda \frac{1}{2} \underbrace{\|\theta\|_2^2} \rightarrow \text{Regularization term}$$

$$\theta = \{w_1, w_2, \dots\}$$

Original loss

(e.g. minimize square error, cross entropy ...)

L2 regularization:

$$\|\theta\|_2^2 = (w_1)^2 + (w_2)^2 + \dots$$

(usually not consider biases)

# Regularization

L2 regularization:

$$\|\theta\|_2 = (w_1)^2 + (w_2)^2 + \dots$$

- New loss function to be minimized

$$L'(\theta) = L(\theta) + \lambda \frac{1}{2} \|\theta\|_2^2 \quad \text{Gradient: } \frac{\partial L'}{\partial w} = \frac{\partial L}{\partial w} + \lambda w$$

Update:

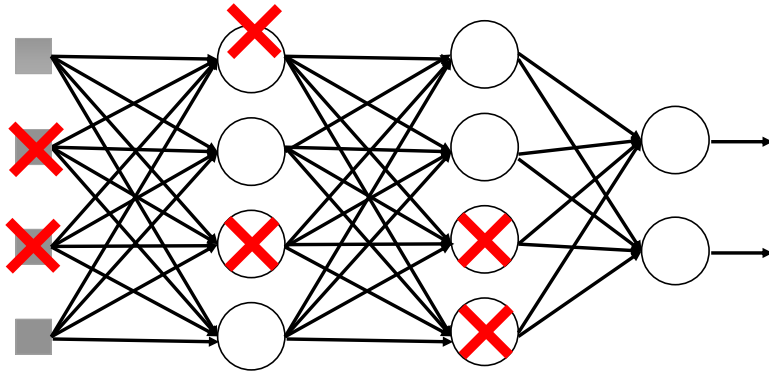
$$= w^t - \eta \left( \frac{\partial L}{\partial w} + \lambda w^t \right)$$

Weight Decay

—  
↓  
Closer to zero

# Dropout

## Training:



### ➤ Each time before updating the parameters

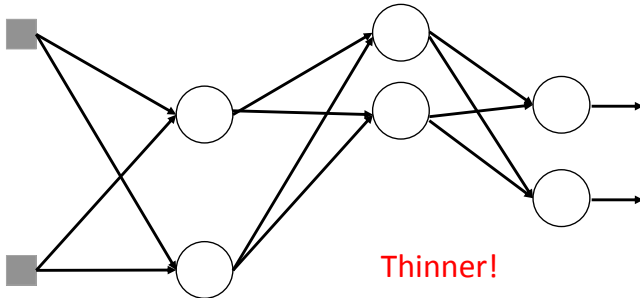
- Each neuron has  $p\%$  to dropout



**The structure of the network is changed.**

- Using the new network for training

For each mini-batch, we resample the dropout neurons



## Testing:

### ➤ No dropout

- If the dropout rate at training is  $p\%$ , all the weights times  $1-p\%$
- Assume that the dropout rate is 50%. If a weight  $w = 1$  by training, set  $w = 0.5$  for testing.

