# My Last Lecture
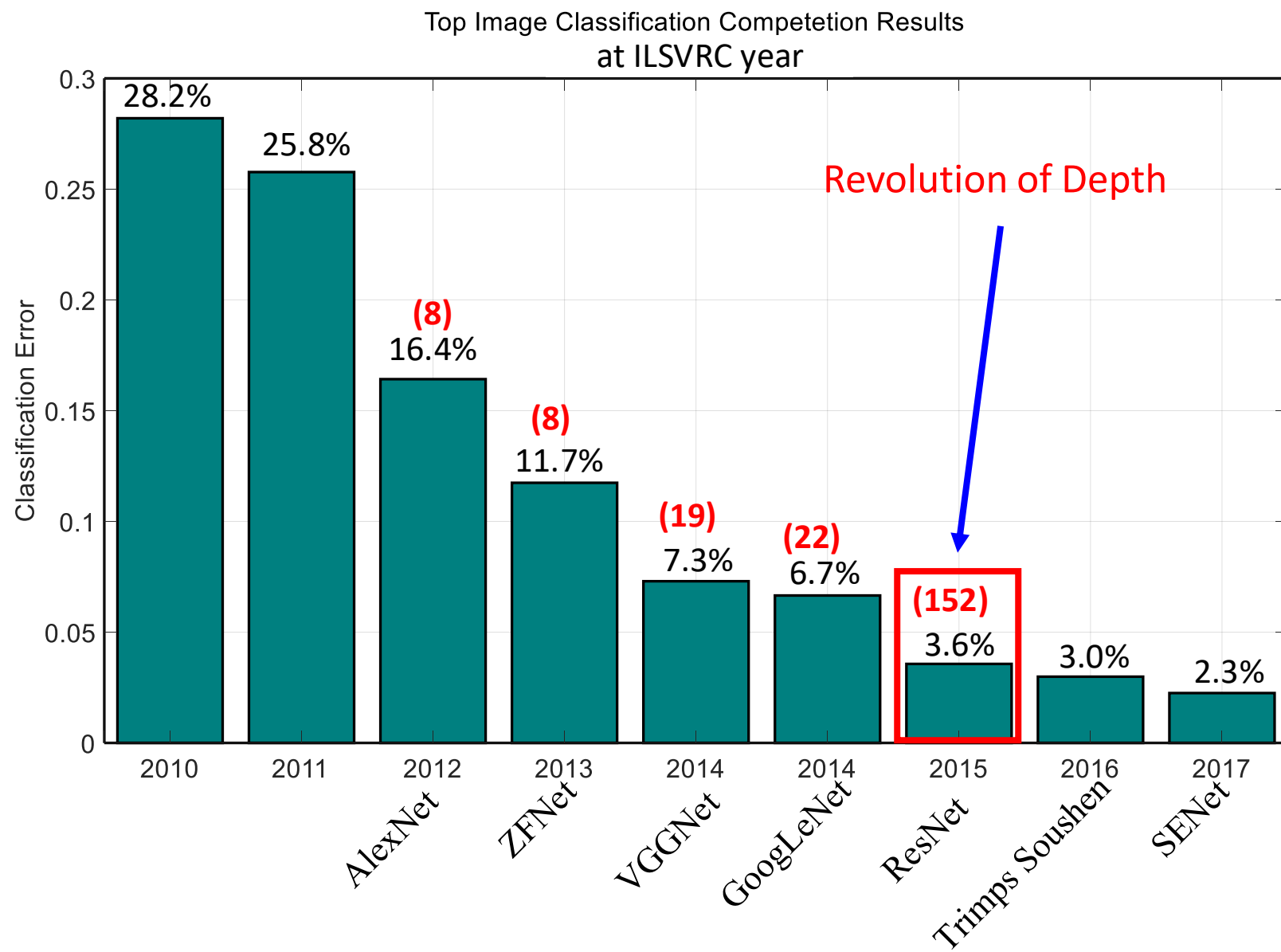# Typical CNN Architectures

Case Studies
- AlexNet
- VGG
- GoogLeNet
- ResNet

about other architecture…

- NIN (Network in Network)
- Wide ResNet
- ResNeXT
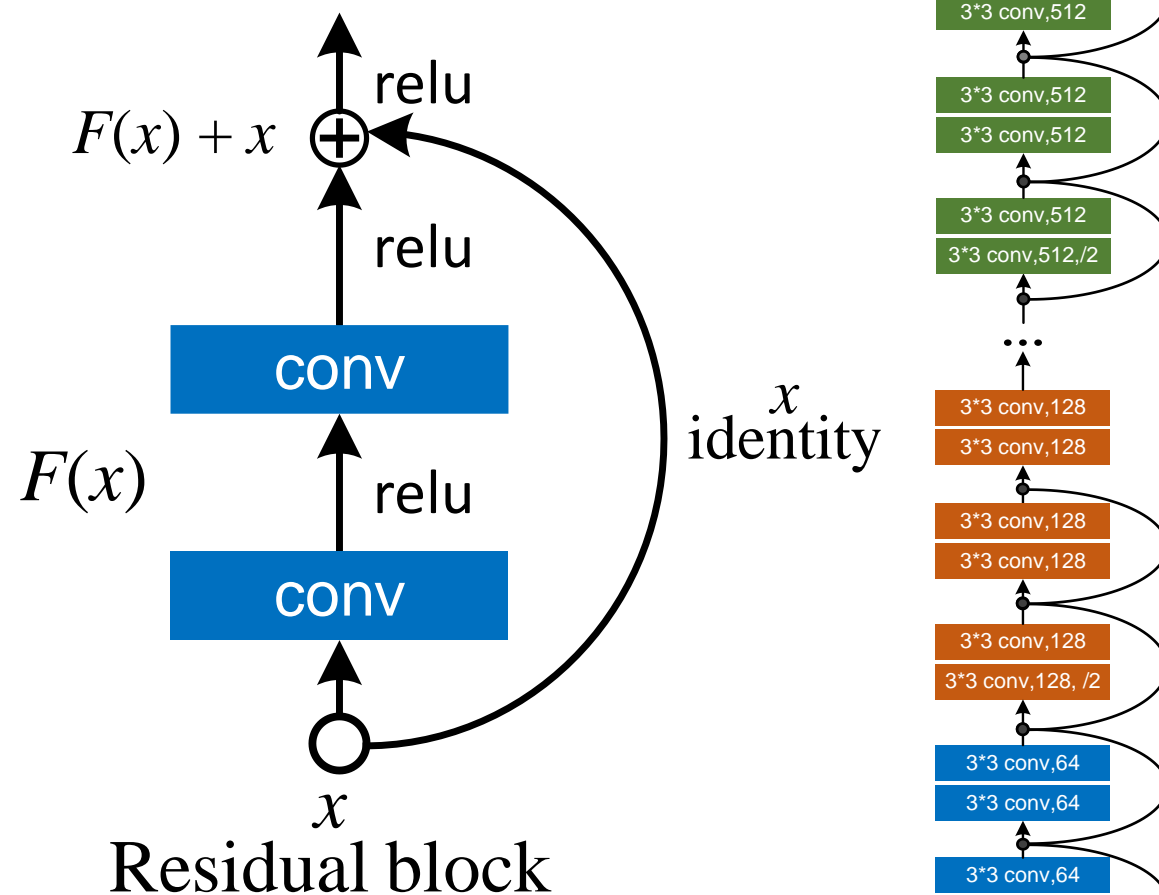- Stochastic Depth

- DenseNet
- SENet
- FractalNet
- SqueezeNet

Top Image Classification Competetion Results
at ILSVRC year

Revolution of Depth

Li Liu et al., Deep Learning for Generic Object Detection: A Survey, IJCV, 2019.

# Case Study: ResNet

[He et al., 2015]

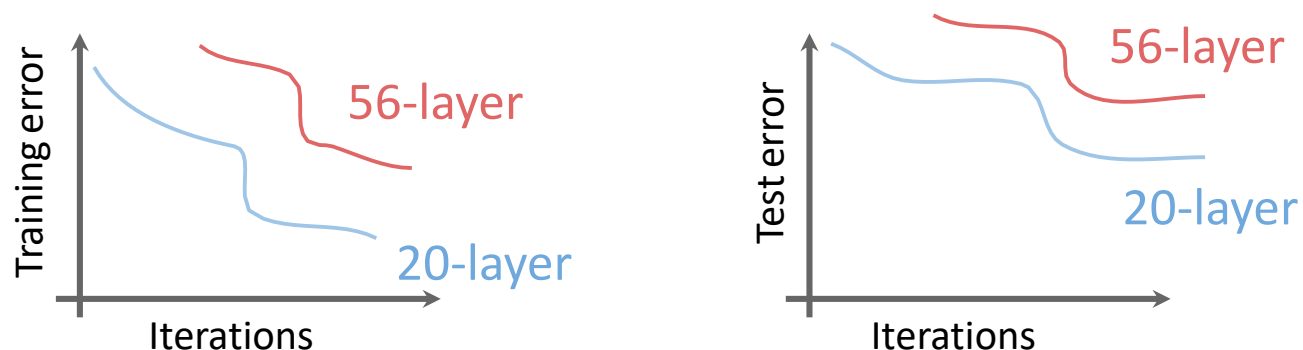**Very deep networks using residual  connections**

- 152-layer model for ImageNet
- ILSVRC'15 classification winner (3.57% top 5 error)
- Swept all classification and detection competitions in ILSVRC'15 and COCO'15!

$F(x) + x$ $\oplus$ relu

relu

conv

$F(x)$ relu

conv

$x$ identity

$x$

**Residual block**

Softmax

FC 1000

Ave. Pool

3*3 conv,512

3*3 conv,512

3*3 conv,512

3*3 conv,512

3*3 conv,512

3*3 conv,512,/2

...

3*3 conv,128

3*3 conv,128

3*3 conv,128

3*3 conv,128

3*3 conv,128

3*3 conv,128, /2

3*3 conv,64

3*3 conv,64

3*3 conv,64

3*3 conv,64

3*3 conv,64

3*3 conv,64

Max Pool, /2

7*7 conv, 64/2

Input

# Case Study: ResNet

[He et al., 2015]

What happens when we continue stacking deeper layers on a "plain" convolutional neural network?



Q: What's strange about these training and test curves?
[Hint: look at the order of the curves]

56-layer model performs worse on both training and test error.
→The deeper model performs worse, but it's not caused by overfitting!

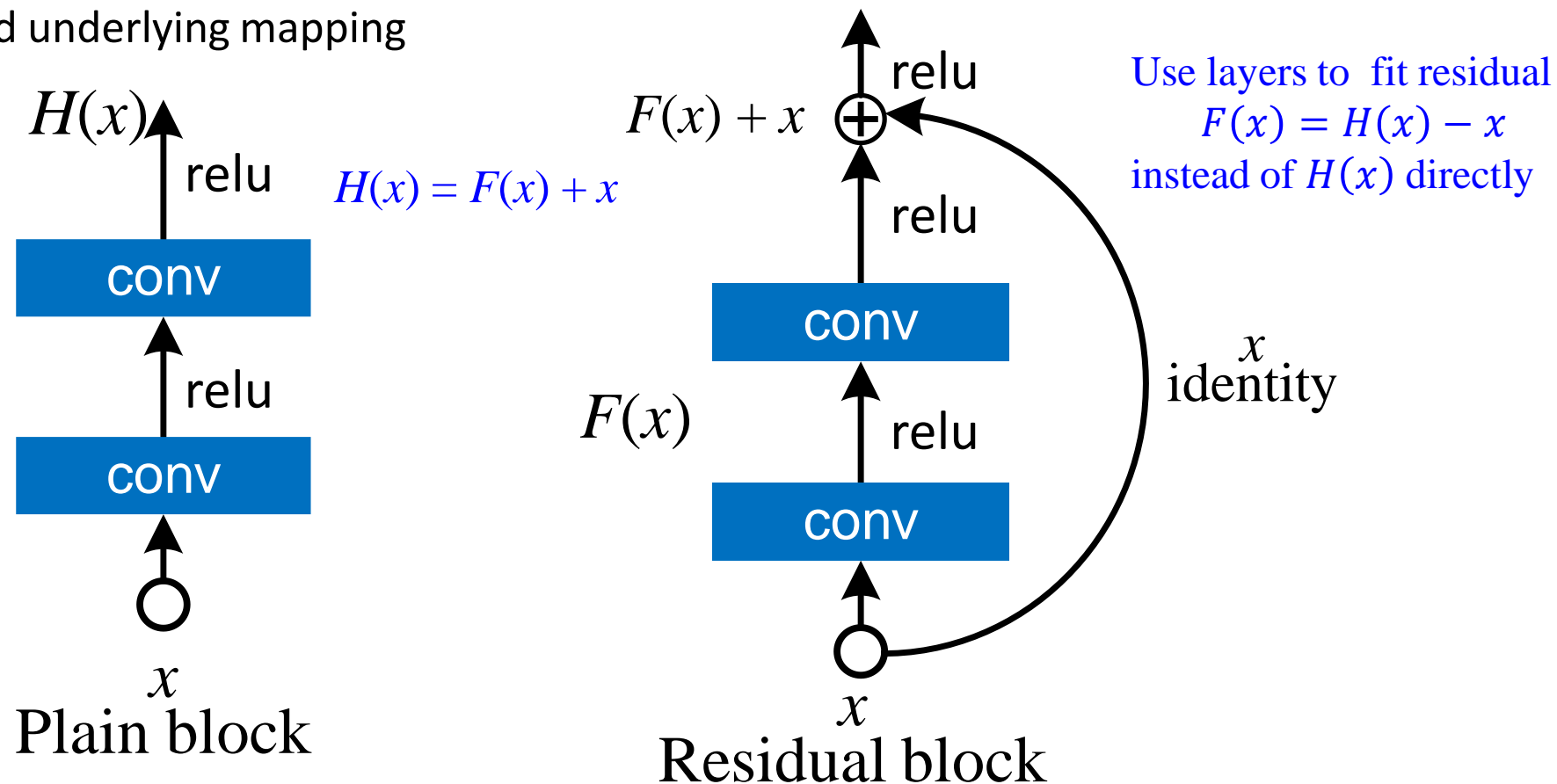# Case Study: ResNet

[He et al., 2015]

Hypothesis: the problem is an *optimization* problem, deeper models are harder to optimize

- The deeper model should be able to perform at least as well as the shallower model.
- A solution by construction is copying the learned layers from the shallower model and setting additional layers to identity mapping.

# Case Study: ResNet

[He et al., 2015]

Solution: Use network layers to fit a residual mapping instead of directly trying to fit a desired underlying mapping
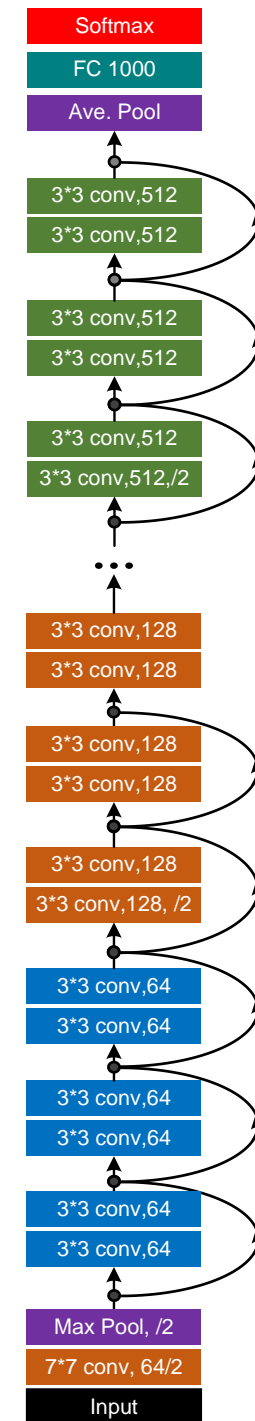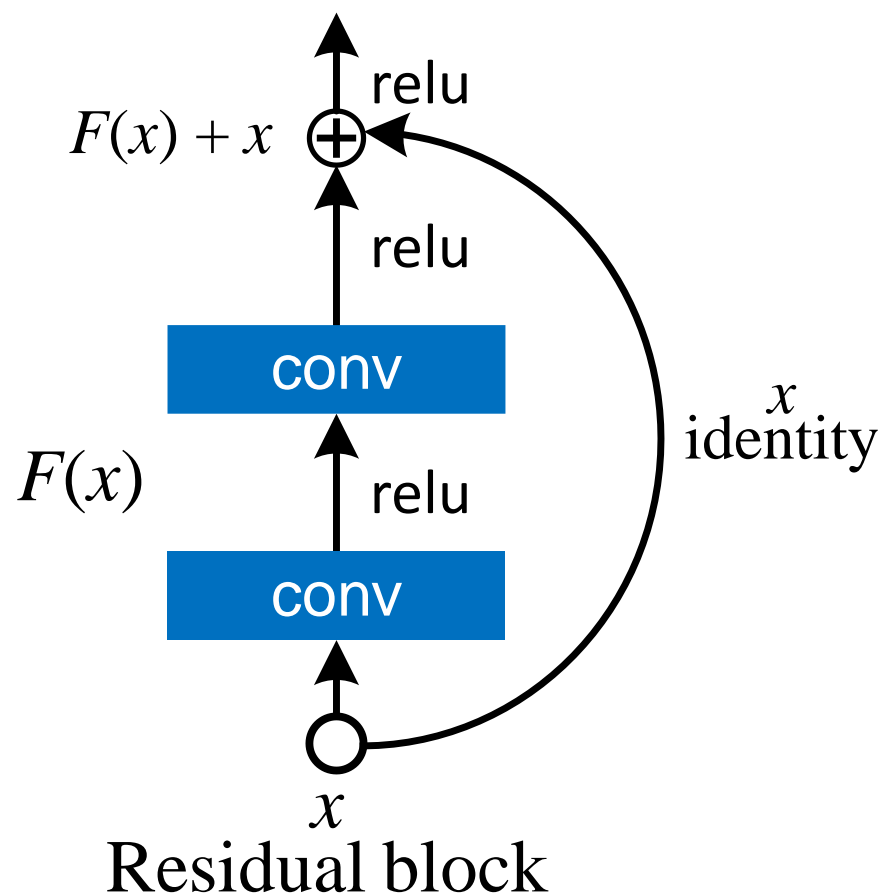
$H(x)$

relu

conv

relu

conv

$x$

Plain block

$H(x) = F(x) + x$

$F(x) + x$

relu

relu

conv

$F(x)$

relu

conv

$x$

$x$

identity

Residual block

Use layers to fit residual
$F(x) = H(x) - x$
instead of $H(x)$ directly

# Case Study: ResNet

[He et al., 2015]

Full ResNet architecture:

- Stack residual blocks
- Every residual block has two 3x3 conv layers

- Periodically, double the number of filters and subsample spatially using stride 2

- Additional conv layer at the beginning

- No FC layers at the end (only FC 1000 to output classes)

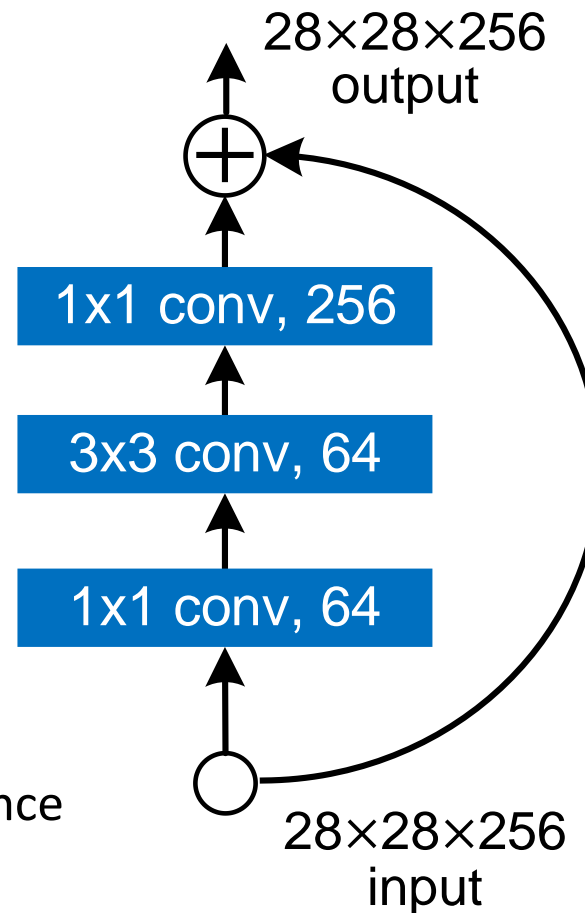Total depths of 34, 50, 101, or 152 layers for ImageNet



relu

$F(x) + x$ ⊕ ← relu

relu

conv

$F(x)$ relu

conv

$x$ identity

$x$

Residual block

Softmax
FC 1000
Ave. Pool
3*3 conv,512
3*3 conv,512
3*3 conv,512
3*3 conv,512
3*3 conv,512
3*3 conv,512,/2
...
3*3 conv,128
3*3 conv,128
3*3 conv,128
3*3 conv,128
3*3 conv,128
3*3 conv,128, /2
3*3 conv,64
3*3 conv,64
3*3 conv,64
3*3 conv,64
3*3 conv,64
3*3 conv,64
Max Pool, /2
7*7 conv, 64/2
Input

# Case Study: ResNet

[He et al., 2015]

For deeper networks (ResNet50+), use "bottleneck" layer to improve efficiency  (similar to GoogLeNet)

ILSVRC 2015 classification winner (3.6%  top 5 error)→better than "human  performance"!

1x1 conv, 256 filters projects back to 256 feature maps (28x28x256)

3x3 conv operates over only 64 feature maps

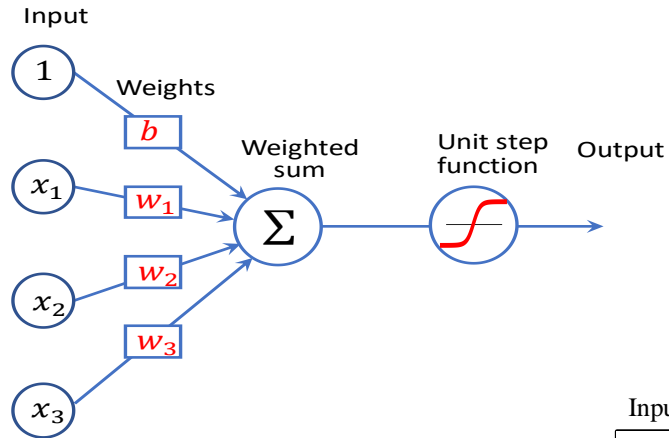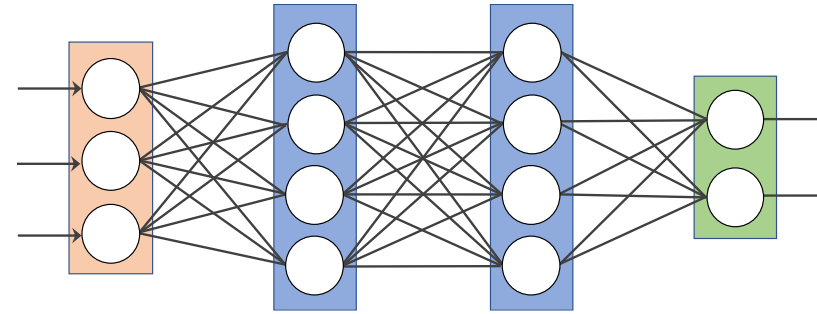1x1 conv, 64 filters  to project to  28x28x64

28×28×256 output

1x1 conv, 256

3x3 conv, 64

1x1 conv, 64

28×28×256 input

Experimental Results
- Able to train very deep  networks without performance degrading  (152 layers on ImageNet, 1202  on Cifar)
- Deeper networks now achieve  low training error as expected
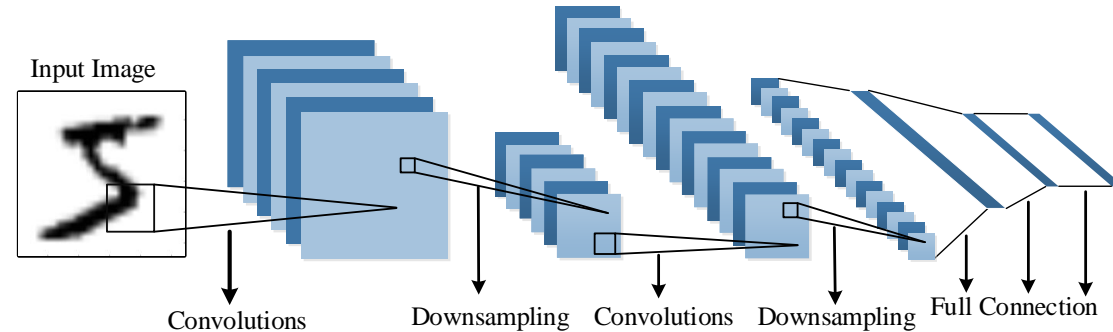- Swept 1st place in all ILSVRC  and COCO 2015 competitions

# In this Course
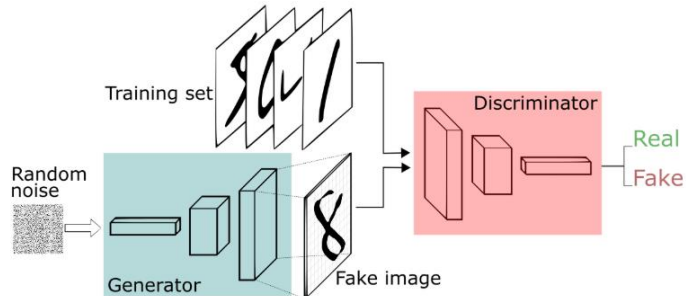
## 1. DL basics, linear regression, logistic regression etc.

Input

1

Weights

$b$

$x_1$

$w_1$

Weighted sum

Unit step function

Output

$x_2$

$w_2$

$\Sigma$

$f$

$w_3$

$x_3$

## 2. Multilayer neural networks, backpropagation

## 3. Convolutional Neural Networks and Applications

Input Image

Convolutions          Downsampling    Convolutions    Downsampling    Full Connection

## 4. Generative Adversarial Networks

Training set

Discriminator

Real

Fake

Random noise

Generator

Fake image

## 5. Recurrent networks and applications

$h_t$

A

=

$h_0$     $h_1$     $h_2$     $h_t$

A     A     A     A

$x_t$     $x_0$     $x_1$     $x_2$     ...     $x_t$

# For this Final Lecture Today: Tips for Training Deep Learning

# Case Study: AlexNet

[Krizhevsky *et al.* 2012]



Full (simplified) AlexNet architecture:
[227x227x3] INPUT
[55x55x96] CONV1: 96 11x11 filters at stride 4, pad 0
[27x27x96] MAX POOL1: 3x3 filters at stride 2
[27x27x96] NORM1: Normalization layer
[27x27x256] CONV2: 256 5x5 filters at stride 1, pad 2
[13x13x256] MAX POOL2: 3x3 filters at stride 2
[13x13x256] NORM2: Normalization layer
[13x13x384] CONV3: 384 3x3 filters at stride 1, pad 1
[13x13x384] CONV4: 384 3x3 filters at stride 1, pad 1
[13x13x256] CONV5: 256 3x3 filters at stride 1, pad 1
[6x6x256] MAX POOL3: 3x3 filters at stride 2
[4096] FC6: 4096 neurons
[4096] FC7: 4096 neurons
[1000] FC8: 1000 neurons (class scores)

**Details:**
- heavy data augmentation
- first use of ReLU
- used Norm layers (not common anymore)
- dropout 0.5
- batch size 128
- SGD Momentum 0.9
- Learning rate 0.01, reduced by 10 manually when val accuracy plateaus
- L2 weight decay 0.0005

7 CNN ensemble: 18.2%→ 15.4%

# Data Augmentation
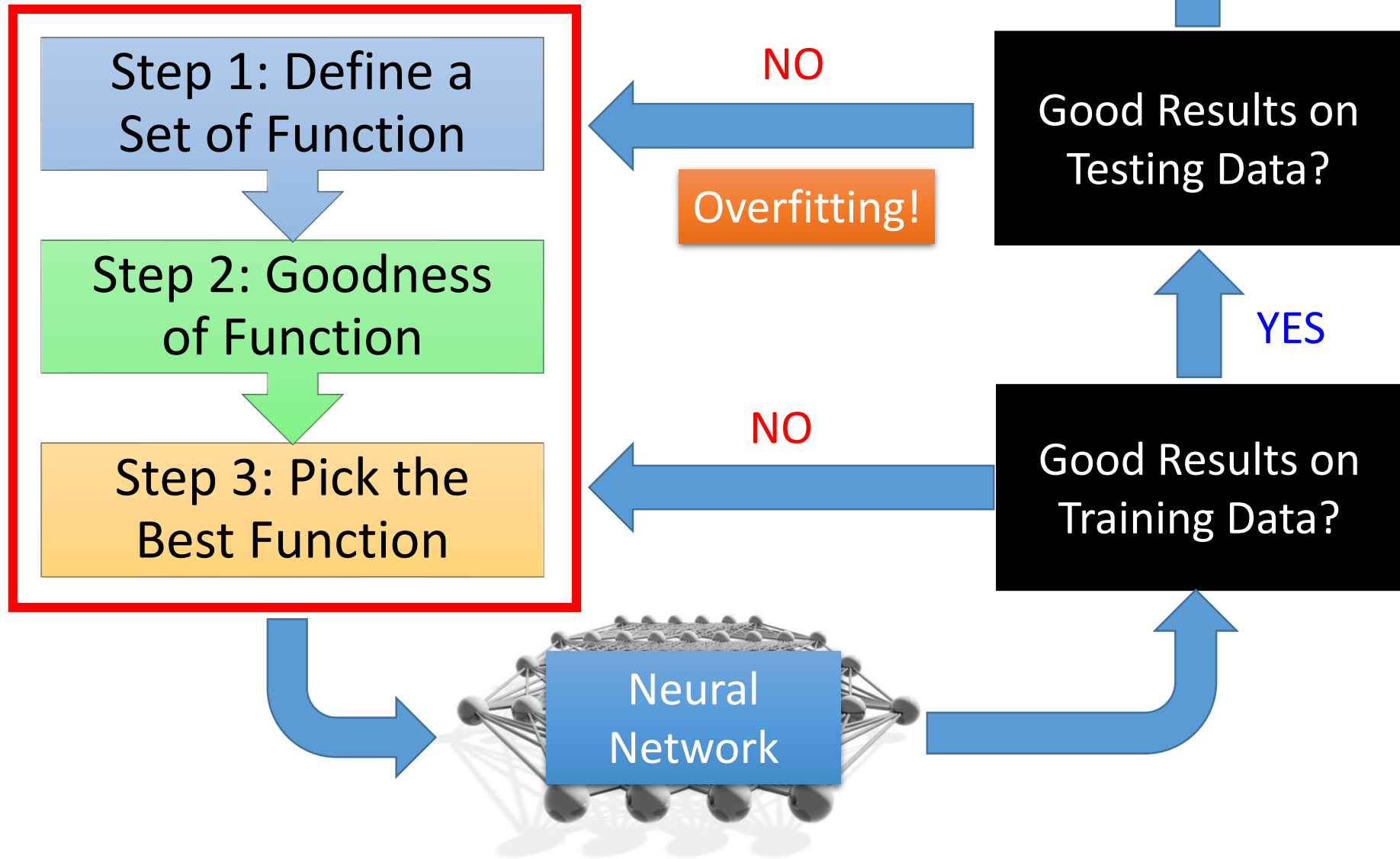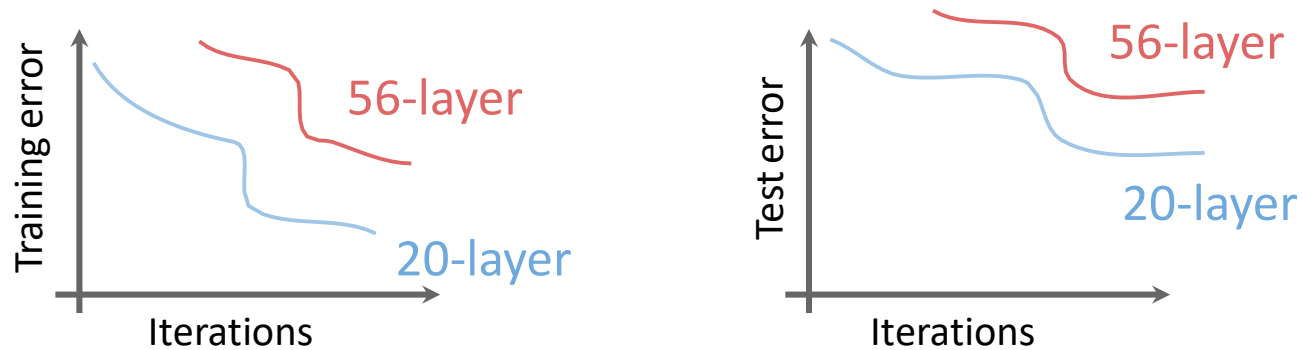
Crop:

Flip:

Scale:

Rotate:

Translation:

Noise:

# Recipe of Deep Learning

# Do not always blame Overfitting



Please refer to:
Kaiming He, Xiangyu Zhang, Shaoqing Ren, Jian Sun, Deep Residual Learning for Image Recognition, CVPR, 2016.

# Recipe of Deep Learning

# Recipe of Deep Learning

# Hard to get the power of Deep ...



**Handwriting Digit Classification**

Results on Training Data

Deeper usually does not imply better.

# Vanishing Gradient Problem

Input

Layer 1

Layer 2

Layer $L-1$

Layer L

$x_1^0$

$x_2^0$

$x_{n^0}^0$

$x_1^1$
$x_2^1$
$x_3^1$
$x_{n^1}^1$

$x_1^2$
$x_2^2$
$x_3^2$
$x_4^2$
$x_{n^2}^2$

$x_1^{L-1}$
$x_2^{L-1}$
$x_3^{L-1}$
$x_{n^L}^{L-1}$

$y_1$
$y_2$
$y_3$
$y_4$
$y_c$

$\{W_i$

$\{W_{ij}^2\}$

$\{W_{ij}^L\}$

Smaller gradients

Learn very slow

Almost random

Larger gradients

Learn very fast

Already converge

converge based on random!?

# Vanishing Gradient Problem



Intuitive way to compute the derivatives ...

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

# ReLU
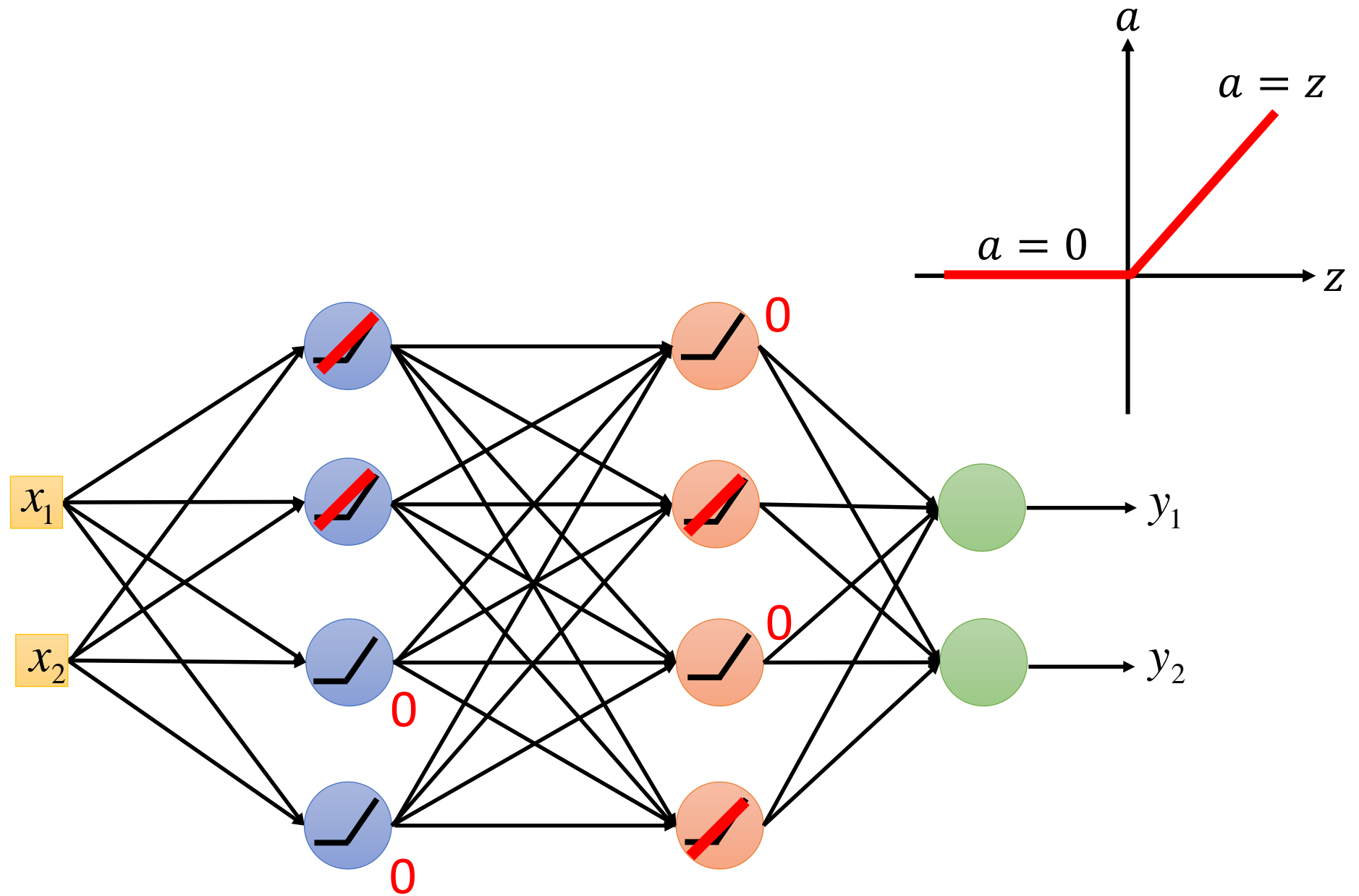
- Rectified Linear Unit (ReLU)
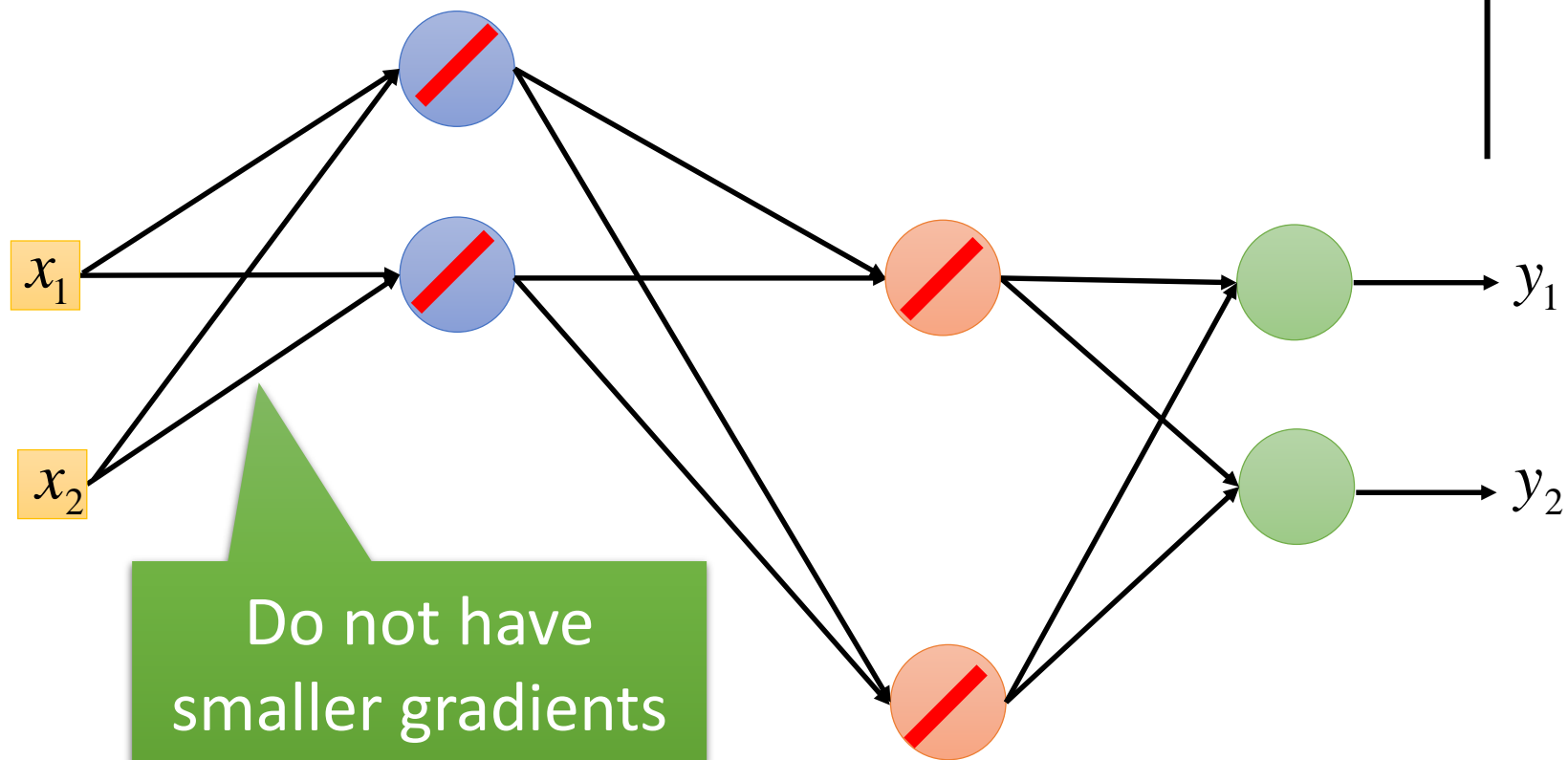
$\sigma(z)$



**Reason:**

1. Fast to compute

2. Biological reason

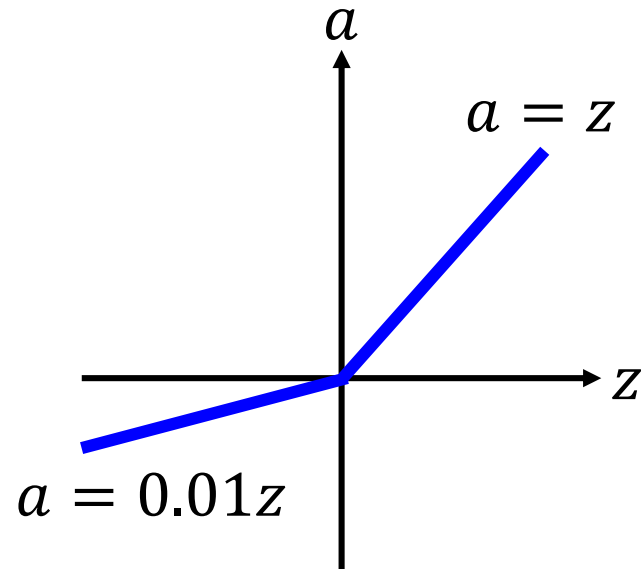3. Infinite sigmoid with different biases

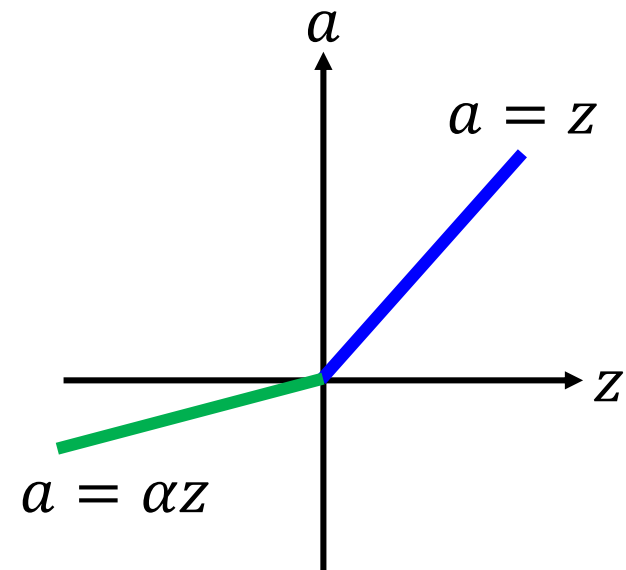4. Vanishing gradient problem

ReLU

# ReLU

A Thinner linear network

Do not have smaller gradients

$a$

$a = z$

$a = 0$

$z$

$x_1$

$x_2$

$y_1$

$y_2$

# Variants of ReLU

## Leaky ReLU

$a$

$a = z$

$z$

$a = 0.01z$

## Parametric ReLU

$a$

$a = z$

$z$

$a = \alpha z$
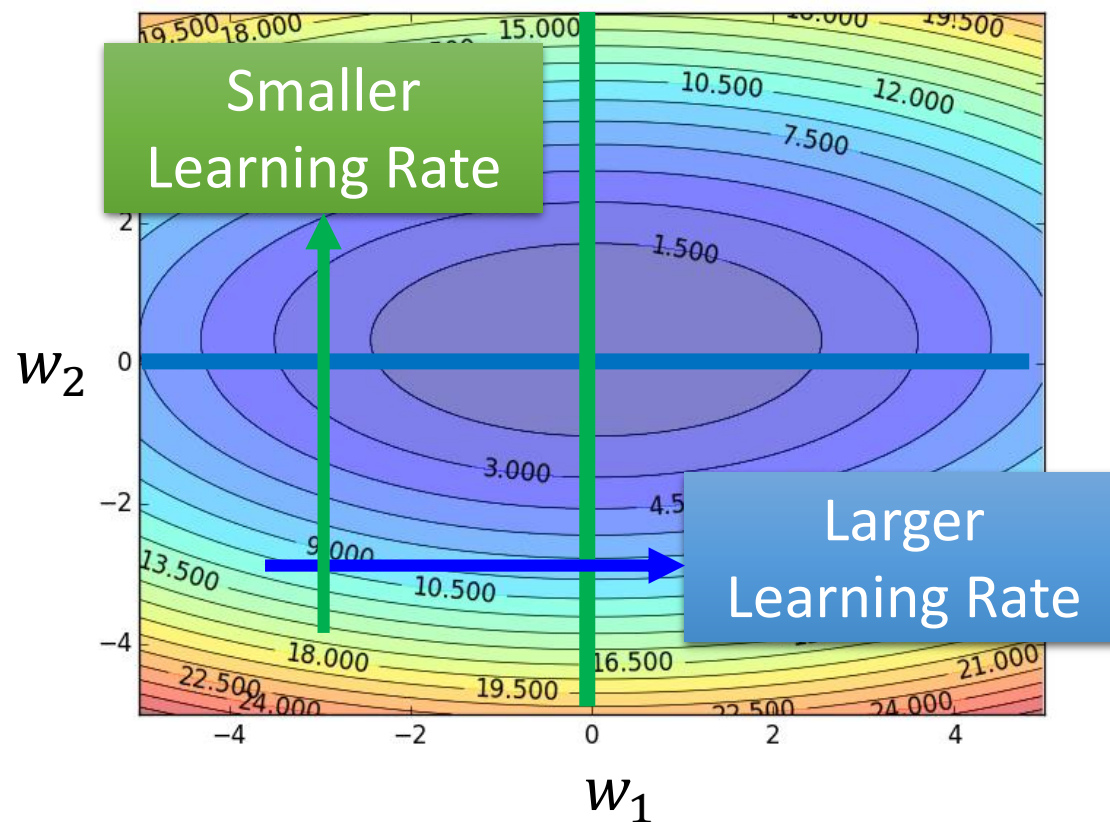
α also learned by
gradient descent

# Maxout

- Learnable activation function [Ian J. Goodfellow, ICML'13]

# Recipe of Deep Learning
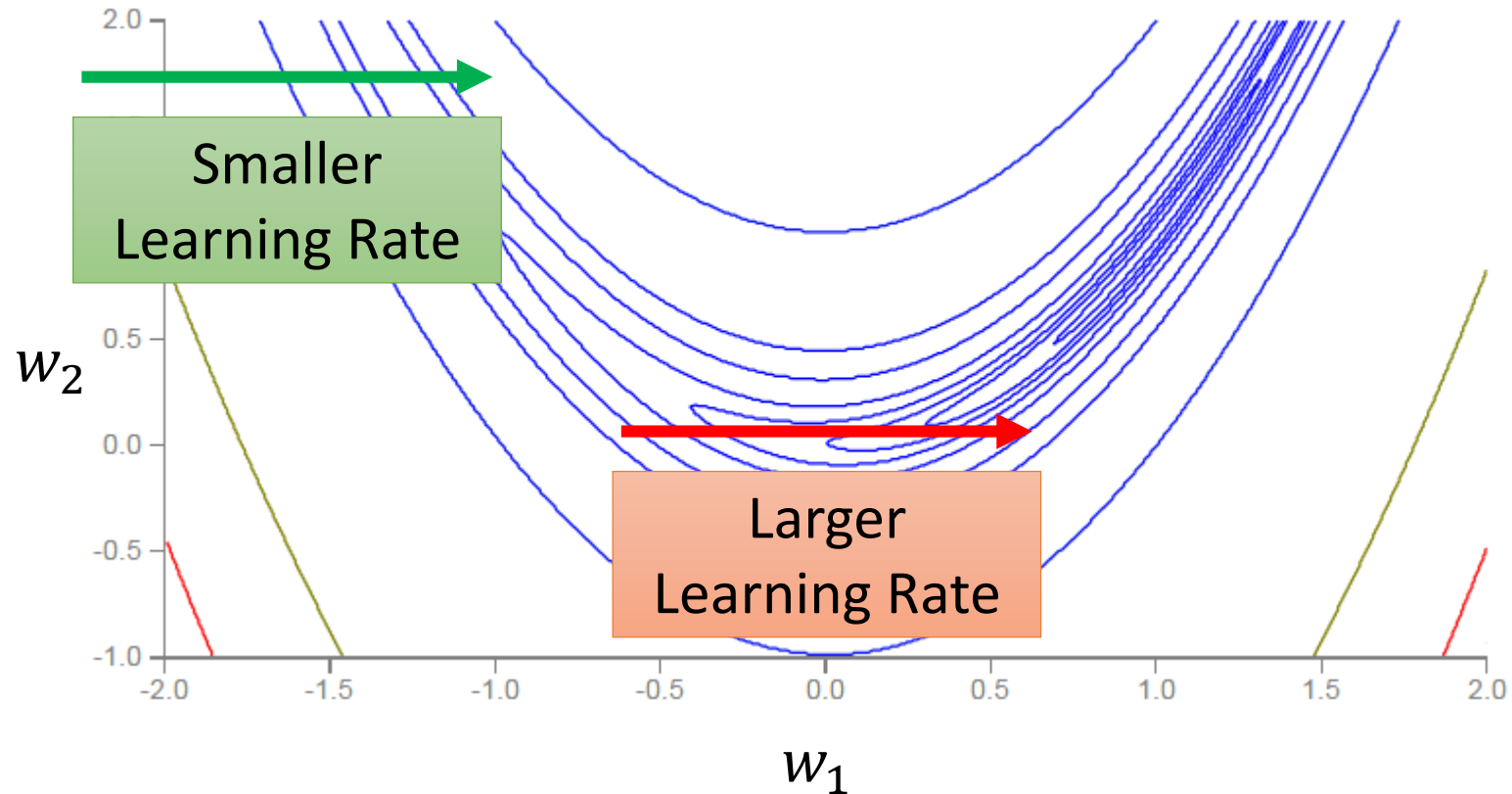
# Review



*AdaGrad*

$$w^{(t+1)} \leftarrow w^{(t)} - \frac{\eta}{\sqrt{\sum_{i=0}^{t}(g^{(i)})^2}} g^{(t)}$$

Use first derivative to estimate second derivative

# RMSProp

Error Surface can be very complex when training NN.

# RMSProp

$$w^{(1)} \leftarrow w^{(0)} - \frac{\eta}{\sigma^{(0)}} g^{(0)}$$

$$w^{(2)} \leftarrow w^{(1)} - \frac{\eta}{\sigma^{(1)}} g^{(1)}$$

$$w^{(3)} \leftarrow w^{(2)} - \frac{\eta}{\sigma^{(2)}} g^{(2)}$$

$$\vdots$$

$$w^{(t+1)} \leftarrow w^{(t)} - \frac{\eta}{\sigma^{(t)}} g^{(t)}$$

$$w^{(t+1)} \leftarrow w^{(t)} - \frac{\eta}{\sqrt{\sum_{i=0}^{t} (g^{(i)})^2}} g^{(t)}$$

$$AdaGrad$$

$$\sigma^{(0)} = g^{(0)}$$

$$\sigma^{(1)} = \sqrt{\alpha(\sigma^{(0)})^2 + (1-\alpha)(g^{(1)})^2}$$

$$\sigma^{(2)} = \sqrt{\alpha(\sigma^{(1)})^2 + (1-\alpha)(g^{(2)})^2}$$

$$\vdots$$
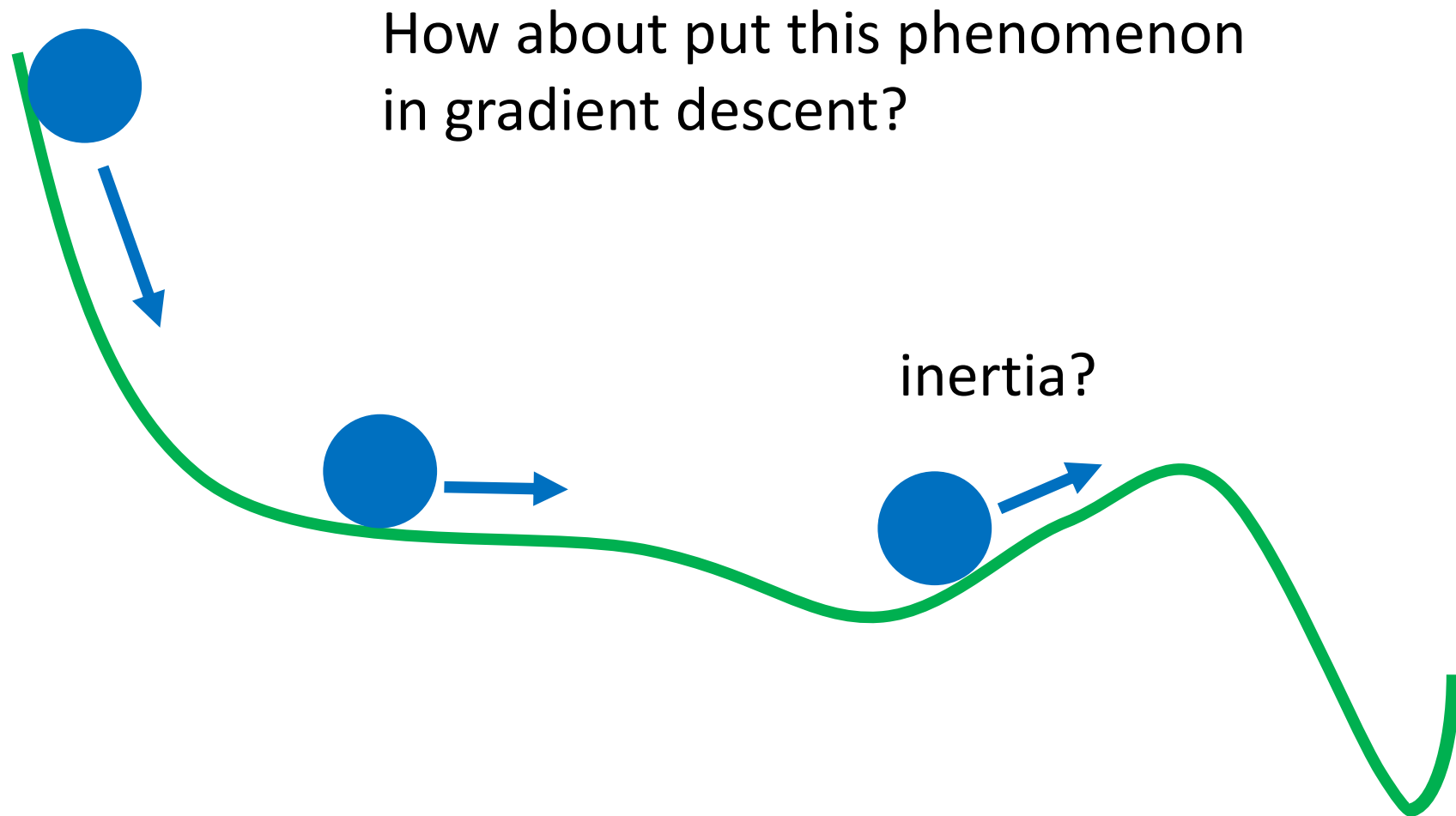
$$\sigma^{(t)} = \sqrt{\alpha(\sigma^{(t-1)})^2 + (1-\alpha)(g^{(t)})^2}$$

Root Mean Square of the gradients
with previous gradients being decayed

# Hard to find optimal network parameters



Very slow at the **plateau**

Stuck at saddle point

Stuck at local minima

Loss

$$\frac{\partial L}{\partial w} \approx 0$$

$$\frac{\partial L}{\partial w} = 0$$

$$\frac{\partial L}{\partial w} = 0$$

The value of the parameter w

# In physical world ......

- Momentum

How about put this phenomenon in gradient descent?

inertia?

# Review: Vanilla Gradient Descent

Gradient: the normal direction of the contour of loss function



$\nabla L(\theta^{(0)})$

$\theta^{(0)}$

$\nabla L(\theta^{(1)})$

$\theta^{(1)}$

$\nabla L(\theta^{(2)})$

$\theta^{(2)}$

$\nabla L(\theta^{(3)})$

$\theta^{(3)}$

→ Gradient

→ Movement

Start at position $\theta^{(0)}$

Compute gradient at $\theta^{(0)}$

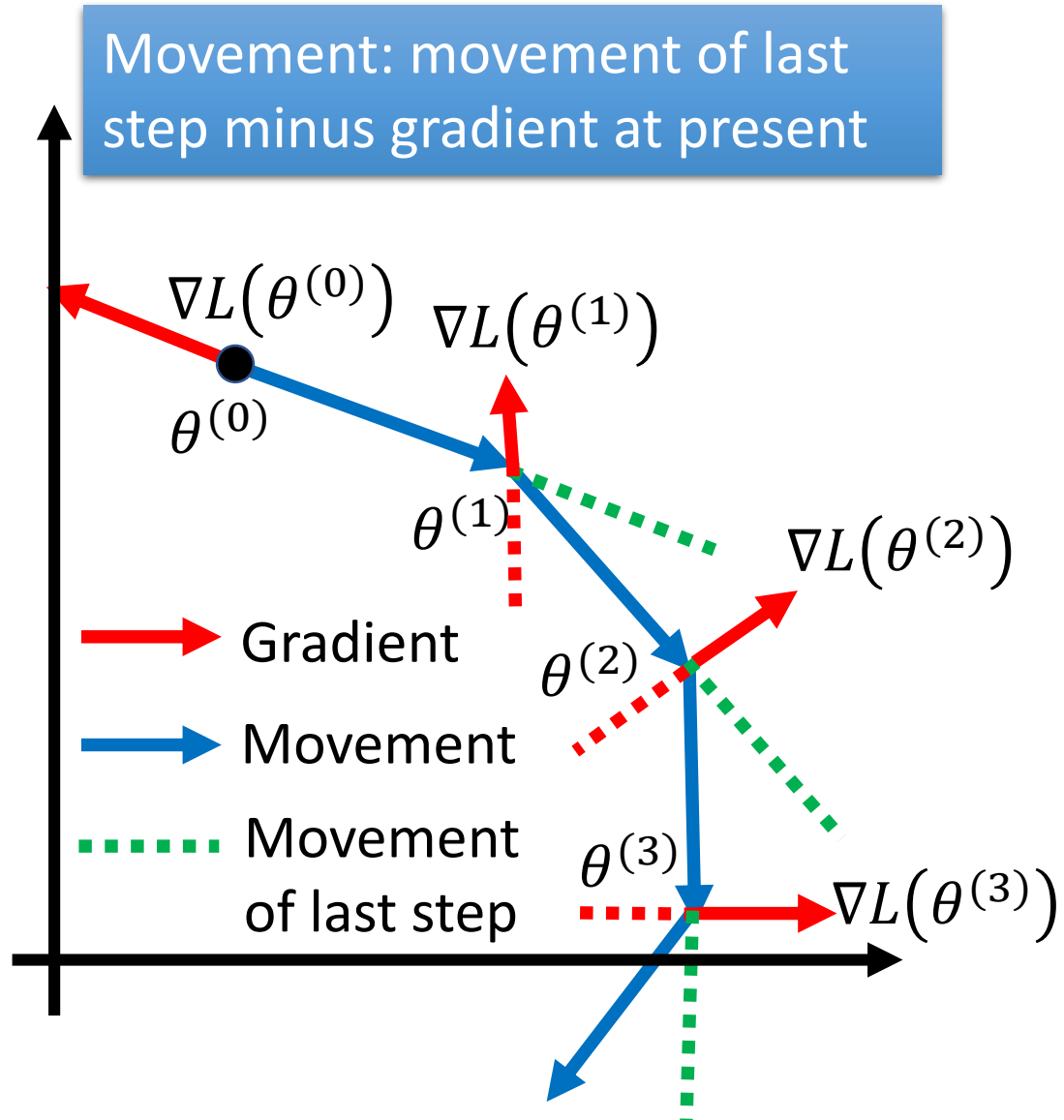Move to $\theta^{(1)} = \theta^{(0)} - \eta \nabla L(\theta^{(0)})$

Compute gradient at $\theta^{(1)}$

Move to $\theta^{(2)} = \theta^{(1)} - \eta \nabla L(\theta^{(1)})$

⋮

Stop until $\nabla L(\theta^{(t)}) \approx 0$

# Momentum



Movement: movement of last step minus gradient at present

$\nabla L(\theta^{(0)})$    $\nabla L(\theta^{(1)})$

$\theta^{(0)}$

$\theta^{(1)}$

$\nabla L(\theta^{(2)})$

$\theta^{(2)}$

→ Gradient

→ Movement

⋯ Movement of last step

$\theta^{(3)}$

$\nabla L(\theta^{(3)})$

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 $\theta^{(1)} = \theta^{(0)} + v^{(1)}$

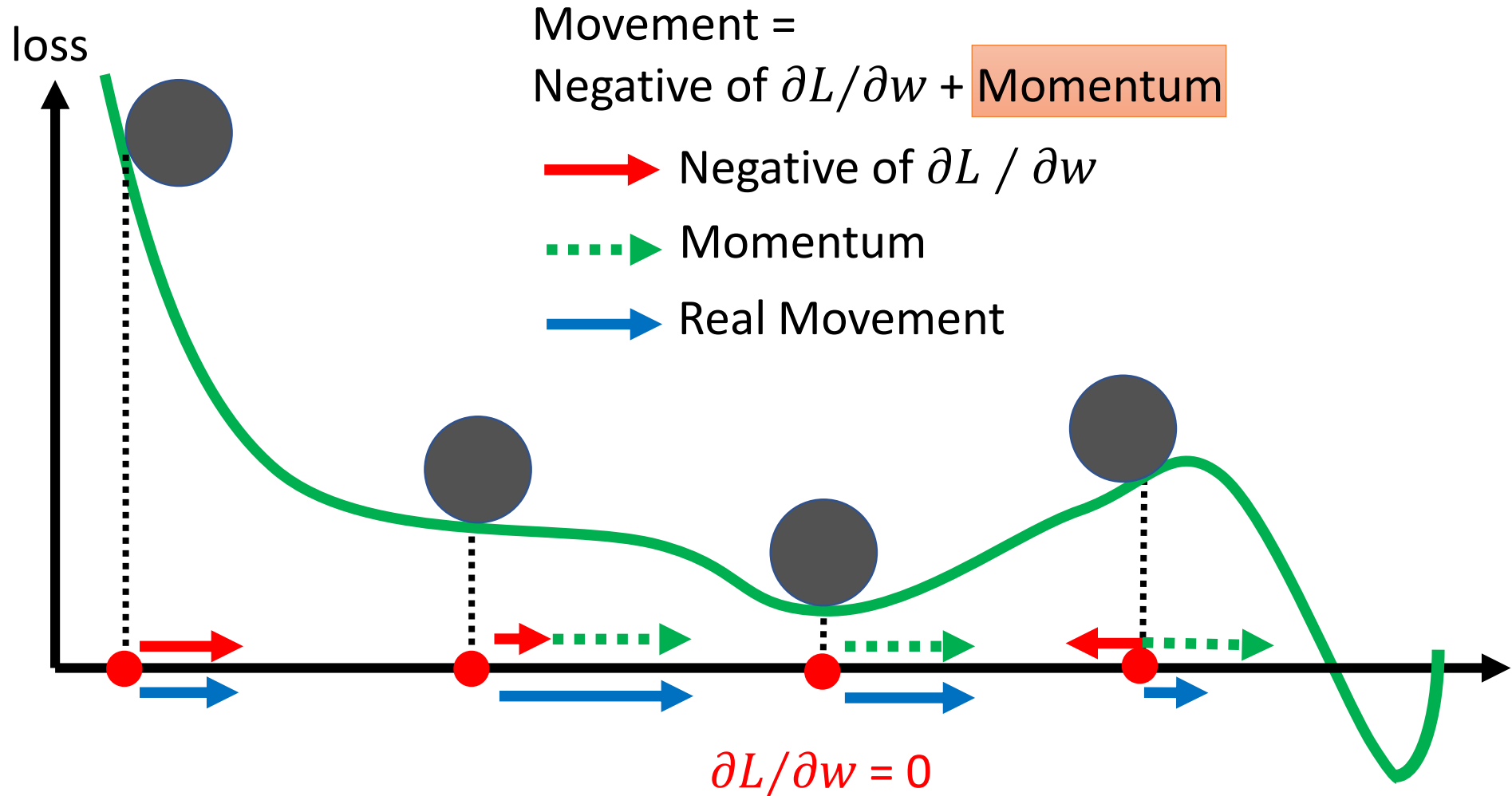Compute gradient at $\theta^{(1)}$

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.

# Momentum

Movement =
Negative of $\partial L / \partial w$ + Momentum

→ Negative of $\partial L / \partial w$

⇢ Momentum

→ Real Movement

loss

$\partial L / \partial w = 0$

# Adam

**Algorithm 1:** *Adam*, our proposed algorithm for stochastic optimization. See section 2 for details, and for a slightly more efficient (but less clear) order of computation. $g_t^2$ indicates the elementwise square $g_t \odot g_t$. Good default settings for the tested machine learning problems are $\alpha = 0.001$, $\beta_1 = 0.9$, $\beta_2 = 0.999$ and $\epsilon = 10^{-8}$. All operations on vectors are element-wise. With $\beta_1^t$ and $\beta_2^t$ we denote $\beta_1$ and $\beta_2$ to the power $t$.

**Require:** $\alpha$: Stepsize
**Require:** $\beta_1, \beta_2 \in [0, 1)$: Exponential decay rates for the moment estimates
**Require:** $f(\theta)$: Stochastic objective function with parameters $\theta$
**Require:** $\theta_0$: Initial parameter vector
  $m_0 \leftarrow 0$ (Initialize 1$^{\text{st}}$ moment vector) $\longrightarrow$ for momentum
  $v_0 \leftarrow 0$ (Initialize 2$^{\text{nd}}$ moment vector) $\longrightarrow$ for RMSprop
  $t \leftarrow 0$ (Initialize timestep)
  **while** $\theta_t$ not converged **do**
    $t \leftarrow t + 1$
    $g_t \leftarrow \nabla_\theta f_t(\theta_{t-1})$ (Get gradients w.r.t. stochastic objective at timestep $t$)
    $m_t \leftarrow \beta_1 \cdot m_{t-1} + (1 - \beta_1) \cdot g_t$ (Update biased first moment estimate)
    $v_t \leftarrow \beta_2 \cdot v_{t-1} + (1 - \beta_2) \cdot g_t^2$ (Update biased second raw moment estimate)
    $\widehat{m}_t \leftarrow m_t / (1 - \beta_1^t)$ (Compute bias-corrected first moment estimate)
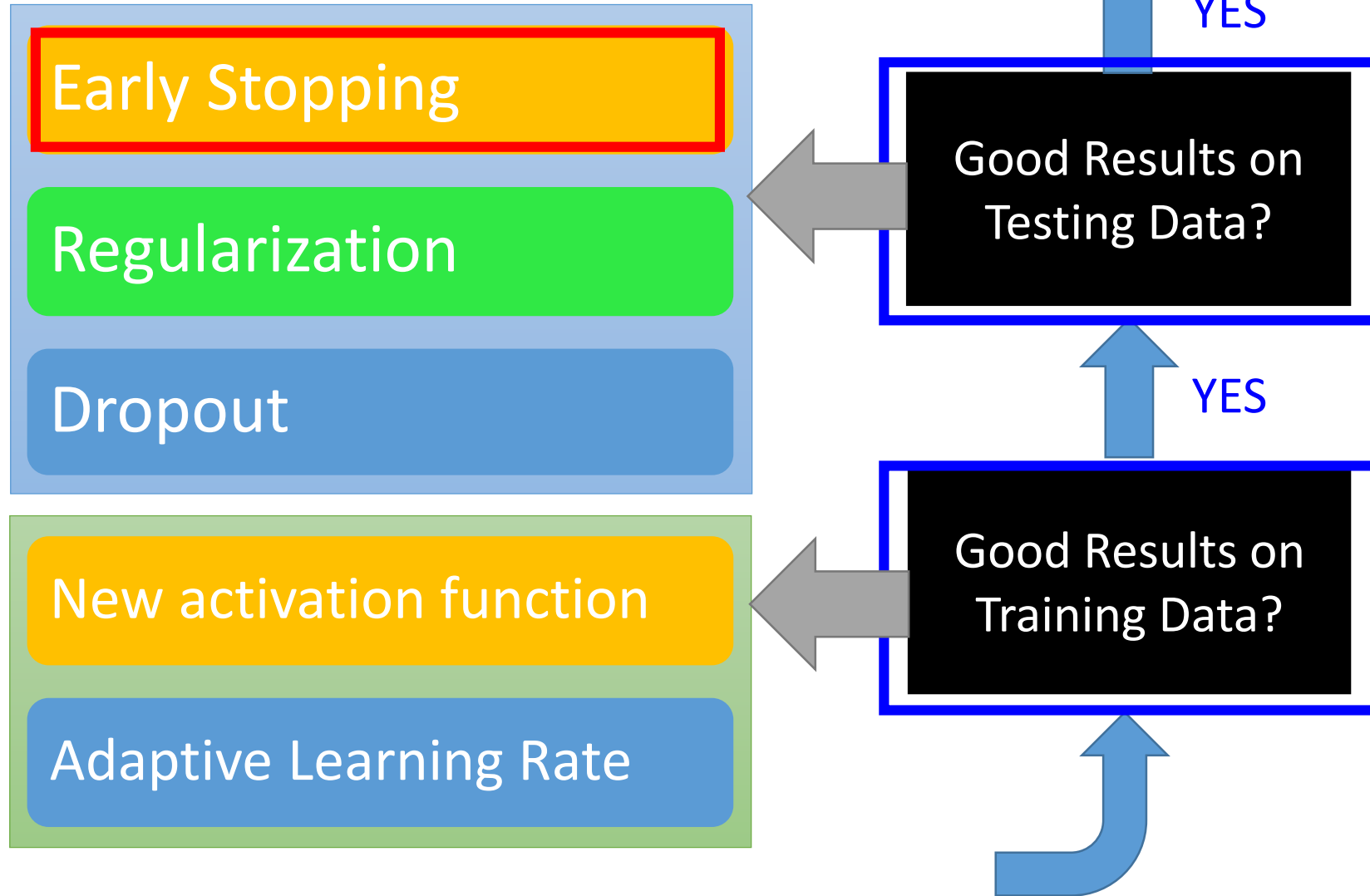    $\widehat{v}_t \leftarrow v_t / (1 - \beta_2^t)$ (Compute bias-corrected second raw moment estimate)
    $\theta_t \leftarrow \theta_{t-1} - \alpha \cdot \widehat{m}_t / (\sqrt{\widehat{v}_t} + \epsilon)$ (Update parameters)
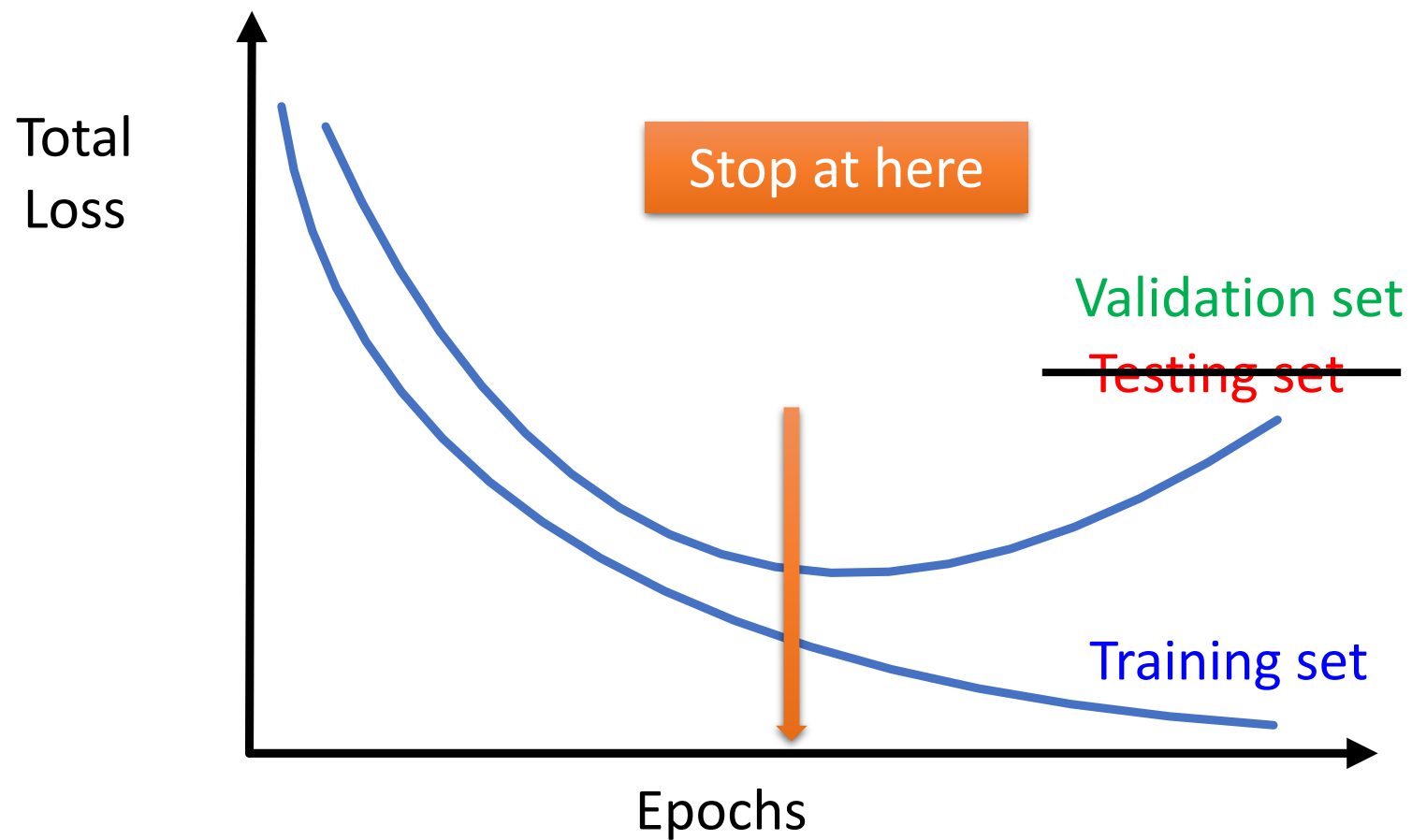  **end while**
  **return** $\theta_t$ (Resulting parameters)

RMSProp + Momentum

# Recipe of Deep Learning

# Early Stopping

# Recipe of Deep Learning

Early Stopping

Regularization

Dropout

New activation function

Adaptive Learning Rate

Good Results on Testing Data?

YES

Good Results on Training Data?

YES
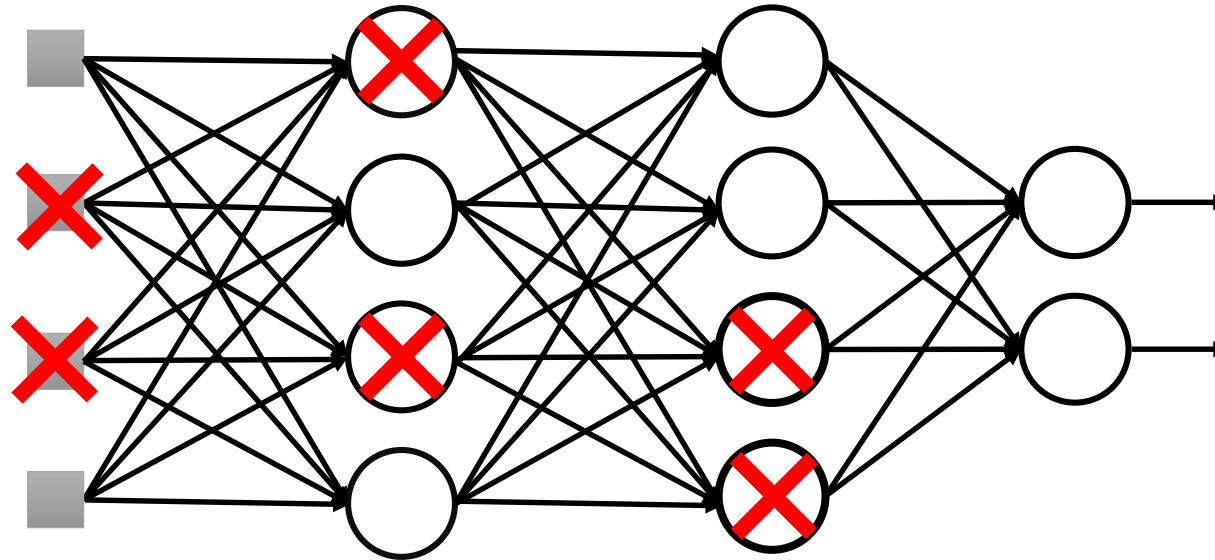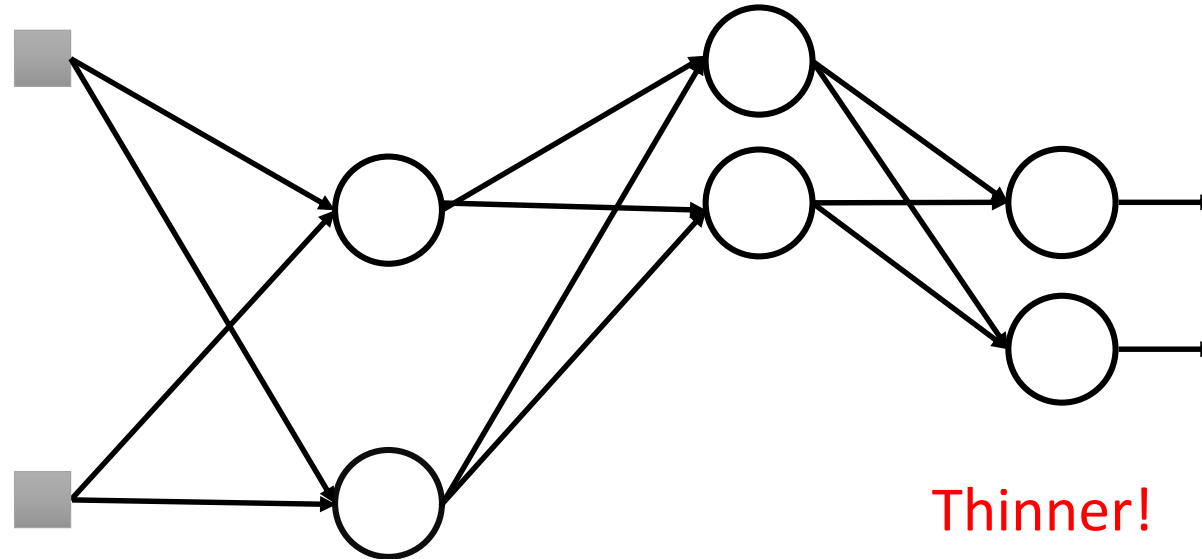
# Dropout

**Training:**



□ **Each time before updating the parameters**

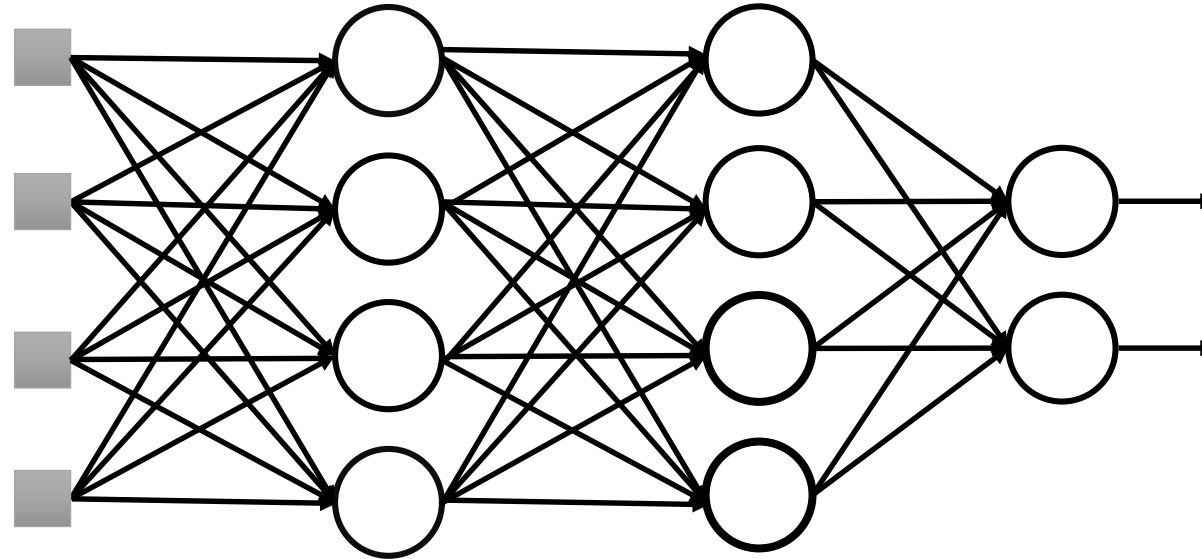• Each neuron has $p\%$ to be preserved, *i.e.* $1 - p\%$ to dropout

# Dropout



*Training:*

Thinner!

□ **Each time before updating the parameters**

- Each neuron has $p\%$ to be preserved

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

- Using the new network for training

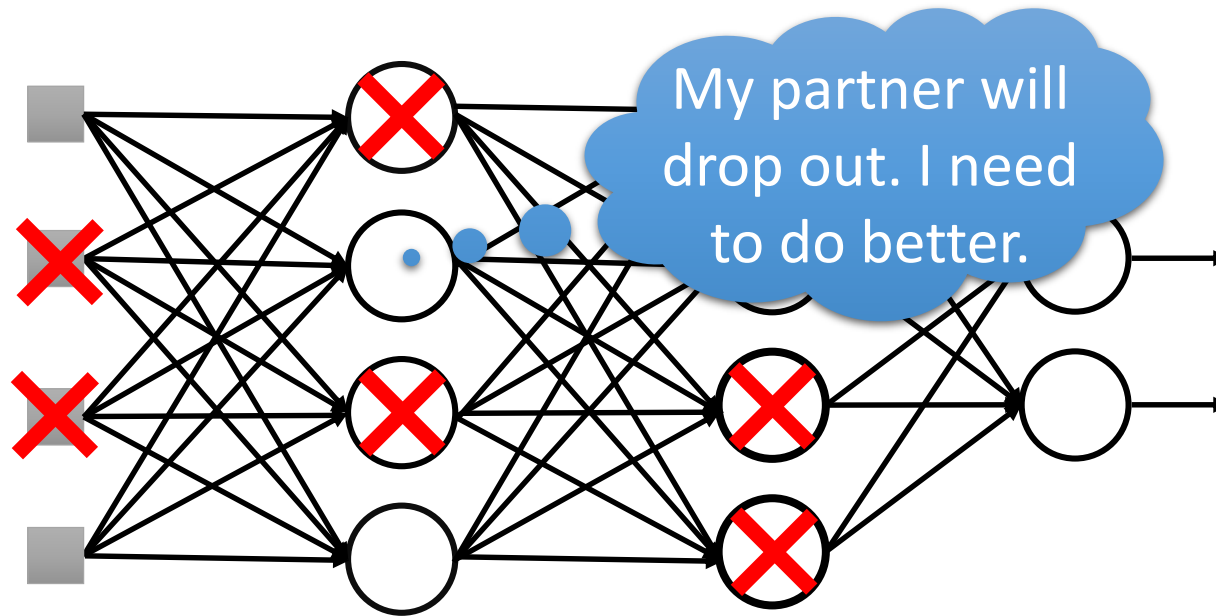For each minibatch, we resample the dropout neurons

# Dropout

**No dropout**

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

# Dropout: Intuitive Reason

Nitish Srivastava                                    NITISH@CS.TORONTO.EDU
Geoffrey Hinton                                      HINTON@CS.TORONTO.EDU
Alex Krizhevsky                                        KRIZ@CS.TORONTO.EDU
Ilya Sutskever                                         ILYA@CS.TORONTO.EDU
Ruslan Salakhutdinov                              RSALAKHU@CS.TORONTO.EDU
*Department of Computer Science*
*University of Toronto*
*10 Kings College Road, Rm 3302*
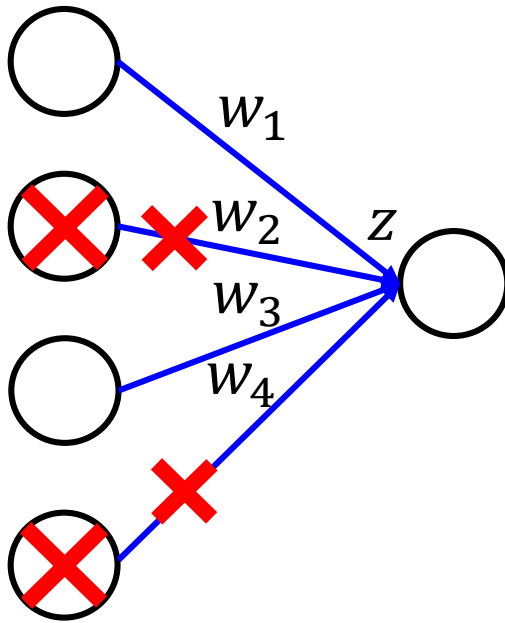*Toronto, Ontario, M5S 3G4, Canada.*

- When teams up, if everyone expect the partner will do the work, nothing will be done finally.

- However, if you know your partner will dropout, you will do better.

- When testing, no one dropout actually, so obtaining good results eventually.

# Dropout: Intuitive Reason

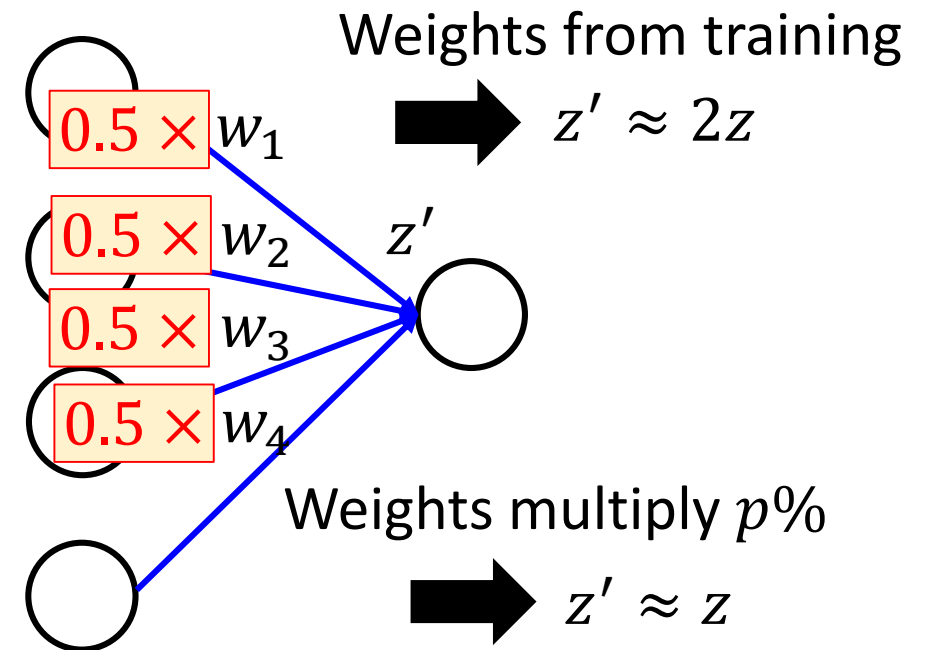- Why the weights should multiply $p\%$ ($p\%$ is the keep rate) when testing?

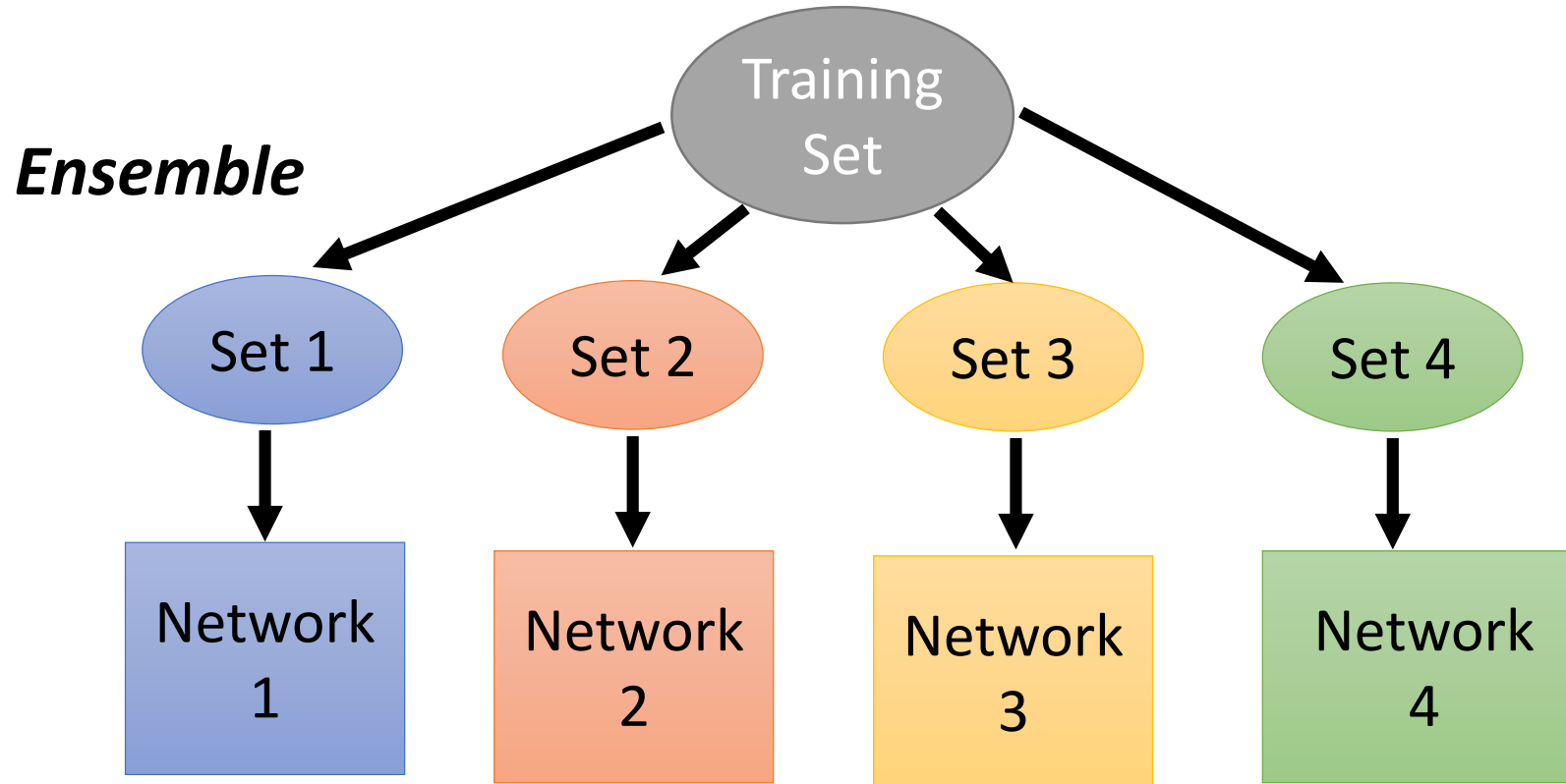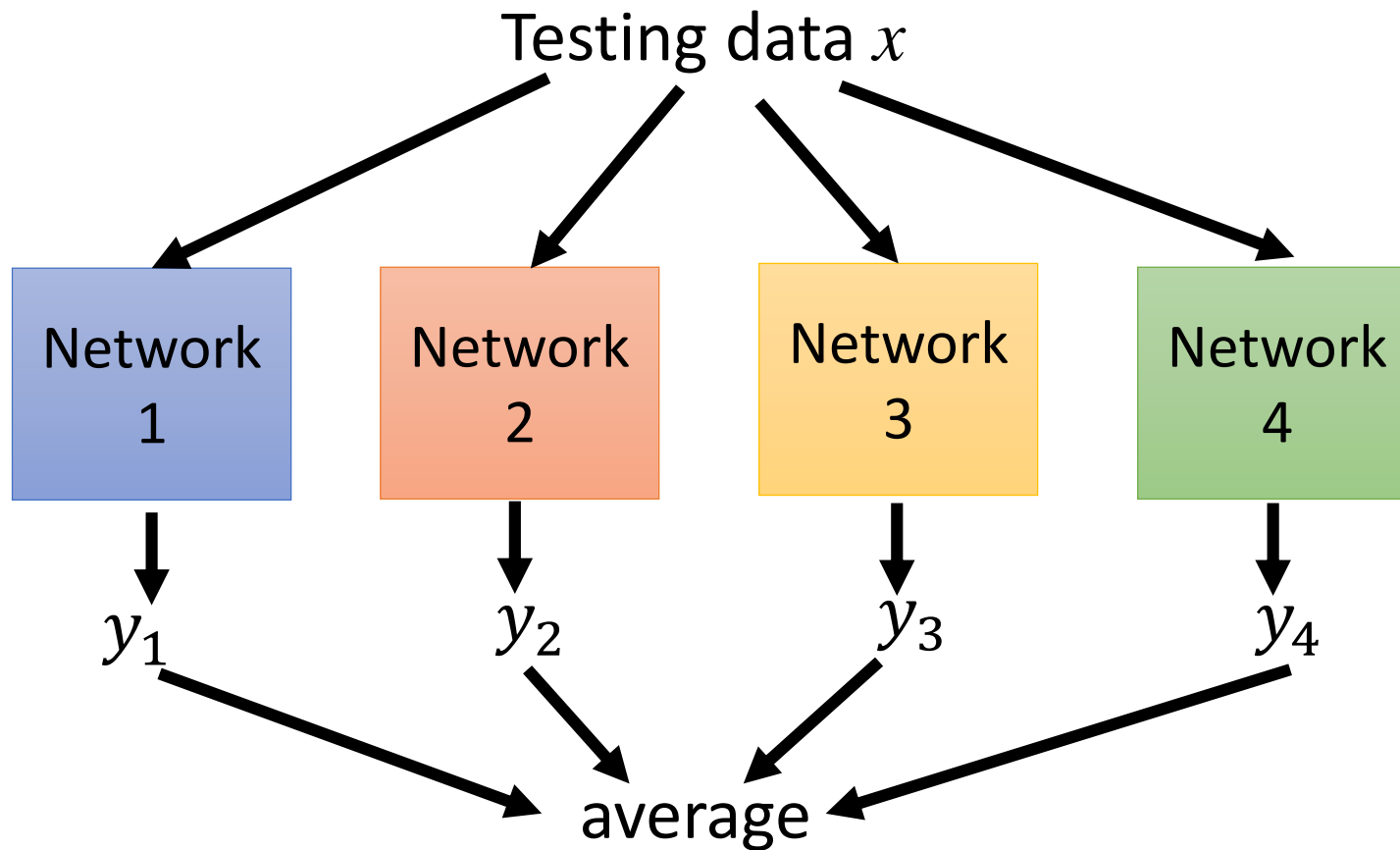**Training of Dropout**

Assume keep rate is 50%

$w_1$

$w_2$

$w_3$

$w_4$

$z$

**Testing of Dropout**

No dropout

$0.5 \times$ $w_1$

$0.5 \times$ $w_2$

$0.5 \times$ $w_3$

$0.5 \times$ $w_4$

$z'$

Weights from training

$z' \approx 2z$

Weights multiply $p\%$

$z' \approx z$

# Dropout is a kind of ensemble.
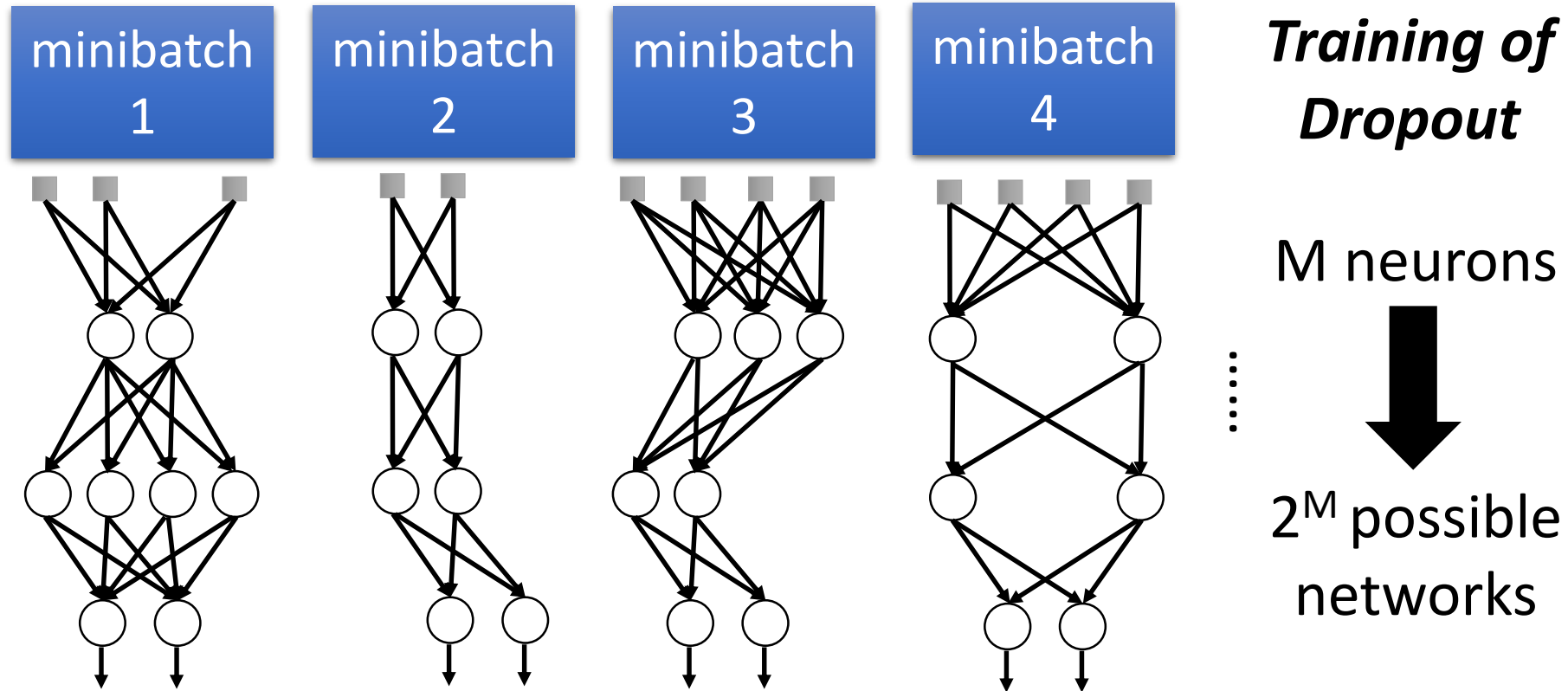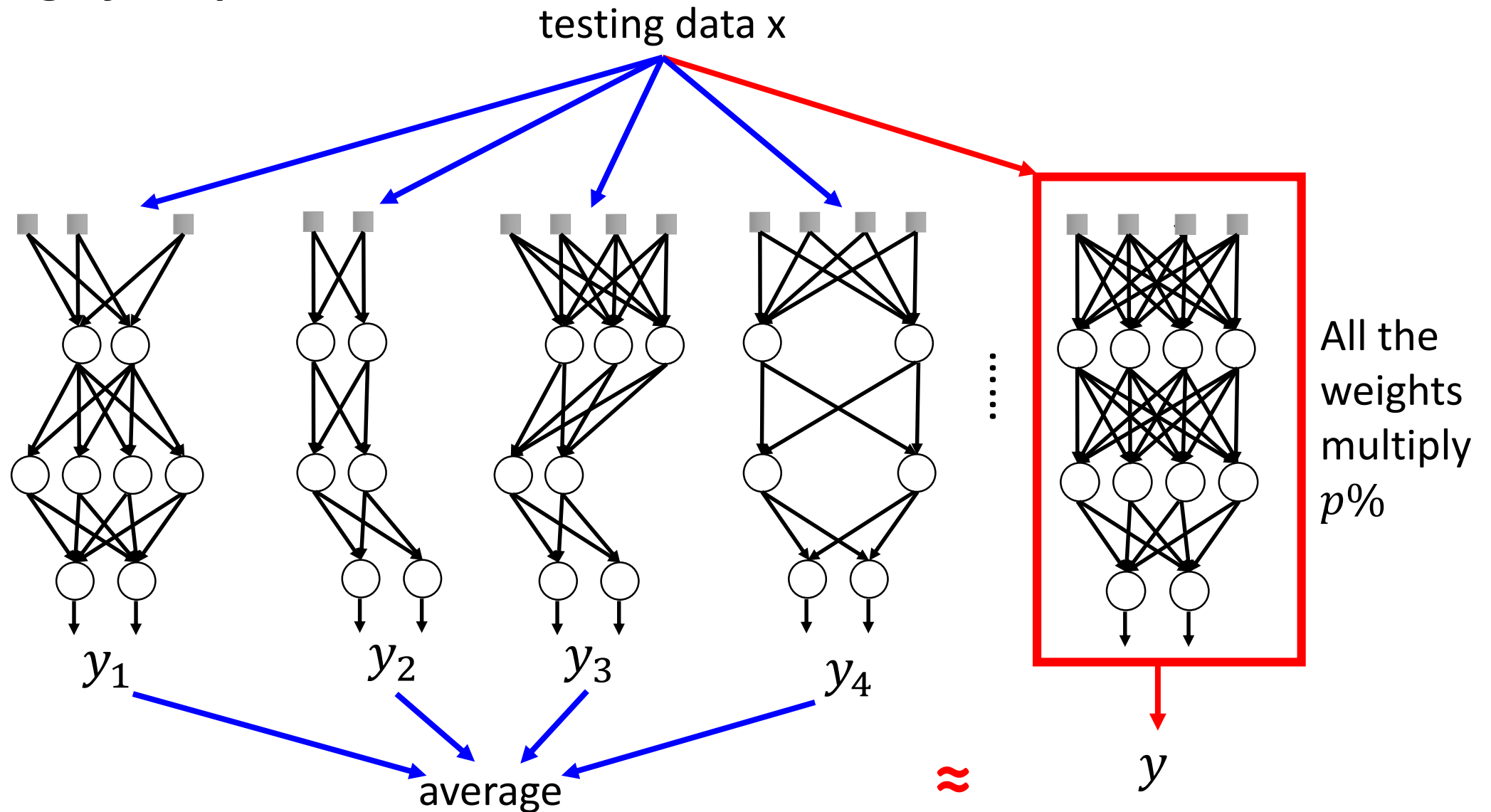


**Ensemble**

Train a bunch of networks with different structures

# Dropout is a kind of ensemble.
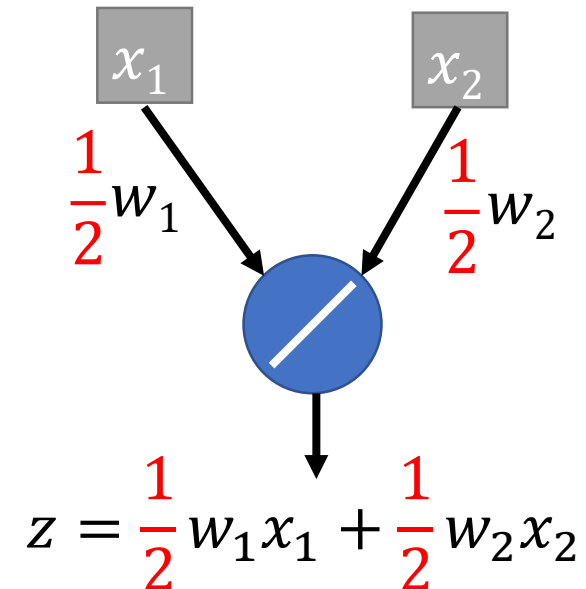
**_Ensemble_**

# Dropout is a kind of ensemble.



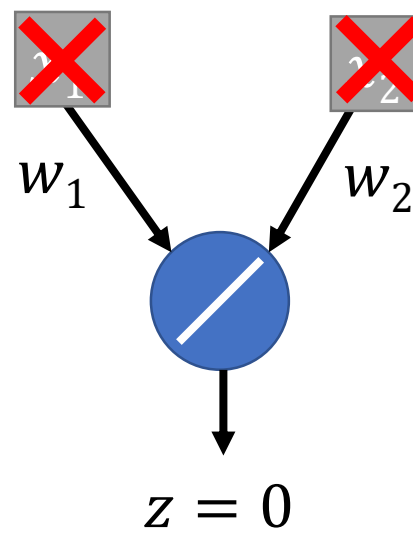*Training of Dropout*
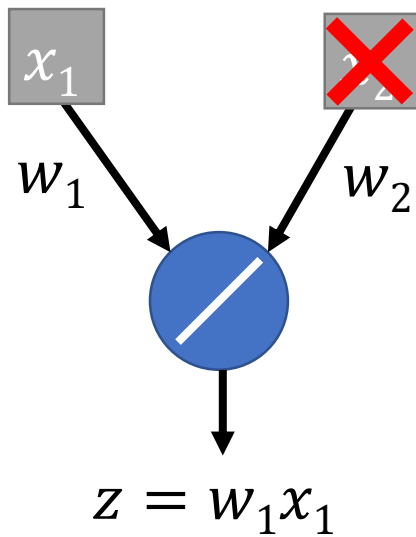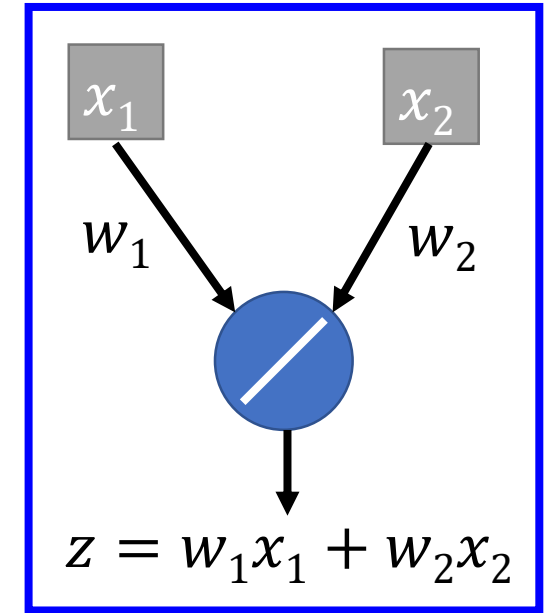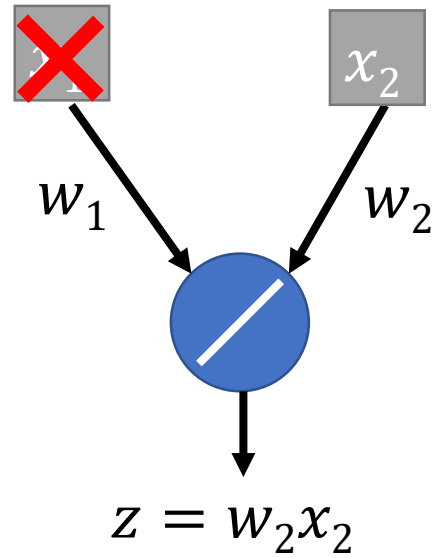
M neurons

$2^M$ possible networks

- Using one minibatch to train one network
- Some parameters in the network are shared

# Dropout is a kind of ensemble.

**_Testing of Dropout_**



testing data x

$y_1$    $y_2$    $y_3$    $y_4$

average

All the weights multiply $p\%$

$\approx$

$y$

**Testing of Dropout**



$z = w_1 x_1 + w_2 x_2$

$z = w_2 x_2$

$z = w_1 x_1 + w_2 x_2$

$z = w_1 x_1$

$z = 0$

$z = \frac{1}{2} w_1 x_1 + \frac{1}{2} w_2 x_2$
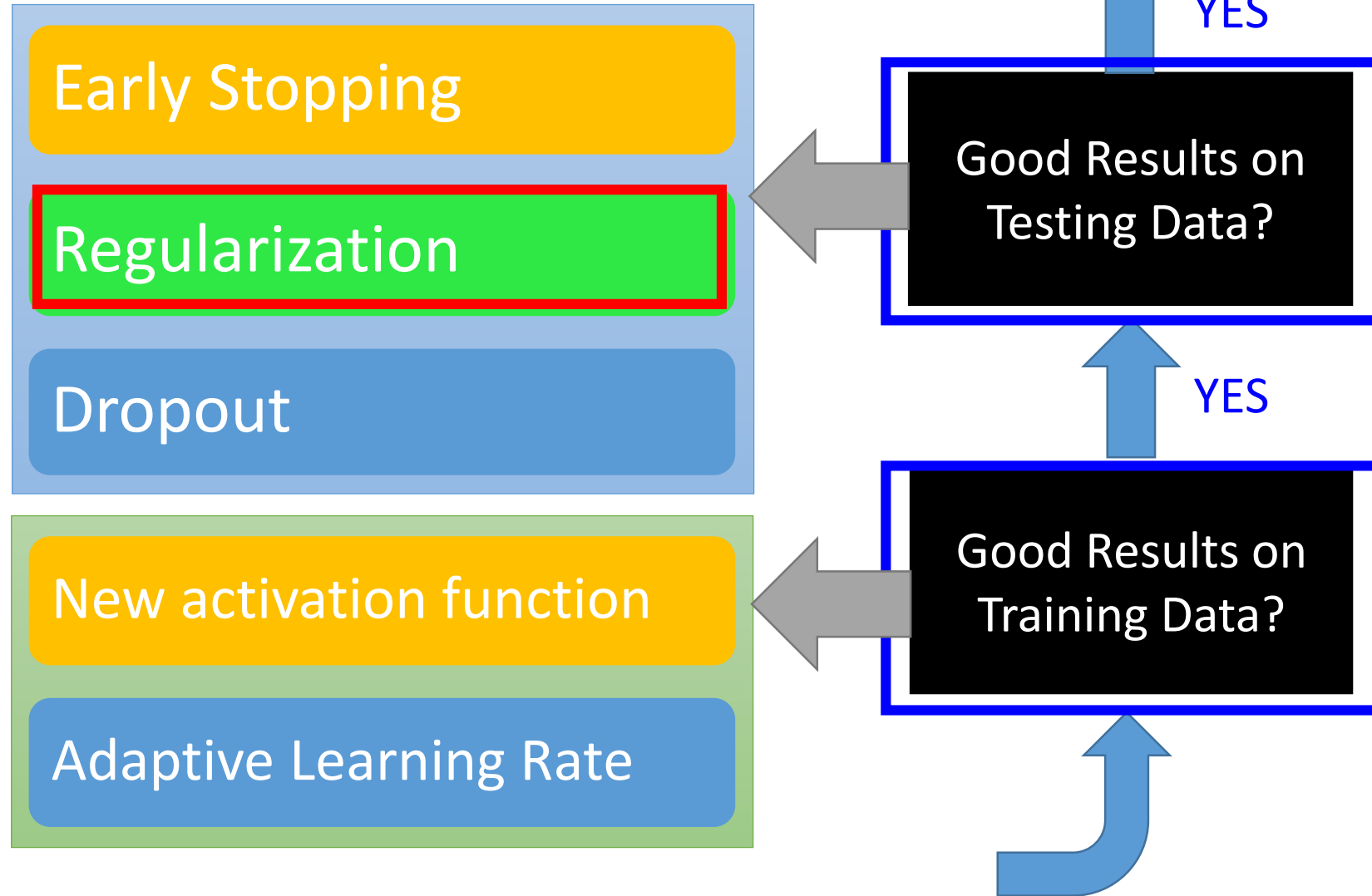
# Recipe of Deep Learning

# Regularization

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

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

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

Original loss
(*e.g.* minimize square error, cross entropy ...)

L2 regularization:
$$\|\theta\|_2 = (w_1)^2 + (w_2)^2 + \cdots$$

(usually not consider biases)

# Regularization

- New loss function to be minimized

$$L'(\theta) = L(\theta) + \lambda \frac{1}{2} \|\theta\|_2$$

Gradient: $\dfrac{\partial L'}{\partial w} = \dfrac{\partial L}{\partial w} + \lambda w$

Update: $w^{(t+1)} \rightarrow w^{(t)} - \eta \dfrac{\partial L'}{\partial w} = w^{(t)} - \eta \left( \dfrac{\partial L}{\partial w} + \lambda w^{(t)} \right)$

$$= \boxed{(1 - \eta\lambda)w^{(t)}} - \eta \dfrac{\partial L}{\partial w}$$

Weight Decay

Closer to zero

# Regularization

- New loss function to be minimized

$$L'(\theta) = L(\theta) + \lambda \frac{1}{2} \|\theta\|_1 \qquad \frac{\partial L'}{\partial w} = \frac{\partial L}{\partial w} + \lambda \operatorname{sgn}(w)$$

Update:

$$w^{(t+1)} \to w^{(t)} - \eta \frac{\partial L'}{\partial w} \qquad = w^{(t)} - \eta \left( \frac{\partial \mathrm{L}}{\partial w} + \lambda \operatorname{sgn}(w^t) \right)$$

$$= w^{(t)} - \eta \frac{\partial \mathrm{L}}{\partial w} - \underline{\eta\lambda \operatorname{sgn}(w^{(t)})} \text{ Always delete}$$

$$w^{(t+1)} = (1 - \eta\lambda)w^{(t)} - \eta \frac{\partial \mathrm{L}}{\partial w} \quad \text{...... L2}$$

# Regularization: L1 vs. L2

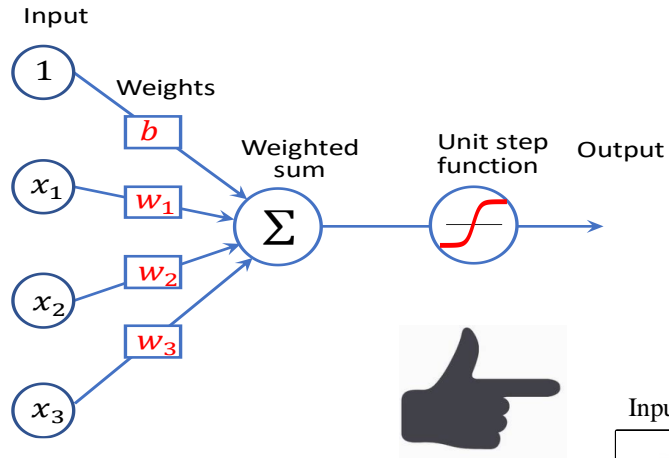L1 regularization: $\|\theta\|_1 = |w_1| + |w_2| + \cdots$

$$w^{(t+1)} = w^{(t)} - \eta \frac{\partial L}{\partial w} - \eta\lambda \operatorname{sgn}\left(w^{(t)}\right)$$

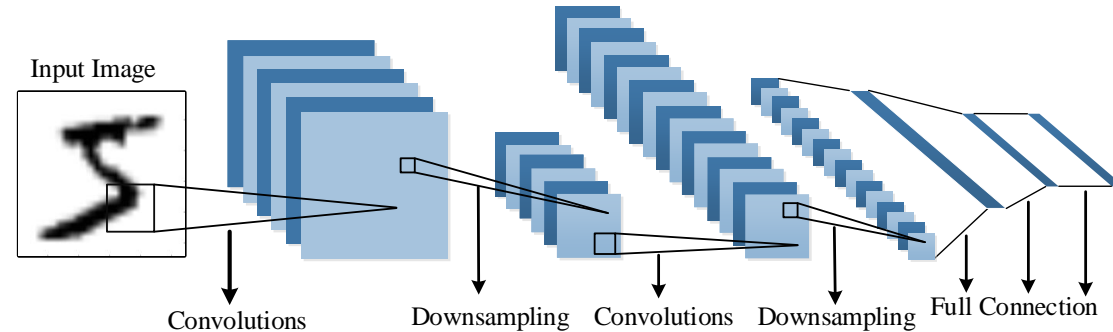L2 regularization: $\|\theta\|_2 = (w_1)^2 + (w_2)^2 + \cdots$

$$w^{(t+1)} = (1 - \eta\lambda)w^{(t)} - \eta \frac{\partial L}{\partial w}$$
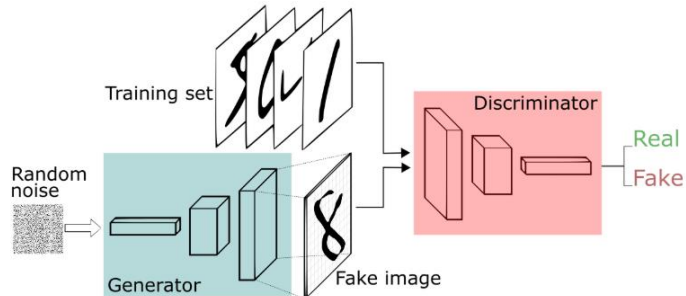
# In this Course

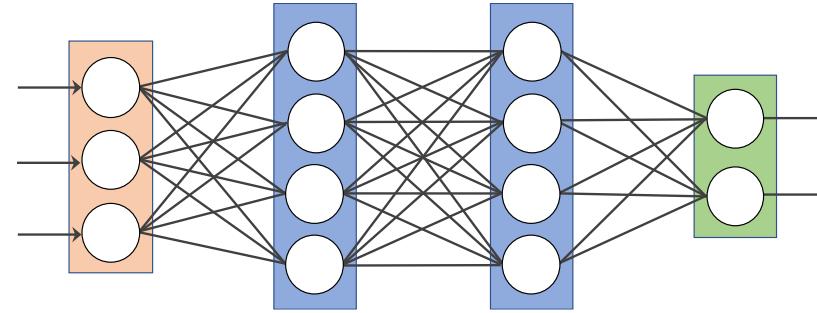**1. DL basics, linear regression, logistic regression etc.**

Input

Weights

1 — b

Weighted sum — Unit step function — Output

$x_1$ — $w_1$

$x_2$ — $w_2$

$x_3$ — $w_3$

$\Sigma$

**2. Multilayer neural networks, backpropagation**

**3. Convolutional Neural Networks and Applications**

Input Image

Convolutions     Downsampling     Convolutions     Downsampling     Full Connection

**4. Generative Adversarial Networks**

Training set

Discriminator

Real

Fake

Random noise

Generator     Fake image

**5. Recurrent networks and applications**

$h_t$

A

$x_t$

=

$h_0$     $h_1$     $h_2$     $h_t$

A  →  A  →  A  →  A

$x_0$     $x_1$     $x_2$     ...     $x_t$