# BIG DATA and AI
## for business
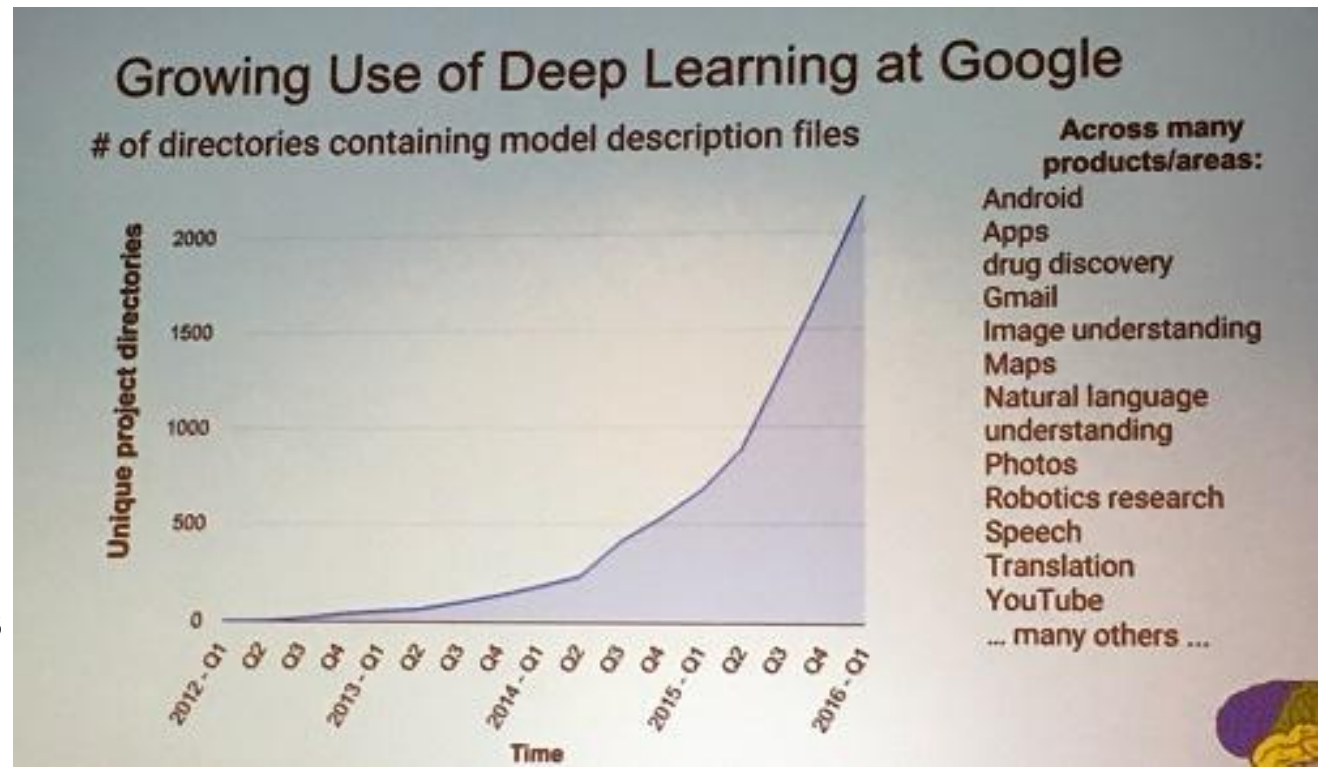
**Deep Learning (1)**

**Decisions, Operations & Information Technologies**
**Robert H. Smith School of Business**
**Fall, 2020**

# Deep learning
# attracts lots of attention.

- I believe you have seen lots of exciting results before.



**Growing Use of Deep Learning at Google**

# of directories containing model description files

**Across many products/areas:**
Android
Apps
drug discovery
Gmail
Image understanding
Maps
Natural language understanding
Photos
Robotics research
Speech
Translation
YouTube
… many others …

Deep learning trends at Google. Source: SIGMOD/Jeff Dean

We mainly focus on the basic techniques.

# Introduction to Deep Learning

# Outline

Introduction of Deep Learning

"Hello World" for Deep Learning

Tips for Deep Learning

# Machine Learning
# ≈ Looking for a Function

- Speech Recognition

$$f\left( \text{} \right) = \text{"How are you"}$$

- Image Recognition

$$f\left( \text{} \right) = \text{"Cat"}$$

- Playing Go

$$f\left( \text{} \right) = \text{"5-5"} \quad \text{(next move)}$$

- Dialogue System

$$f\left( \underset{\text{(what the user said)}}{\text{"Hi"}} \right) = \underset{\text{(system response)}}{\text{"Hello"}}$$

# Framework

$$f(\ \text{[image]}\ ) = \text{"cat"}$$

A set of function

**Model**

$f_1, f_2 \cdots$

$f_1(\ \text{[image]}\ ) = \text{"cat"}$

$f_2(\ \text{[image]}\ ) = \text{"money"}$

$f_1(\ \text{[image]}\ ) = \text{"dog"}$

$f_2(\ \text{[image]}\ ) = \text{"snake"}$

# Image Recognition:

# Framework

$$f( \text{[image of cat]} ) = \text{"cat"}$$

A set of function

Model

$f_1, f_2 \cdots$

$f_1( \text{[image]} ) = \text{"cat"}$    $f_2( \text{[image]} ) = \text{"money"}$

Better!

$f_1( \text{[image]} ) = \text{"dog"}$    $f_2( \text{[image]} ) = \text{"snake"}$

Goodness of function f

Training Data

Supervised Learning

function input:    [monkey image]  [cat image]  [dog image]

function output:    "monkey"    "cat"    "dog"

# Framework

Image Recognition:

$$f\left( \text{🐱} \right) = \text{"cat"}$$



**Step 1** — **Model**: A set of function $f_1, f_2 \cdots$

**Training**  **Testing**

**Step 2** — Goodness of function f

Training Data

Pick the "Best" Function

**Step 3** — $f^*$

"monkey"   "cat"   "dog"

Using $f^*$

"cat"

# Three Steps for Deep Learning

Step 1: define a set of function

Neural Network

Step 2: goodness of function

Step 3: pick the best function

# Neural Network

## *Neuron*

$$z = a_1 w_1 + \cdots + a_k w_k + \cdots + a_K w_K + b$$



$a_1$

$a_k$

$a_K$

$w_1$

$w_k$

$w_K$

weights

A simple function

$+$

$z$

$\sigma(z)$

$a$

Activation function

$b$ bias

# Neural Network

**_Neuron_**

Sigmoid Function $\sigma(z)$

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



2

1

-1

-2

1

-1

+ → 4 → $\sigma(z)$ → 0.98

Activation function

weights

1 bias

# Neural Network



Different connections lead to different network structures



$+$ $\sigma(z)$

$+$ $\sigma(z)$

$+$ $\sigma(z)$

$+$ $\sigma(z)$

$+$ $\sigma(z)$

The neurons have different values of weights and biases.

Weights and biases are network parameters $\theta$

# Fully Connect Feedforward Network



Sigmoid Function

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

# Fully Connect Feedforward Network
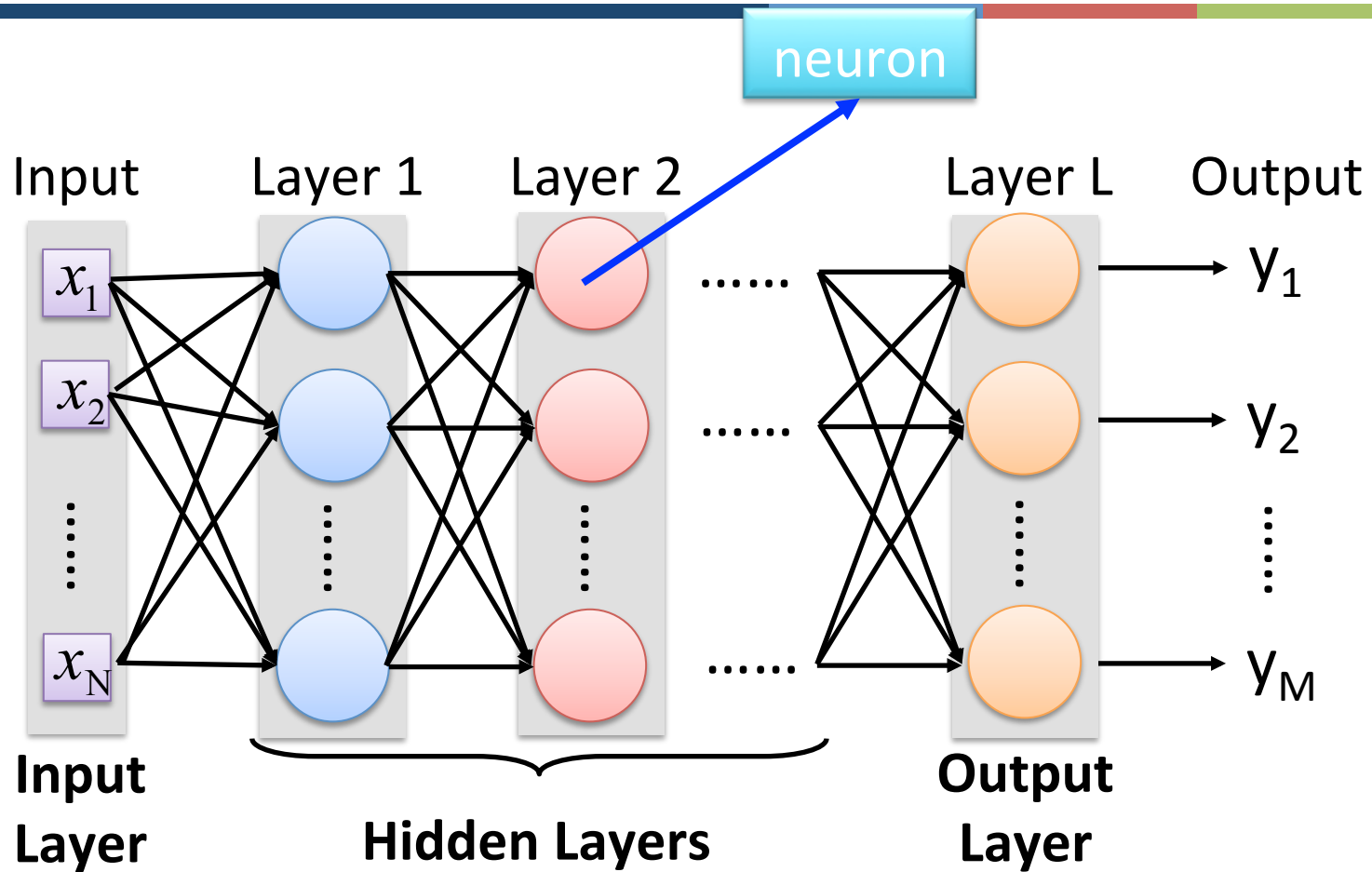
# Fully Connect Feedforward Network



This is a function.

Input vector, output vector

$f\left(\begin{bmatrix} 1 \\ -1 \end{bmatrix}\right) = \begin{bmatrix} 0.62 \\ 0.83 \end{bmatrix}$  $f\left(\begin{bmatrix} 0 \\ 0 \end{bmatrix}\right) = \begin{bmatrix} 0.51 \\ 0.8 \end{bmatrix}$

Given parameters $\theta$, define a function

Given network structure, define **_a function set_**

# Fully Connect Feedforward Network



neuron

Input  Layer 1  Layer 2  Layer L  Output

$x_1$  $x_2$  $x_N$  $y_1$  $y_2$  $y_M$

**Input Layer**

**Hidden Layers**

**Output Layer**

Deep means many hidden layers

# Why Deep? Universality Theorem

Any continuous function f

$$f : R^N \rightarrow R^M$$

Can be realized by a network with one hidden layer

(given **enough** hidden neurons)



Reference for the reason:
http://
neuralnetworksanddeeplearning.
com/chap4.html

Why "Deep" neural network not "Fat" neural network?

# *Why Deep? Analogy*

**Logic circuits**

- Logic circuits consists of **gates**
- **A two layers of logic gates** can represent **any Boolean function.**
- Using multiple layers of logic gates to build some functions are much simpler

**Neural network**

- Neural network consists of **neurons**
- **A hidden layer network** can represent **any continuous function.**
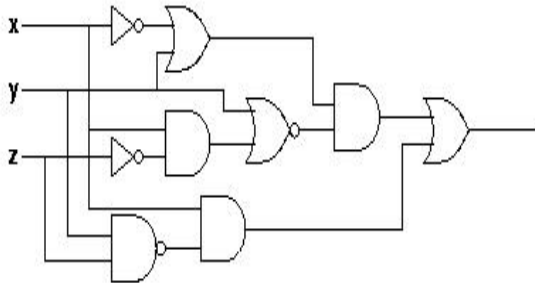- Using multiple layers of neurons to represent some functions are much simpler
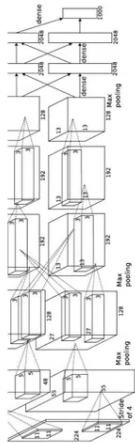
➡ **less gates needed**

➡ **less parameters** ➡ **less data?**



More reason: https://www.youtube.com/watch?v=XsC9byQkUH8&list=PLJV_el3uVTsPy9oCRY30oBPNLCo89yu49&index=13

# Deep = Many hidden layers

http://cs231n.stanford.edu/
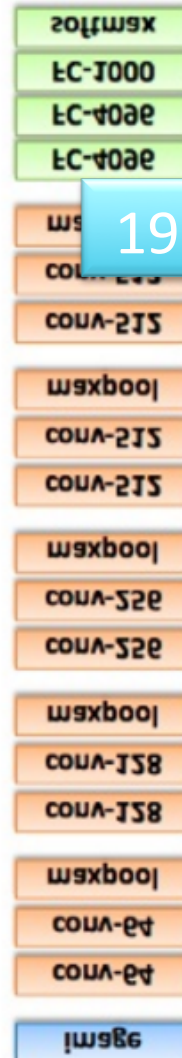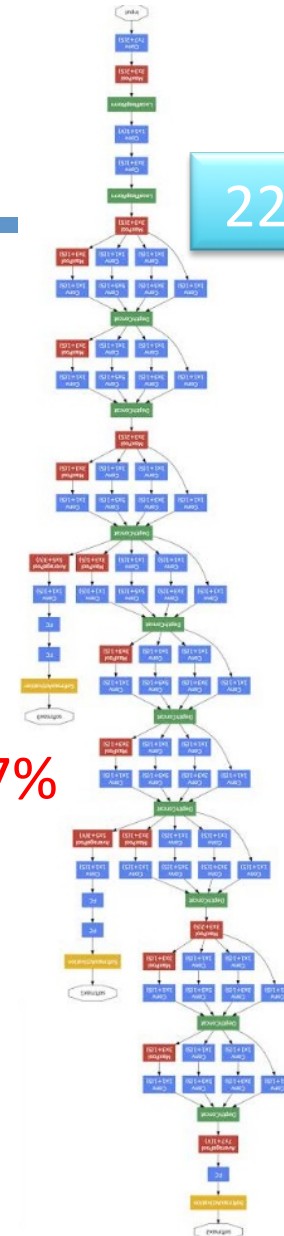slides/
winter1516_lecture8.pdf



22 layers

19 layers

8 layers

6.7%

7.3%

16.4%

AlexNet (2012)

VGG (2014)

GoogleNet (2014)

# Deep = Many hidden layers



Special structure

101 layers

152 layers

3.57%

16.4%

7.3%

6.7%

AlexNet (2012)     VGG (2014)     GoogleNet (2014)     Residual Net (2015)     Taipei 101

# Output Layer

- Softmax layer as the output layer

**_Ordinary Layer_**

$z_1 \longrightarrow \sigma \longrightarrow y_1 = \sigma(z_1)$

$z_2 \longrightarrow \sigma \longrightarrow y_2 = \sigma(z_2)$

$z_3 \longrightarrow \sigma \longrightarrow y_3 = \sigma(z_3)$

In general, the output of network can be any value.
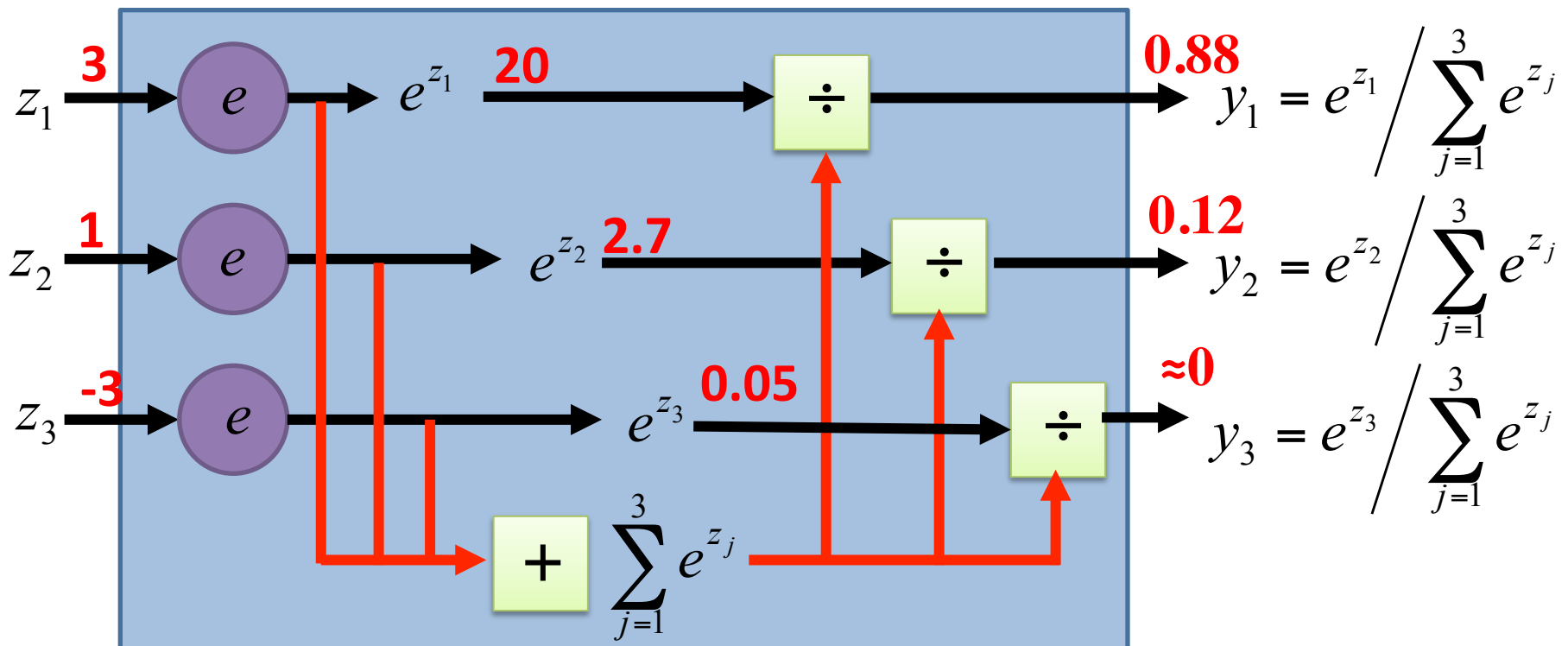
May not be easy to interpret

# Output Layer

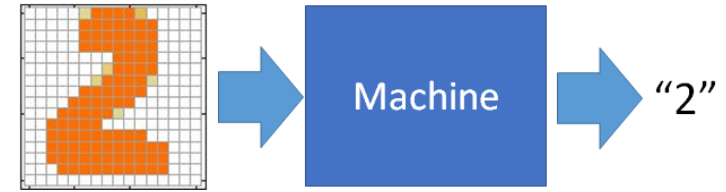- Softmax layer as the output layer

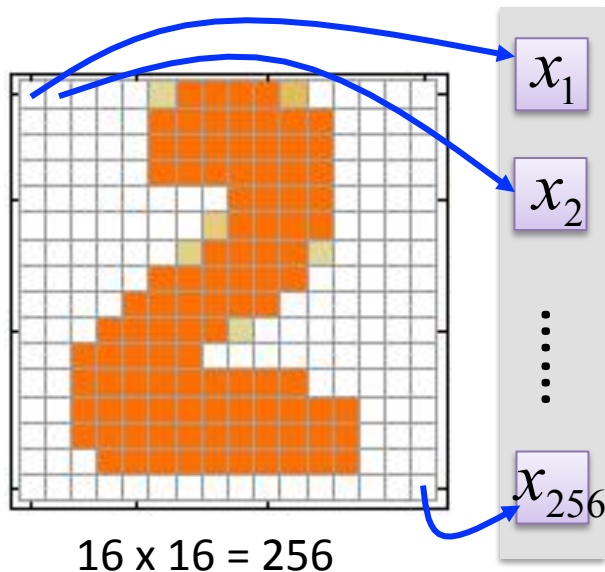**Probability**:
- $1 > y_i > 0$
- $\sum_i y_i = 1$

**Softmax Layer**



$z_1$ → **3** → $e$ → $e^{z_1}$ → **20** → ÷ → **0.88** → $y_1 = e^{z_1} \bigg/ \sum_{j=1}^{3} e^{z_j}$

$z_2$ → **1** → $e$ → $e^{z_2}$ → **2.7** → ÷ → **0.12** → $y_2 = e^{z_2} \bigg/ \sum_{j=1}^{3} e^{z_j}$

$z_3$ → **-3** → $e$ → $e^{z_3}$ → **0.05** → ÷ → **≈0** → $y_3 = e^{z_3} \bigg/ \sum_{j=1}^{3} e^{z_j}$

$+ \quad \sum_{j=1}^{3} e^{z_j}$

# Example Application



"2"

## Input



$x_1$

$x_2$

$\vdots$

$x_{256}$

16 x 16 = 256

Ink → 1
No ink → 0

## Output

0.1 — is 1

0.7 — is 2

$\vdots$

0.2 — is 0

The image is "2"

Each dimension represents the confidence of a digit.

# Example Application

- Handwriting Digit Recognition



$x_1$

$x_2$

$\vdots$

$x_{256}$

Neural Network

What is needed is a function ……

$y_1$ — is 1

$y_2$ — is 2

$\vdots$

$y_{10}$ — is 0

Input:
256-dim vector

output:
10-dim vector

# Example Application



You need to decide the network structure to let a good function in your function set.

- Q: How many layers? How many neurons for each layer?

  Trial and Error + Intuition

- Q: Can we design the network structure?

  Convolutional Neural Network (CNN) in the next lecture

- Q: Can the structure be automatically determined?
  - Yes, but not widely studied yet.

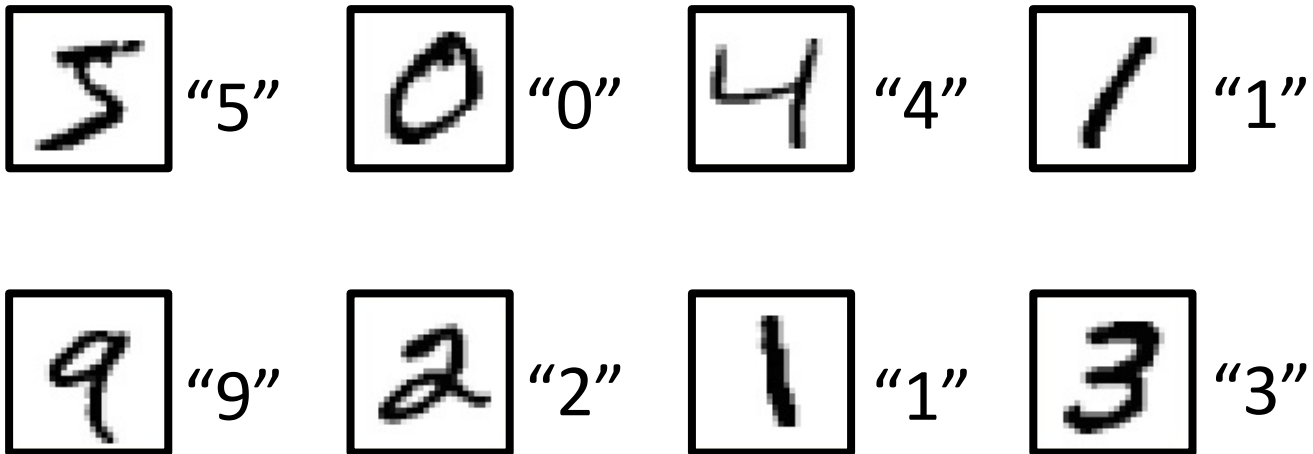# Three Steps for Deep Learning

Step 1: define a set of function

↓

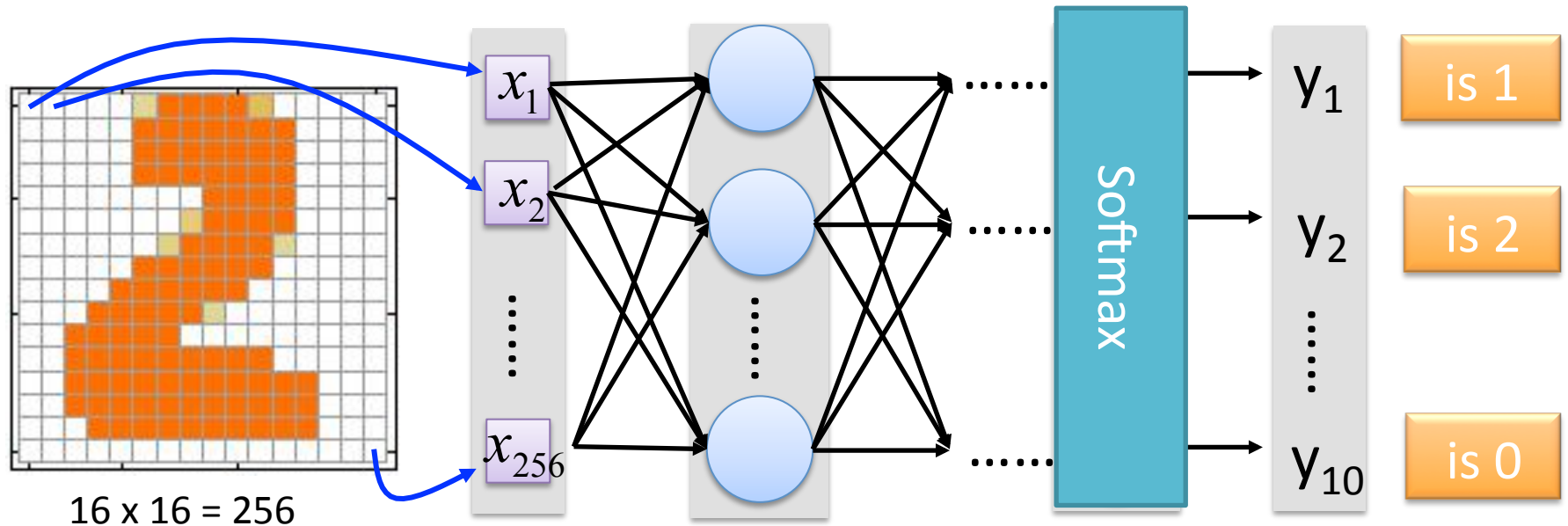**Step 2: goodness of function**

↓

Step 3: pick the best function

# Training Data

- Preparing training data: images and their labels



"5"  "0"  "4"  "1"

"9"  "2"  "1"  "3"

The learning target is defined on the training data.

# Learning Target



16 x 16 = 256

Ink → 1
No ink → 0

The learning target is ......

Input: $y_1$ has the maximum value

Input: $y_2$ has the maximum value

A good function should make the loss of all examples as small as possible.

"1"

$x_1$

$x_2$

$x_{256}$

Given a set of parameters

Softmax

$y_1$

$y_2$

$y_{10}$

As close as possible

Loss $l$

1

0

0

target

Loss can be **square error** or **cross entropy** between the network output and target

# Total Loss

$$L = \sum_{r=1}^{R} l_r$$

For all training data ...

Total Loss:



**For all training data diagram:**

$x^1 \rightarrow$ NN $\rightarrow y^1 \longleftrightarrow \hat{y}^1$, $l_1$

$x^2 \rightarrow$ NN $\rightarrow y^2 \longleftrightarrow \hat{y}^2$, $l_2$

$x^3 \rightarrow$ NN $\rightarrow y^3 \longleftrightarrow \hat{y}^3$, $l_3$

$x^R \rightarrow$ NN $\rightarrow y^R \longleftrightarrow \hat{y}^R$, $l_R$

As small as possible

Find ***a function in function set*** that minimizes total loss L

Find ***the network parameters*** $\theta^*$ that minimize total loss L

# Three Steps for Deep Learning

Step 1: define a set of function

Step 2: goodness of function

Step 3: pick the best function
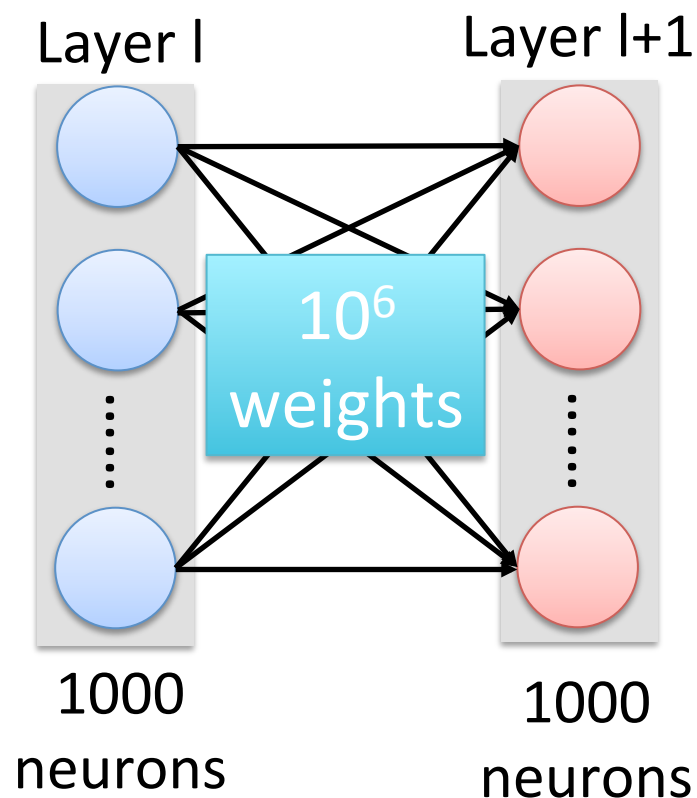
# How to pick the best function

Find **_network parameters_** $\theta^*$ that minimize total loss L

Enumerate all possible values

Network parameters $\theta = \{w_1, w_2, w_3, \cdots, b_1, b_2, \cdots\}$

Millions of parameters

E.g. speech recognition: 8 layers and 1000 neurons each layer

Layer l          Layer l+1

$10^6$ weights

1000 neurons          1000 neurons

# Gradient Descent

Network parameters $\theta = \{w_1, w_2, \ldots, b_1, b_2, \ldots\}$

Find **_network parameters $\theta^*$_** that minimize total loss L
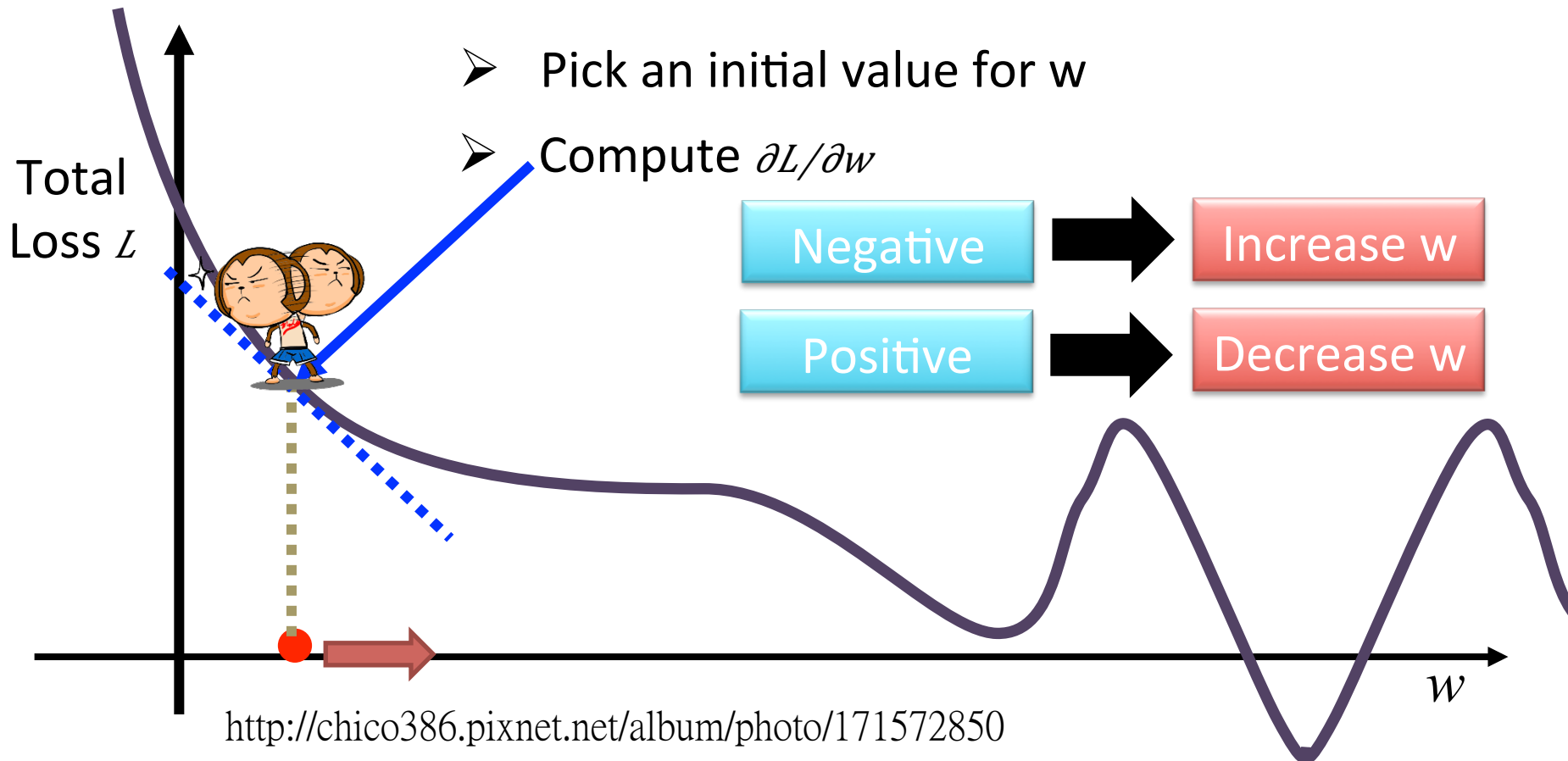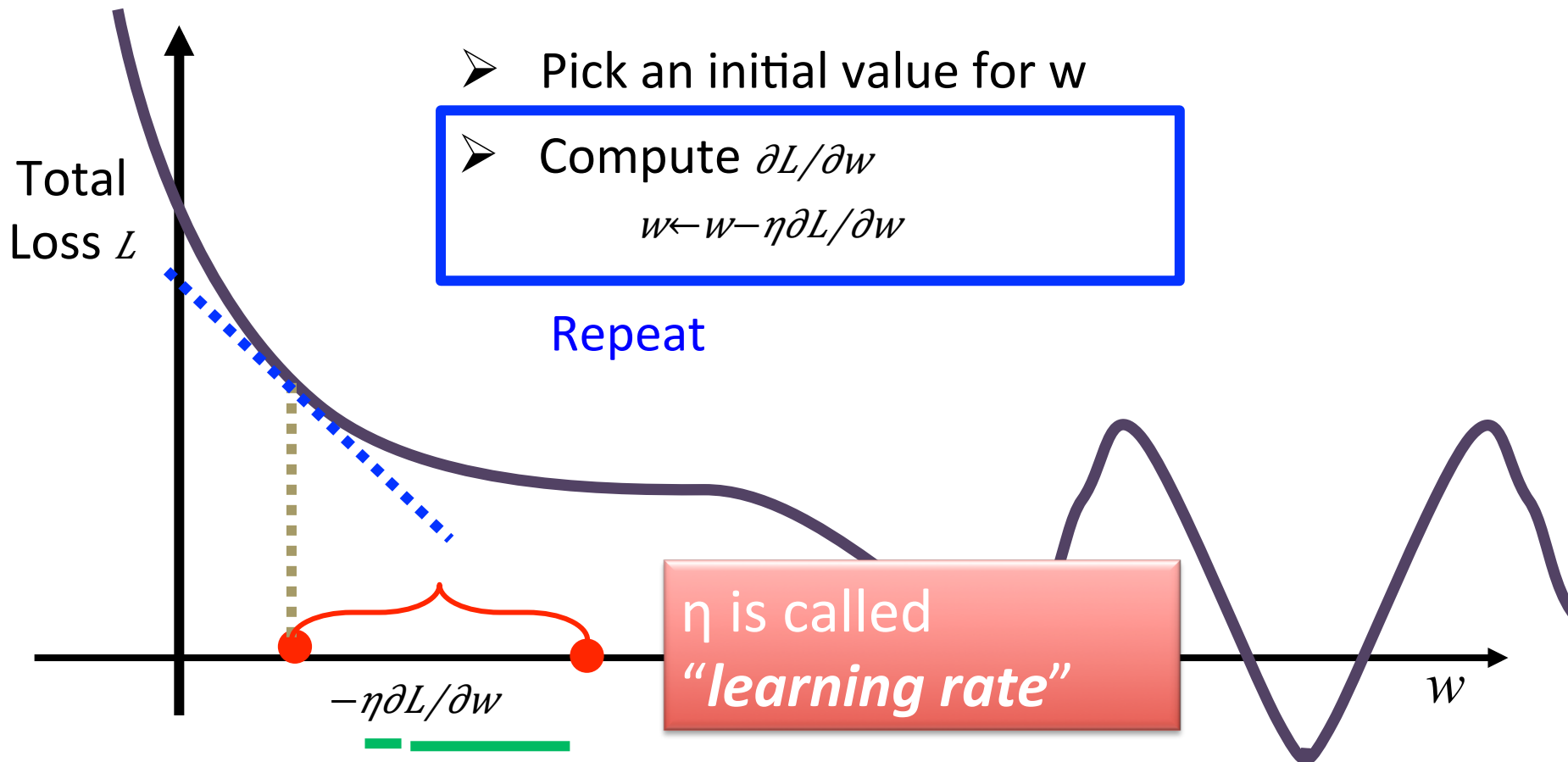
➤ Pick an initial value for w

Random, RBM pre-train

Usually good enough
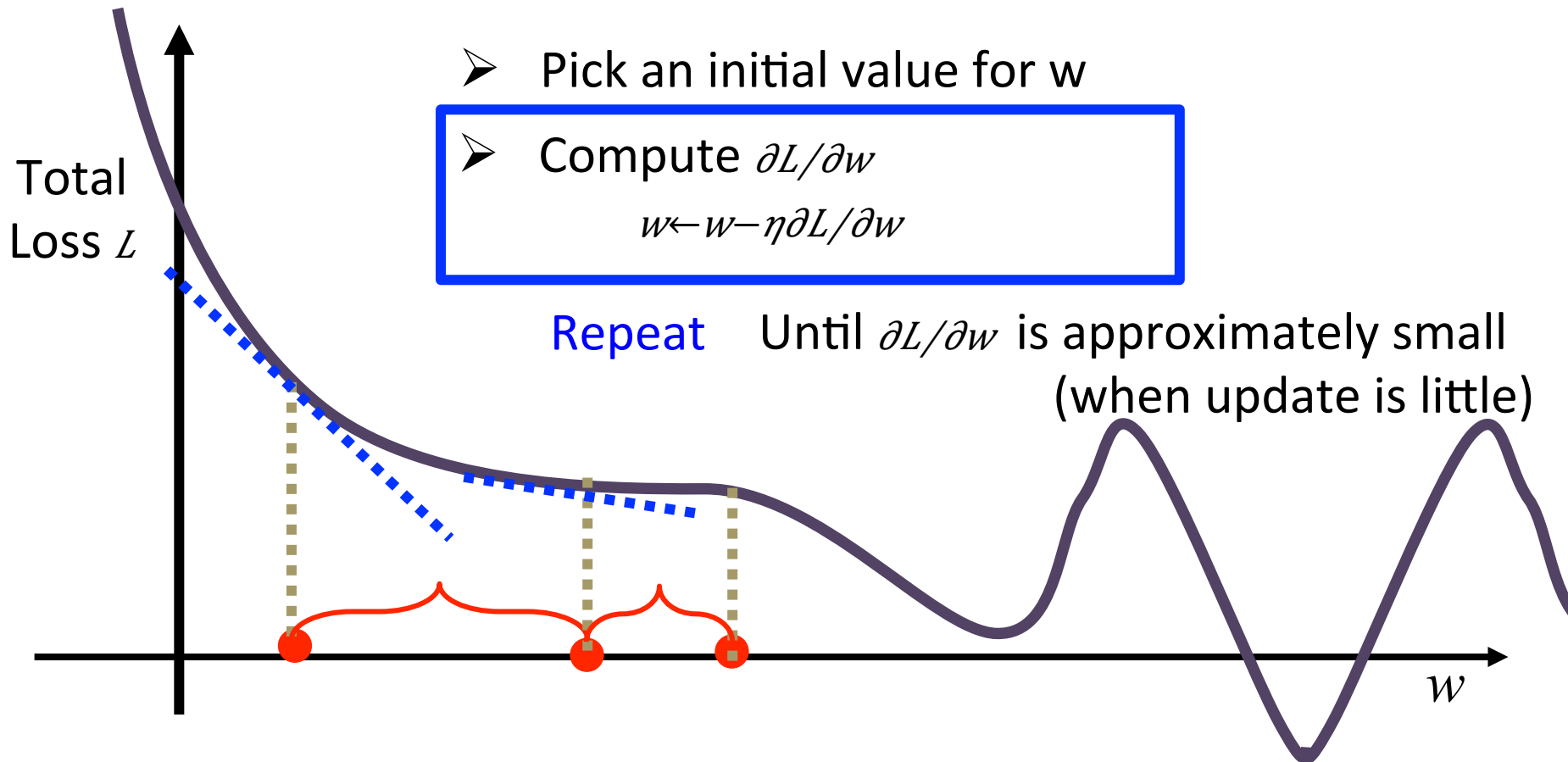
Total Loss $L$

$w$

# Gradient Descent

Network parameters $\theta=\{w_1, w_2, \ldots, b_1, b_2, \ldots\}$

## Find **_network parameters $\theta^*$_** that minimize total loss L

➢ Pick an initial value for w

➢ Compute $\partial L / \partial w$

Total Loss $L$

| Negative | → | Increase w |
| Positive | → | Decrease w |

$w$

http://chico386.pixnet.net/album/photo/171572850

# Gradient Descent

Network parameters $\theta=\{w_1, w_2, \ldots, b_1, b_2, \ldots\}$

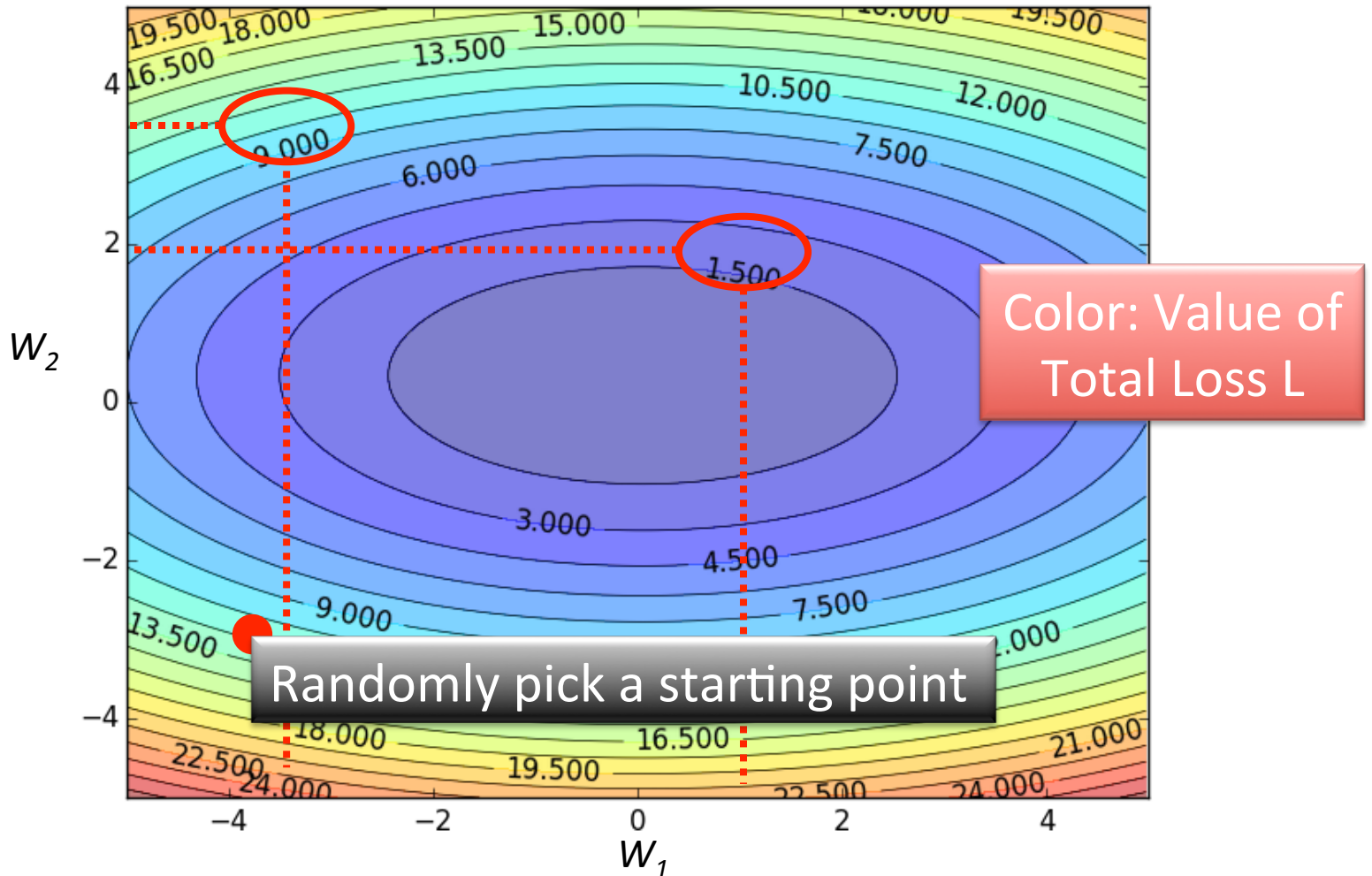Find **_network parameters $\theta^*$_** that minimize total loss L

Total Loss $L$

➢ Pick an initial value for w

➢ Compute $\partial L/\partial w$

$$w \leftarrow w - \eta \partial L/\partial w$$

Repeat

$-\eta \partial L/\partial w$

η is called **_"learning rate"_**

$w$

# Gradient Descent

Network parameters $\theta=\{w_1, w_2, \ldots, b_1, b_2, \ldots\}$

Find **_network parameters $\theta^*$_** that minimize total loss L

Total Loss $L$

➢ Pick an initial value for w

➢ Compute $\partial L/\partial w$

$$w \leftarrow w - \eta \partial L/\partial w$$

Repeat    Until $\partial L/\partial w$ is approximately small
(when update is little)
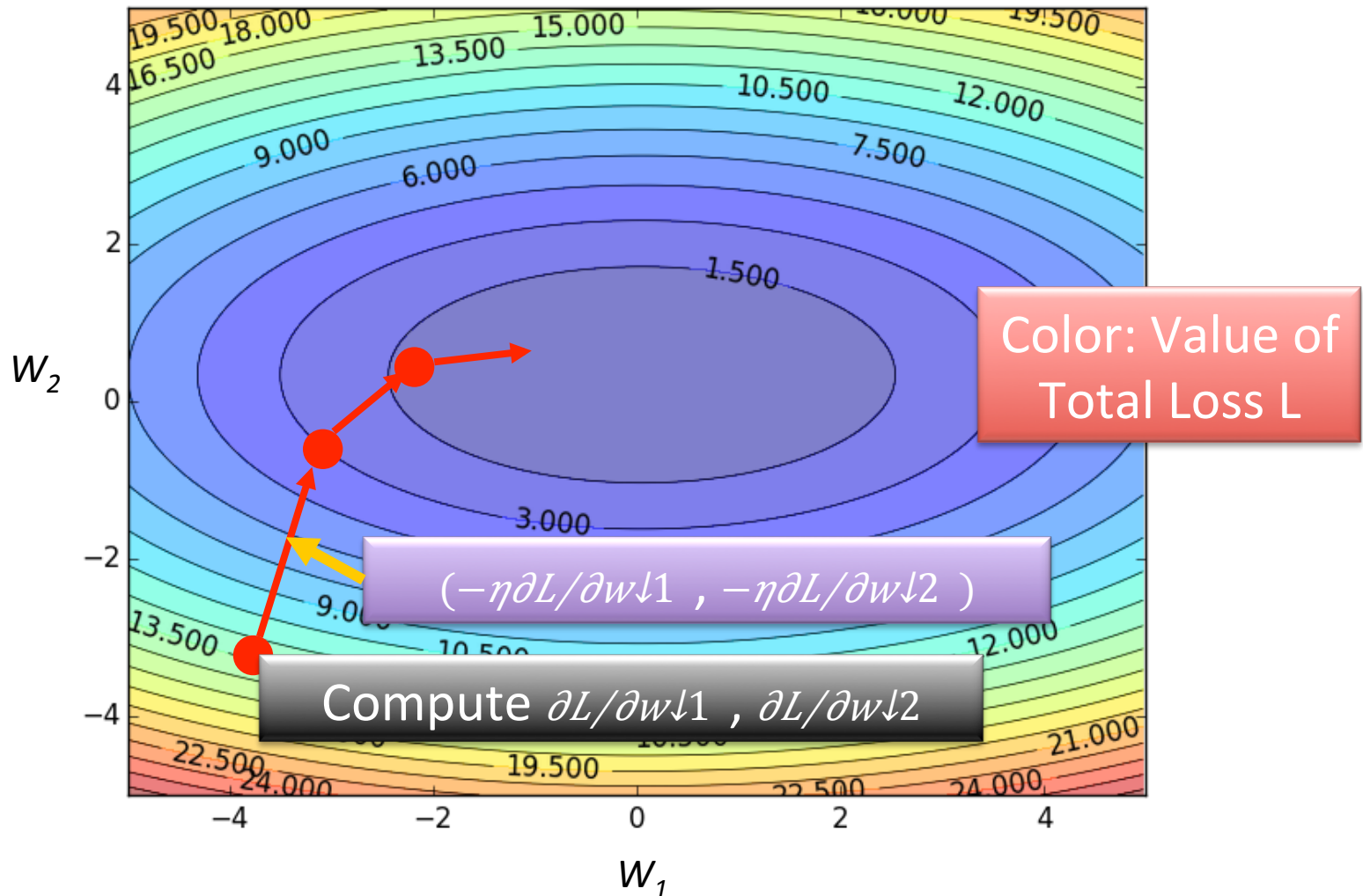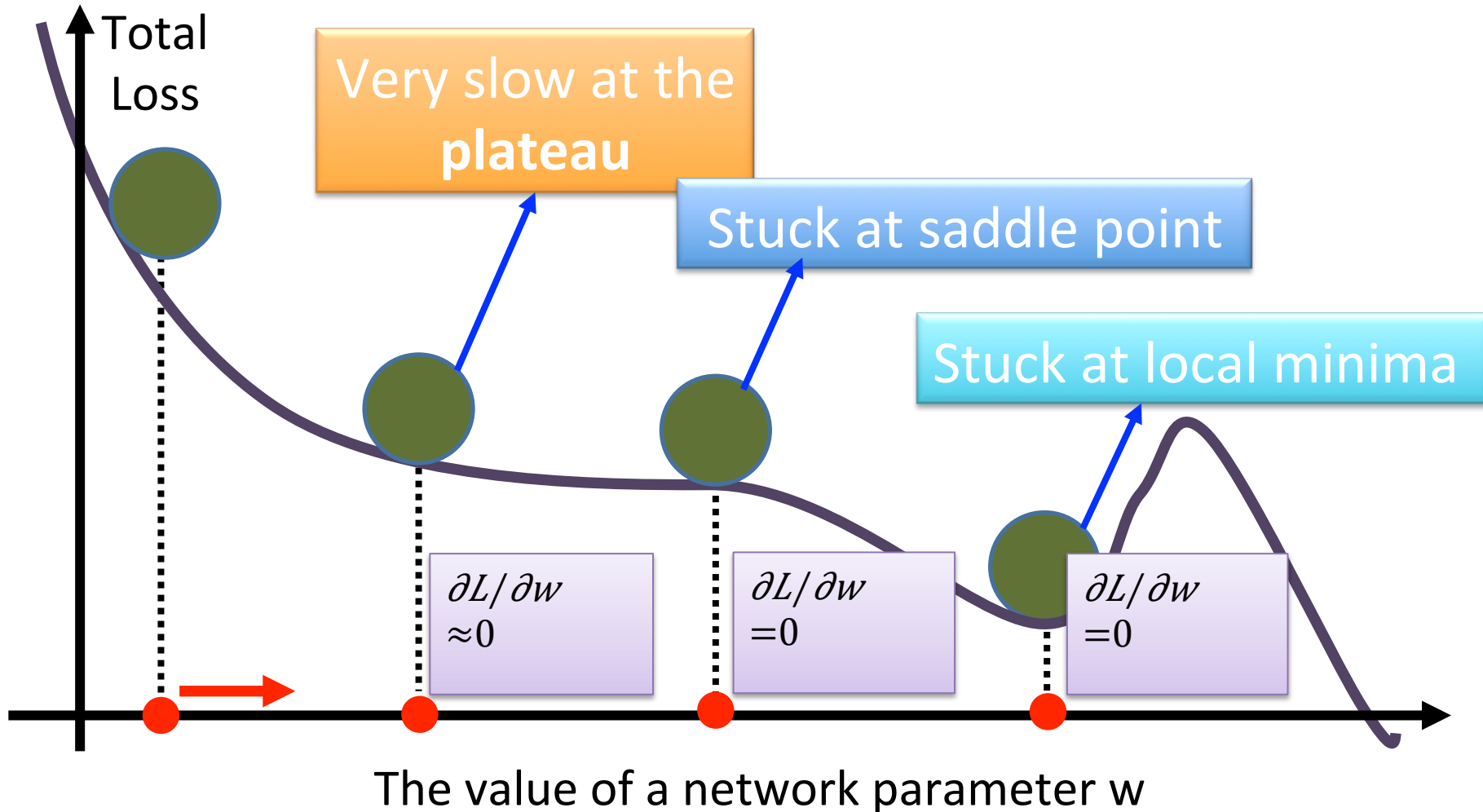
$w$

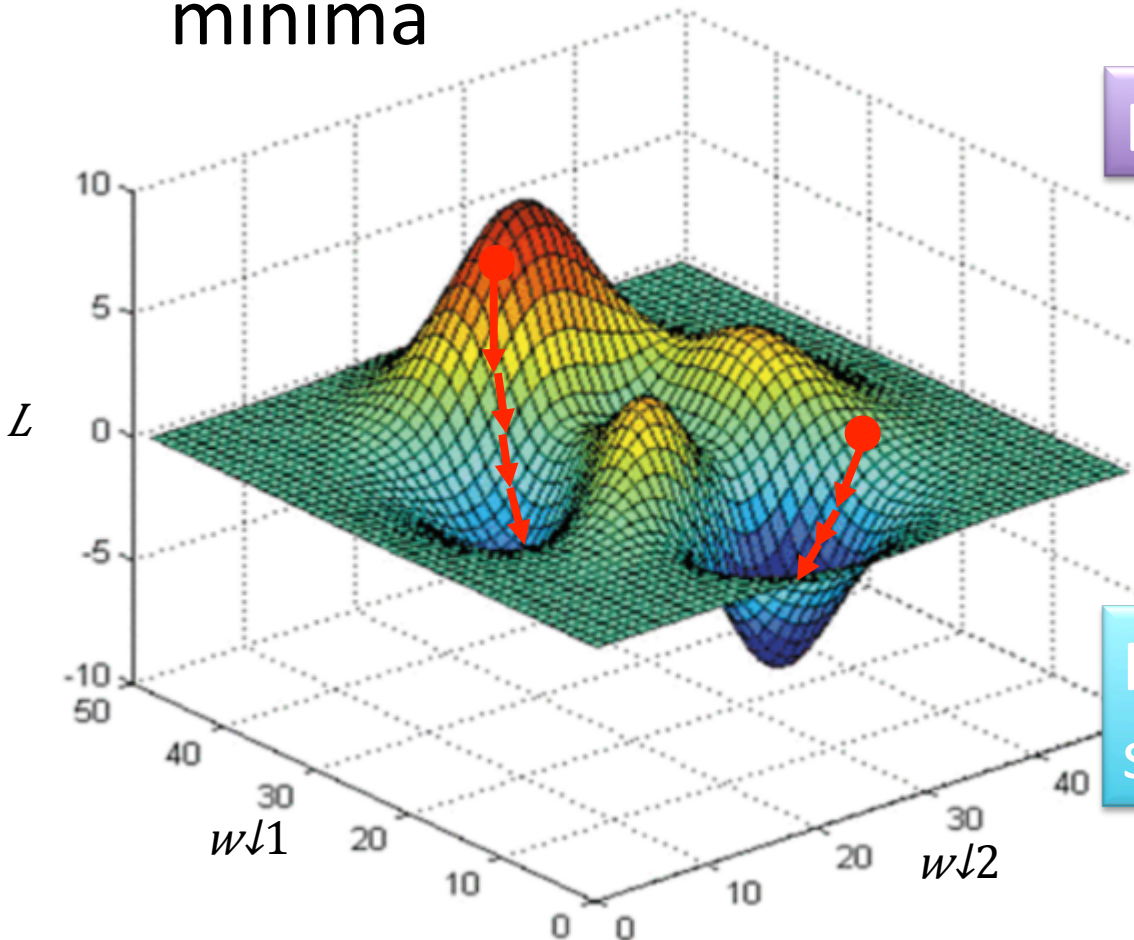# Gradient Descent

# Gradient Descent

Hopfully, we would reach a minima …..



Color: Value of Total Loss L

$(-\eta \partial L / \partial w\!\downarrow\!1 \ , -\eta \partial L / \partial w\!\downarrow\!2 \ )$

Compute $\partial L / \partial w\!\downarrow\!1 \ , \partial L / \partial w\!\downarrow\!2$

# Local Minima



Total Loss

Very slow at the **plateau**

Stuck at saddle point

Stuck at local minima

$\partial L / \partial w \approx 0$

$\partial L / \partial w = 0$

$\partial L / \partial w = 0$

The value of a network parameter w

# Local Minima

- Gradient descent never guarantee global minima



**Different initial point**

↓

**Reach different minima, so different results**

$L$

$w{\downarrow}1$

$w{\downarrow}2$

# Gradient Descent

This is the "learning" of machines in deep learning ......

➡️ Even alpha go using this approach.

People image ......



Actually .....



I hope you are not too disappointed :p

# Backpropagation

- Backpropagation: an efficient way to compute $\partial L / \partial w$ in neural network

# Three Steps for Deep Learning

| Step 1: define a set of function | → | Step 2: goodness of function | → | Step 3: pick the best function |
|---|---|---|---|---|

Deep Learning is so simple ……

Now If you want to find a function

If you have lots of function input/output (?) as training data

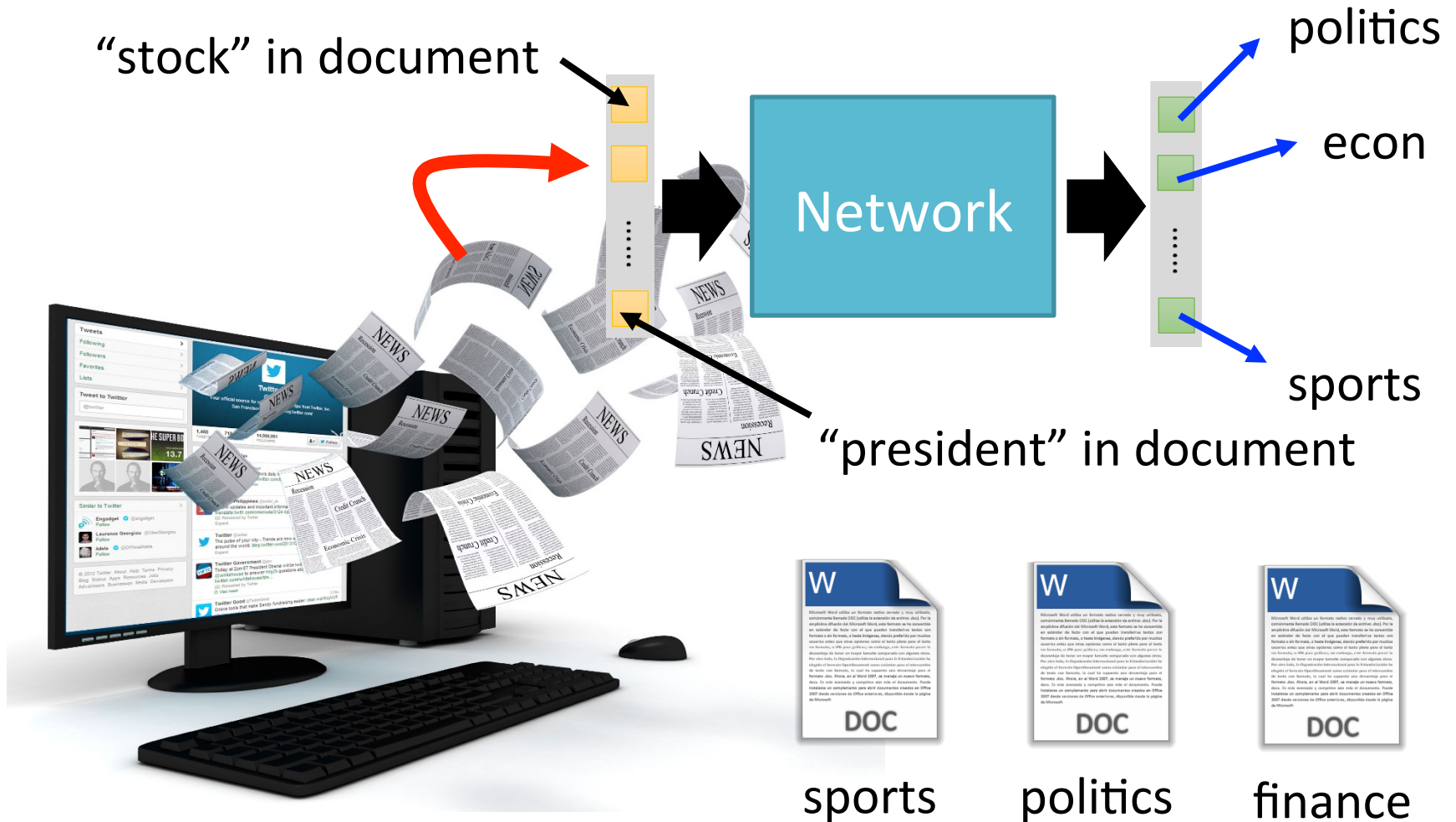➡ You can use deep learning

# For example, you can do .......

- Image Recognition

# For example, you can do …….



"Talk" in e-mail

*Spam filtering*

EMAIL

Network

1/0

(Yes/No)

"free" in e-mail

1 (Yes)

0 (No)

SPAM

SPAM FILTER

(http://spam-filter-review.toptenreviews.com/)

# For example, you can do …….



"stock" in document

"president" in document

Network

politics

econ

sports

http://top-breaking-news.com/

sports    politics    finance

# Outline

Introduction of Deep Learning

"Hello World" for Deep Learning

Tips for Deep Learning

# Keras

TensorFlow or theano

Very flexible

Need some effort to learn

Interface of TensorFlow or Theano

**K**

keras

Easy to learn and use

(still have some flexibility)

You can modify it if you can write TensorFlow or Theano

# Keras

- François Chollet is the author of Keras.

- Keras means *horn* in Greek

- Documentation: [http://keras.io/](http://keras.io/)

- Example: https://github.com/fchollet/keras/tree/master/examples

# PyTorch

- Tutorial: https://pytorch.org/tutorials

# Example Application

- Handwriting Digit Recognition



28 x 28

MNIST Data: http://yann.lecun.com/exdb/mnist/

"Hello world" for deep learning

Keras provides data sets loading function: http://keras.io/datasets/

# Keras

28x28

500

500

Softmax

$y_1$    $y_2$......    $y_{10}$

```
model = Sequential()
```

```
model.add( Dense( input_dim=28*28,
                  output_dim=500 ))
model.add( Activation('sigmoid') )
```

```
model.add( Dense( output_dim=500 ) )
model.add( Activation('sigmoid') )
```

```
model.add( Dense(output_dim=10 ) )
model.add( Activation('softmax') )
```

# Keras

# Keras

Step 1: define a set of function → Step 2: goodness of function → **Step 3: pick the best function**

Step 3.1: Configuration

```
model.compile(loss='mse',
              optimizer=SGD(lr=0.1),
              metrics=['accuracy'])
```
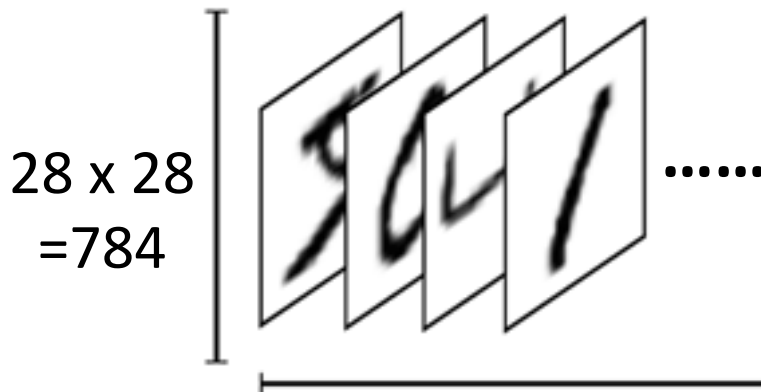
$$w \leftarrow w - \eta \partial L / \partial w$$

0.1

Step 3.2: Find the optimal network parameters

```
model.fit(x_train, y_train, batch_size=100, nb_epoch=20)
```

Training data (Images)

Labels (digits)

# Keras

Step 1: define a set of function → Step 2: goodness of function → **Step 3: pick the best function**

Step 3.2: Find the optimal network parameters

```
model.fit(x_train, y_train, batch_size=100, nb_epoch=20)
```
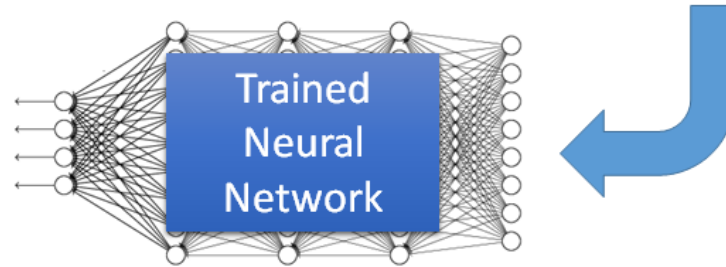
numpy array

numpy array

28 x 28 =784

10

Number of training examples

Number of training examples

Save and load models

http://keras.io/getting-started/faq/#how-can-i-save-a-keras-model

How to use the neural network (testing):

case 1:
```
score = model.evaluate(x_test,y_test)
print('Total loss on Testing Set:', score[0])
print('Accuracy of Testing Set:', score[1])
```

case 2:
```
result = model.predict(x_test)
```
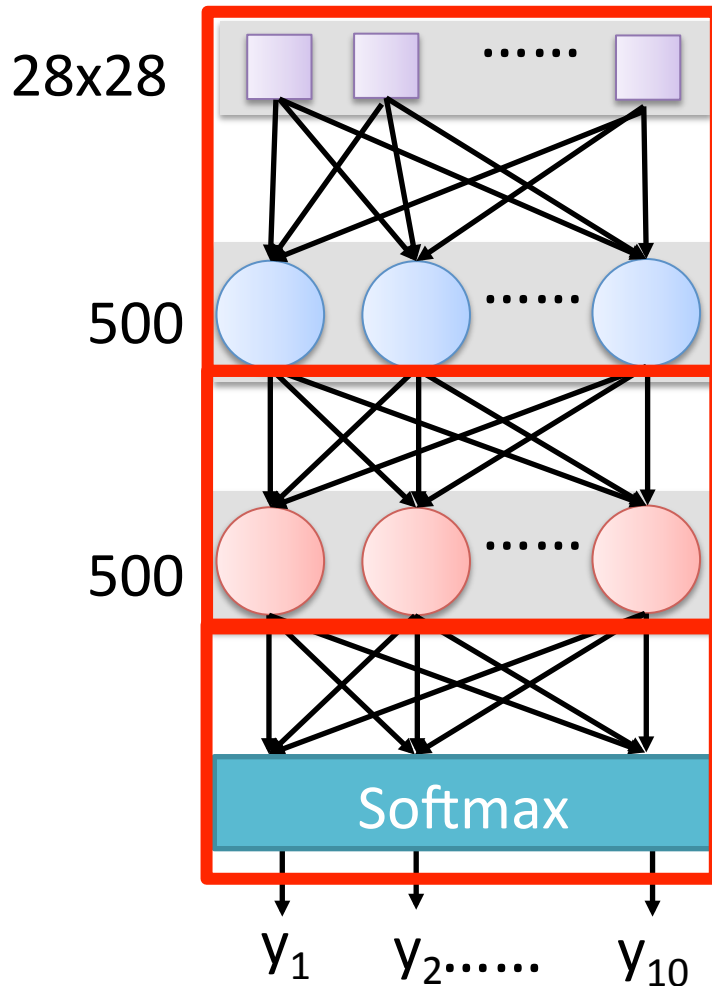
# PyTorch

28x28

500

500

Softmax

$y_1$  $y_2$......  $y_{10}$

```python
import torch.nn as nn
import torch.nn.functional as F

class MyNetwork(nn.Module):
    def __init__(self):
        super(MyNetwork, self).__init__()

        self.fc1 = nn.Linear(28 * 28, 500)
        self.fc2 = nn.Linear(500, 500)
        self.fc3 = nn.Linear(500, 10)

    def forward(self, x):
        x = F.sigmoid(self.fc1(x))
        x = F.sigmoid(self.fc2(x))
        x = self.fc3(x)
        return F.log_softmax(x)
```

PyTorch

Step 1: define a set of function → Step 2: goodness of function → Step 3: pick the best function

"1"

$x_1$ $x_2$ ⋮ $x_{256}$ → Softmax → $y_1$ $y_2$ ⋮ $y_{10}$

difference

Loss $l$

1 0 ⋮ 0

```
net = MyNetwork()
optimizer = torch.optim.SGD(net.parameters(), lr=0.1, momentum=0.9)
criterion = nn.MSELoss()
```

# PyTorch

| Step 1: define a set of function | → | Step 2: goodness of function | → | Step 3: pick the best function |
|---|---|---|---|---|

## Step 3.1: Training

```
for epoch in range(epochs):
    for batch_idx, (data, target) in enumerate(train_loader):
        data, target = Variable(data), Variable(target)
        # resize data from (batch_size, 1, 28, 28) to (batch_size, 28*28)
        data = data.view(-1, 28*28)

        optimizer.zero_grad()
        net_out = net(data)
        loss = criterion(net_out, target)
        loss.backward()
        optimizer.step()

        if batch_idx % log_interval == 0:
            print('Train Epoch: {} [{}/{} ({:.0f}%)]\tLoss: {:.6f}'.format(
                epoch, batch_idx * len(data), len(train_loader.dataset),
                    100. * batch_idx / len(train_loader), loss.data[0]))
```
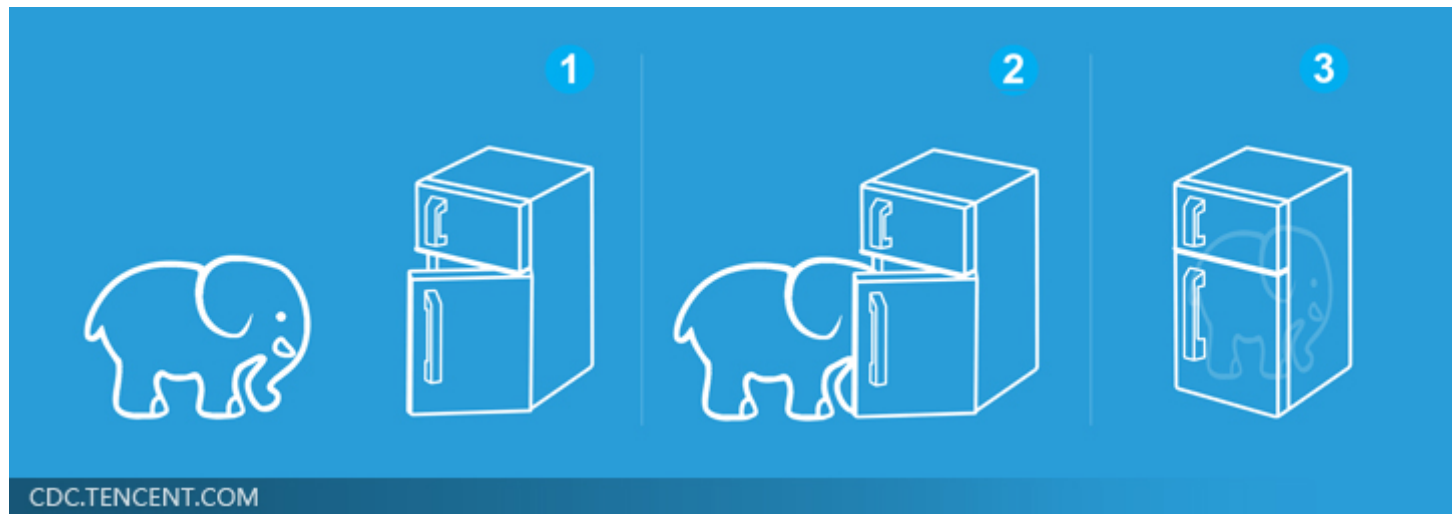
Step 3.2: Performance report

```
test_loss, correct = 0, 0
for data, target in test_loader:
    data, target = Variable(data, volatile=True), Variable(target)
    data = data.view(-1, 28 * 28)
    net_out = net(data)
    # sum up batch loss
    test_loss += criterion(net_out, target).data[0]
    pred = net_out.data.max(1)[1]  # get the index of the max log-probability
    correct += pred.eq(target.data).sum()

test_loss /= len(test_loader.dataset)
print('\nTest set: Average loss: {:.4f}, Accuracy: {}/{} ({:.0f}%)\n'.format(
    test_loss, correct, len(test_loader.dataset),
    100. * correct / len(test_loader.dataset)))
```

# Three Steps for Deep Learning

| Step 1: define a set of function | → | Step 2: goodness of function | → | Step 3: pick the best function |

Deep Learning is so simple ……



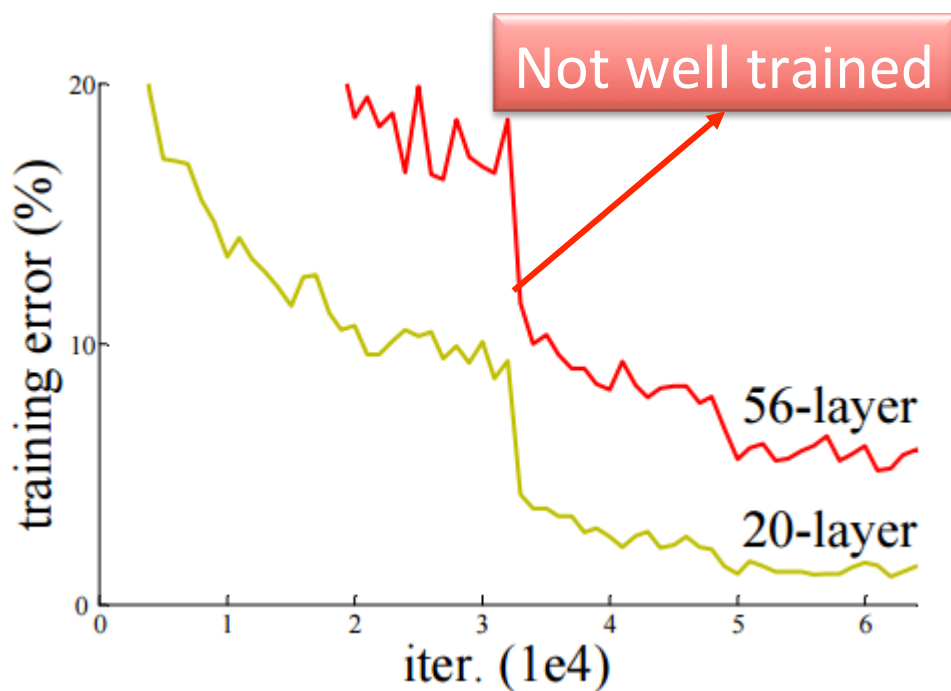CDC.TENCENT.COM

# Outline

Introduction of Deep Learning

"Hello World" for Deep Learning

Tips for Deep Learning

# Recipe of Deep Learning
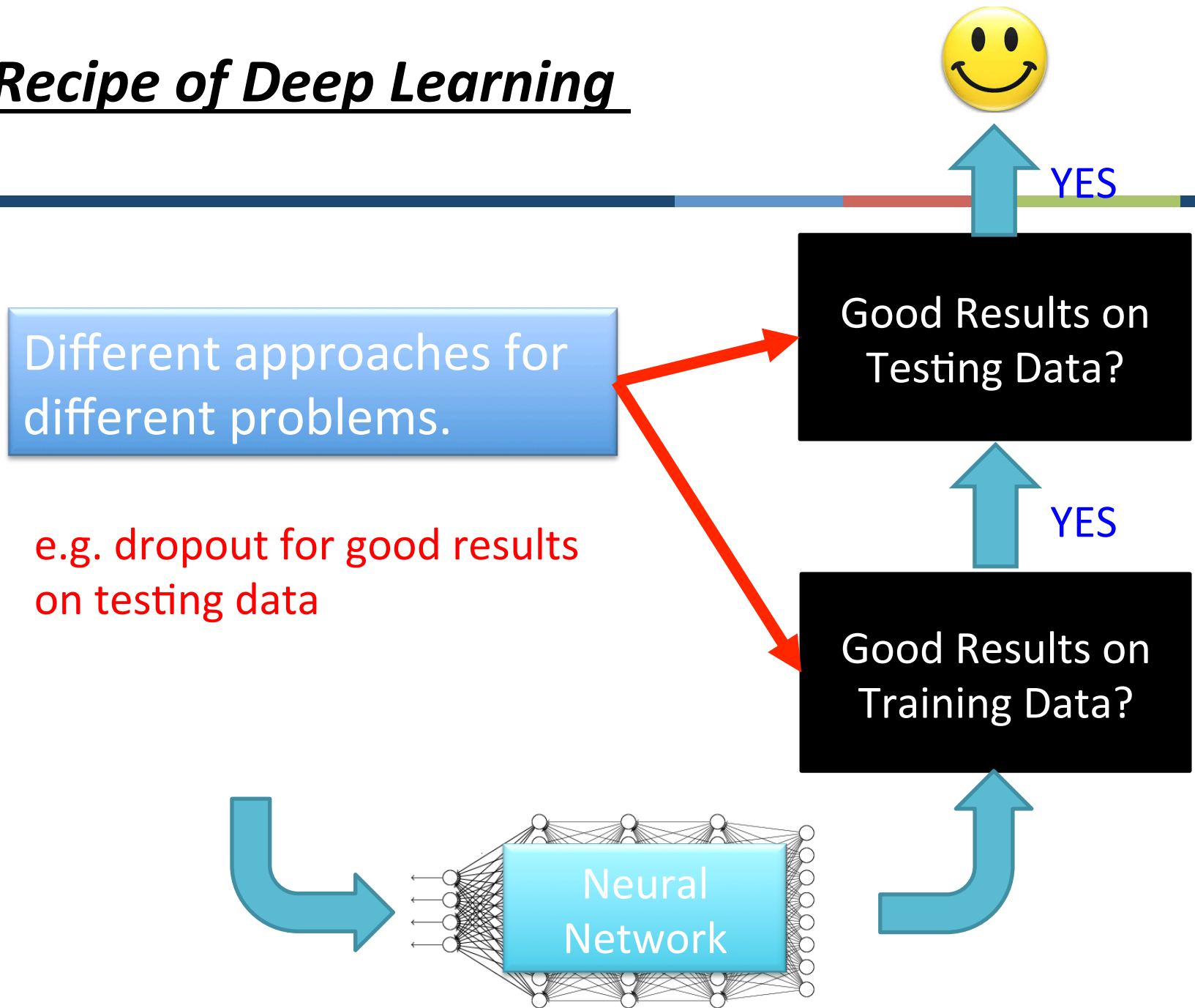
# Do not always blame Overfitting



Not well trained

Overfitting?
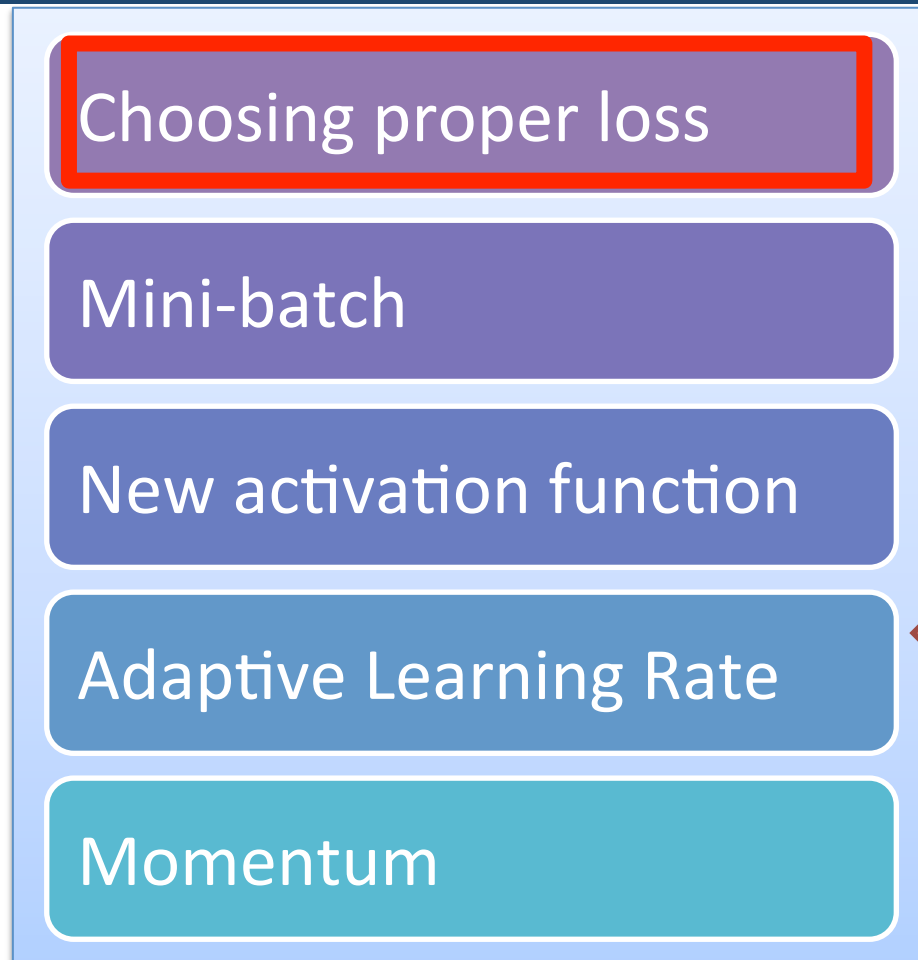
Training Data

Testing Data

Deep Residual Learning for Image Recognition
http://arxiv.org/abs/1512.03385

# *Recipe of Deep Learning*

Different approaches for different problems.

e.g. dropout for good results on testing data

Good Results on Testing Data?

YES

Good Results on Training Data?

YES

Neural Network

# *Recipe of Deep Learning*

Choosing proper loss

Mini-batch

New activation function

Adaptive Learning Rate

Momentum

Good Results on Testing Data?
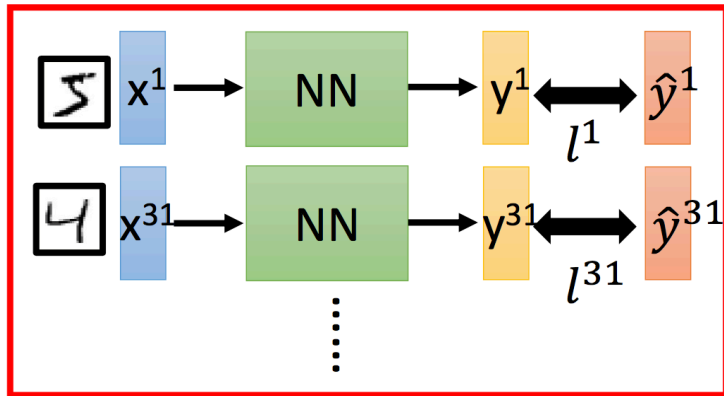
Good Results on Training Data?

YES

YES

# Choosing Proper Loss



"1"

$x_1$  $x_2$  $x$

Softmax

$y_1$  1
$y_2$  0
$y_{10}$  0

**Which one is better?**

loss

$y_{\downarrow 1}$  1
$y_{\downarrow 2}$  0
$y_{\downarrow 10}$  0

target

Square Error  $\sum_{i=1}^{10}(y_i - \hat{y}_i)^2$  **=0**

Cross Entropy  $-\sum_{i=1}^{10}\hat{y}_i \ln y_i$  **=0**

# Demo

## Square Error

```
model.compile(loss='mse',
              optimizer=SGD(lr=0.1),
              metrics=['accuracy'])
```

## Cross Entropy

```
model.compile(loss='categorical_crossentropy',
              optimizer=SGD(lr=0.1),
              metrics=['accuracy'])
```

Several alternatives: https://keras.io/objectives/

# Choosing Proper Loss

When using softmax output layer, choose cross entropy



Total Loss

Cross Entropy

Square Error

http://jmlr.org/proceedings/papers/v9/glorot10a/glorot10a.pdf

# *Recipe of Deep Learning*



Choosing proper loss

Mini-batch

New activation function

Adaptive Learning Rate

Momentum

Good Results on Testing Data?

Good Results on Training Data?

YES

YES

```
model.fit(x_train, y_train, batch_size=100, nb_epoch=20)
```
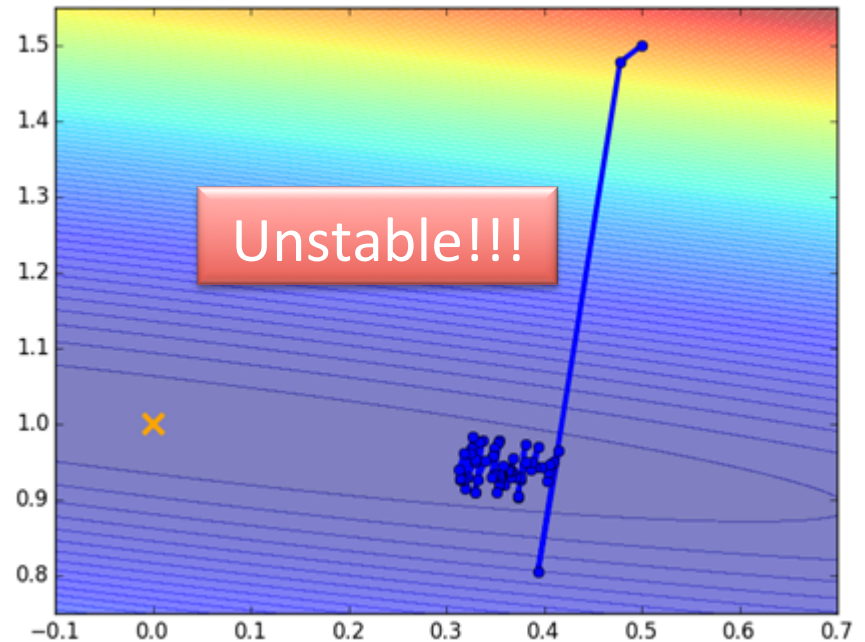
# Mini-batch

➢ Randomly initialize network parameters



➢ Pick the 1st batch

$$L' = l^1 + l^{31} + \cdots$$

Update parameters once

➢ Pick the 2nd batch

$$L'' = l^2 + l^{16} + \cdots$$

Update parameters once

⋮

➢ Until all mini-batches have been picked

one epoch

Repeat the above process

# Mini-batch

```
model.fit(x_train, y_train, batch_size=100, nb_epoch=20)
```



Mini-batch

100 examples in a mini-batch

Repeat 20 times

➢ Pick the 1st batch

$$L' = l^1 + l^{31} + \cdots$$

Update parameters once

➢ Pick the 2nd batch

$$L'' = l^2 + l^{16} + \cdots$$

Update parameters once

⋮

➢ Until all mini-batches have been picked

# Mini-batch

**_Original Gradient Descent_**

**_With Mini-batch_**



Unstable!!!

The colors represent the total loss.

# Mini-batch is Faster

Not always true with parallel computing.

## *Original Gradient Descent*
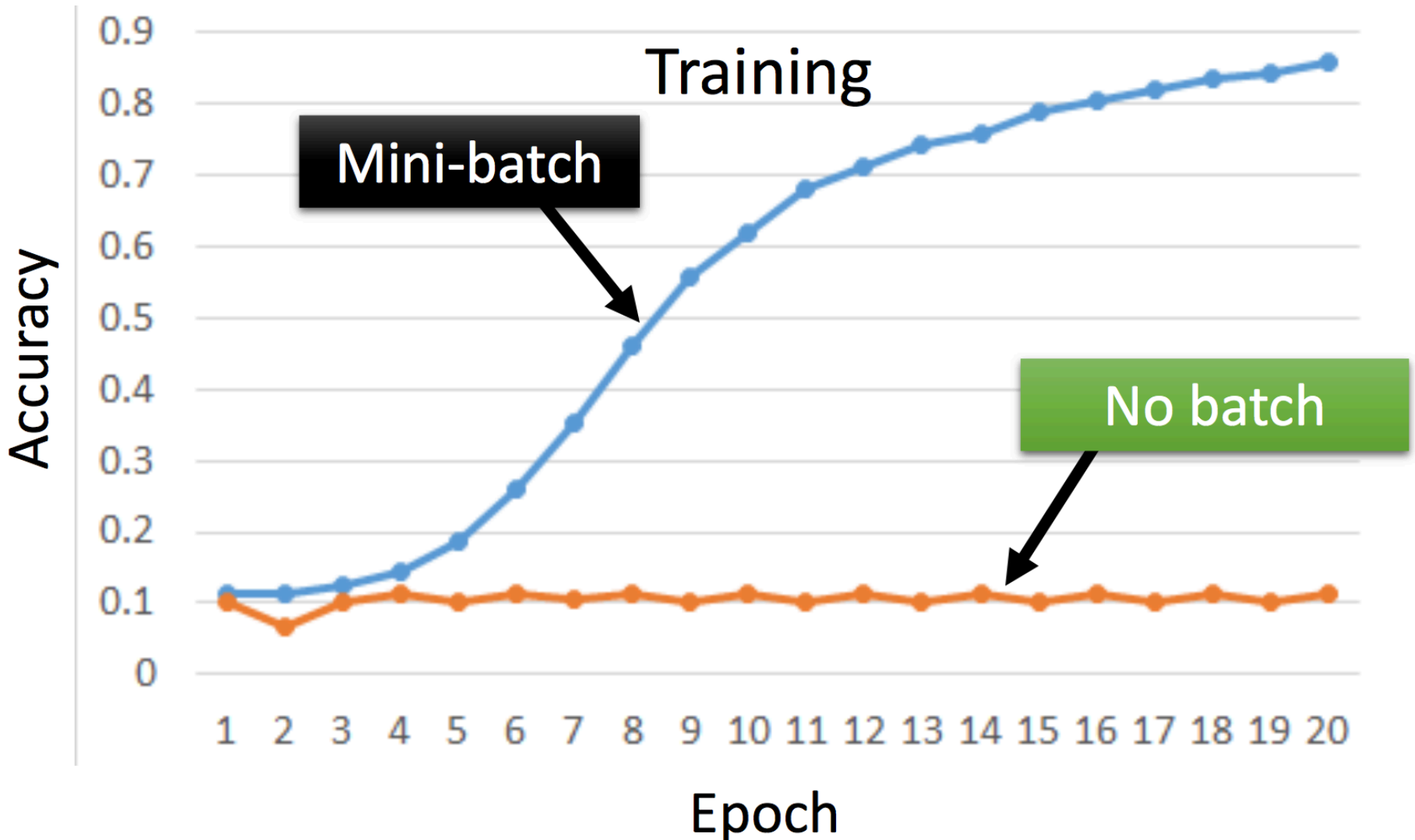
Update after seeing all examples

## *With Mini-batch*

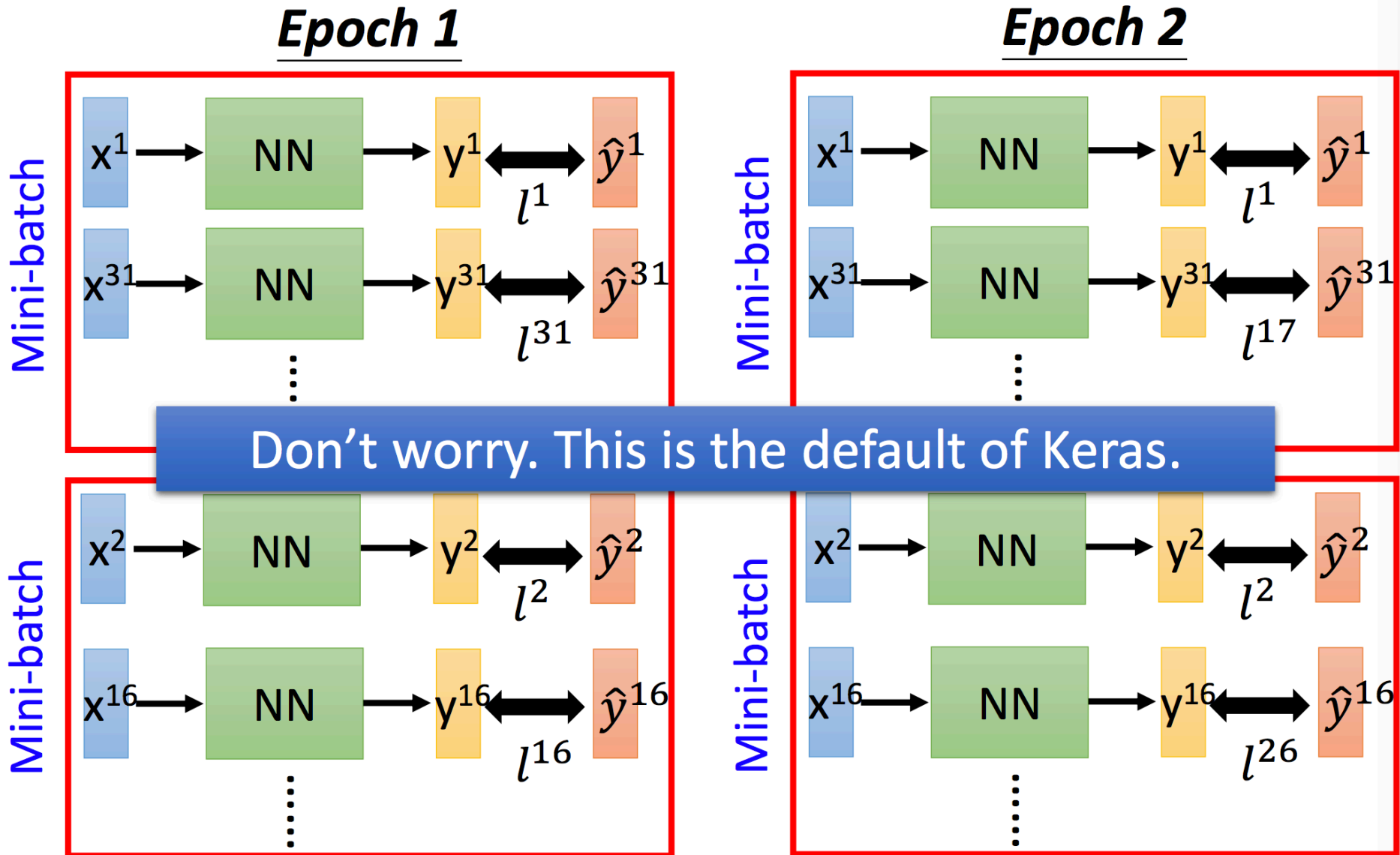If there are 20 batches, update 20 times in one epoch.



See all examples

Can have the same speed (not super large data set)

See only one batch

1 epoch

Mini-batch has better performance!

# Mini-batch is Faster

# *Shuffle the training examples for each epoch*

## *Epoch 1*



## *Epoch 2*



Don't worry. This is the default of Keras.

# Hard to get the power of Deep …



**Handwriting Digit Classification**

Results on Training Data

Deeper usually does not imply better.

# Vanishing Gradient Problem



Smaller gradients

Learn very slow

Almost random

Larger gradients

Learn very fast

Already converge

# Vanishing Gradient Problem



Smaller gradients
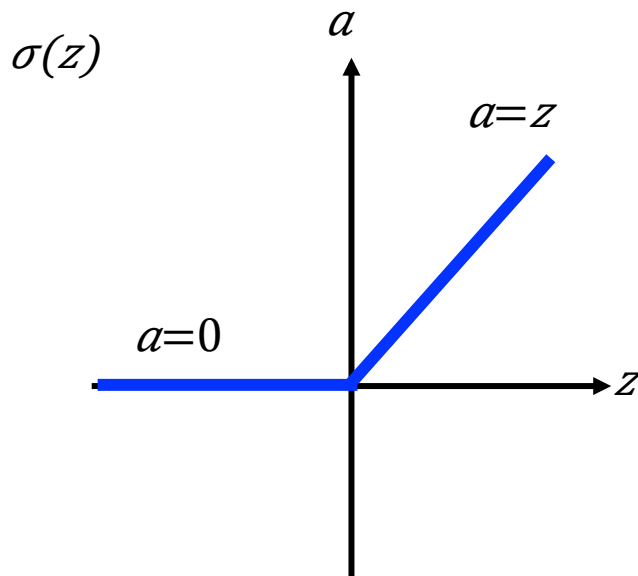
$x_1$ $x_2$ ⋮ $x_N$

$+\Delta w$

Small output

Large input

Intuitive way to compute the derivatives …

$$\partial l/\partial w = ?\ \Delta l/\Delta w$$

# ReLU

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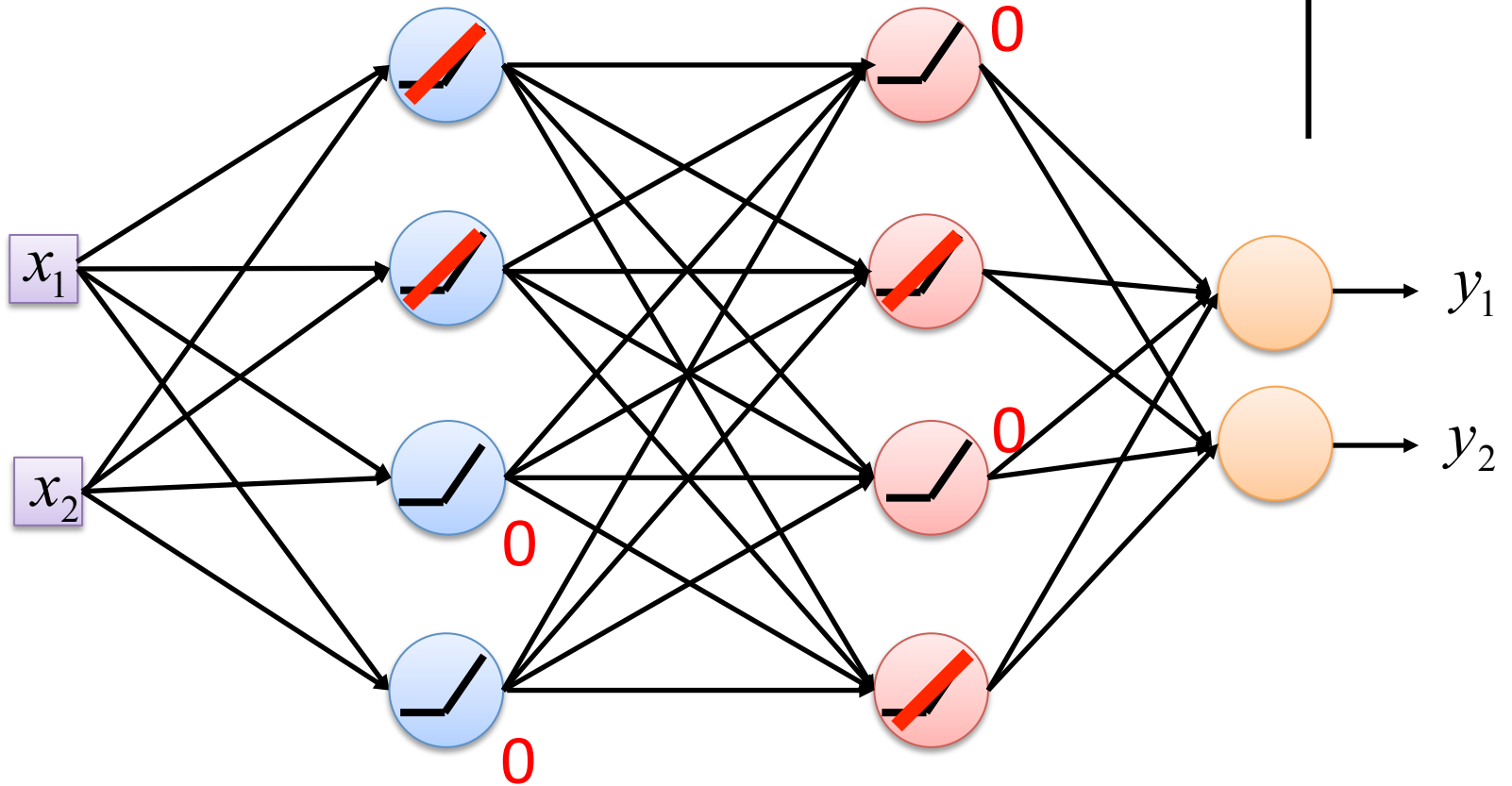
$\sigma(z)$

$a$

$a=z$

$a=0$

$z$

**_Reason:_**

1. Fast to compute

2. Sparsity

3. Vanishing gradient problem

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

# ReLU

# ReLU

# ReLU - variant

*Leaky ReLU*

$a$

$a=z$
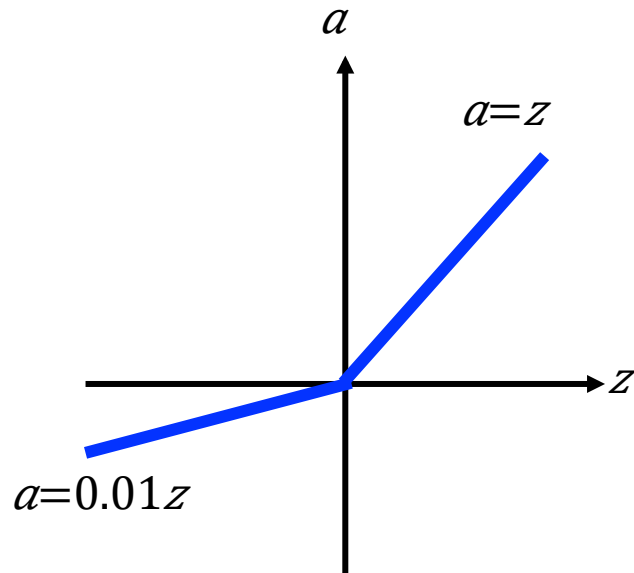
$z$

$a=0.01z$

*Parametric ReLU*

$a$

$a=z$

$z$

$a=\alpha z$

α also learned by gradient descent

# Maxout

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



You can have more than 2 elements in a group.

# 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

# *Recipe of Deep Learning*

Choosing proper loss

Mini-batch

New activation function

**Adaptive Learning Rate**

Momentum

Good Results on Testing Data?

YES

Good Results on Training Data?

YES

# Learning Rates

Set the learning rate η carefully



If learning rate is too large

Total loss may not decrease after each update

# Learning Rates

Set the learning rate η carefully



If learning rate is too large

Total loss may not decrease after each update

If learning rate is too small

Training would be too slow

# Learning Rates

- Popular & Simple Idea: Reduce the learning rate by some factor every few epochs.
  - At the beginning, we are far from the destination, so we use larger learning rate
  - After several epochs, we are close to the destination, so we reduce the learning rate
  - E.g. 1/t decay: $\eta^t = \eta/\sqrt{t+1}$
- Learning rate cannot be one-size-fits-all
  - Giving different parameters different learning rates

# Adagrad

Original: $w \leftarrow w - \eta \partial L / \partial w$

Adagrad: $w \leftarrow w - \boxed{\eta_w} \partial L / \partial w$

Parameter dependent learning rate

$$\eta_w = \frac{\eta}{\sqrt{\sum_{i=0}^{t} (g^i)^2}}$$

constant

$g^i$ is $\partial L / \partial w$ obtained at the i-th update

Summation of the square of the previous derivatives

# Adagrad

$$\eta_w = \frac{\eta}{\sqrt{\sum_{i=0}^{t}(g^i)^2}}$$

$w_1$

| $g^0$ |
|---|
| 0.1 |

Learning rate:

$$\frac{\eta}{\sqrt{0.1^2}} = \frac{\eta}{0.1}$$

$$\frac{\eta}{\sqrt{0.1^2 + 0.2^2}} = \frac{\eta}{0.22}$$

$w_2$

| $g^0$ |
|---|
| 20.0 |

Learning rate:

$$\frac{\eta}{\sqrt{20^2}} = \frac{\eta}{20}$$
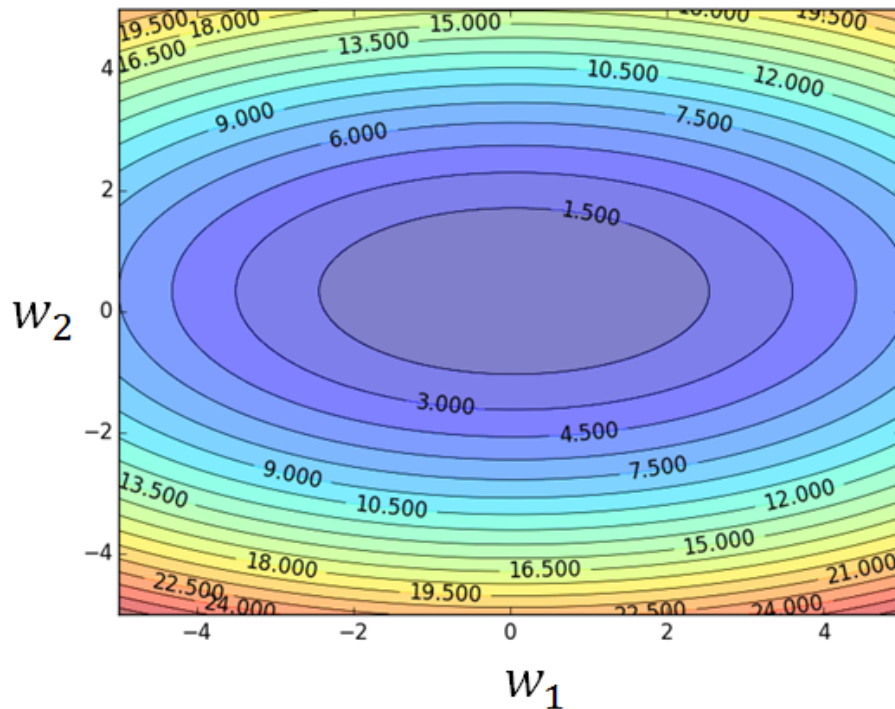
$$\frac{\eta}{\sqrt{20^2 + 10^2}} = \frac{\eta}{22}$$

**_Observation:_** 1. Learning rate is smaller and smaller for all parameters

2. Smaller derivatives, larger learning rate, and vice versa

Why?

Larger derivatives

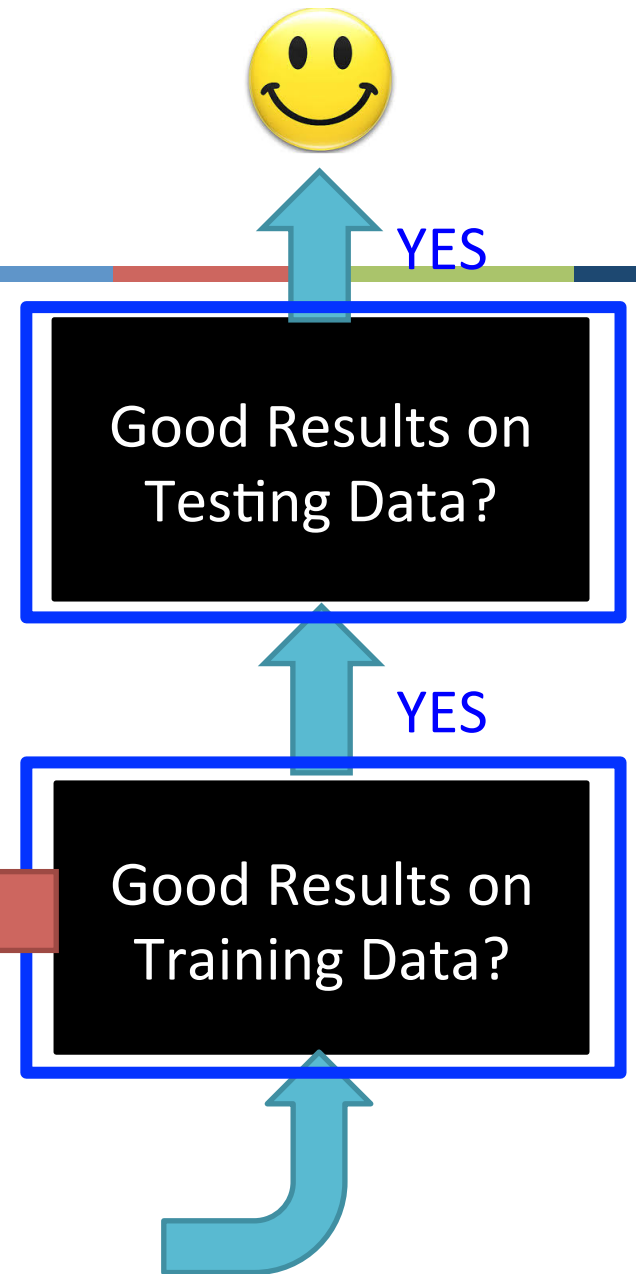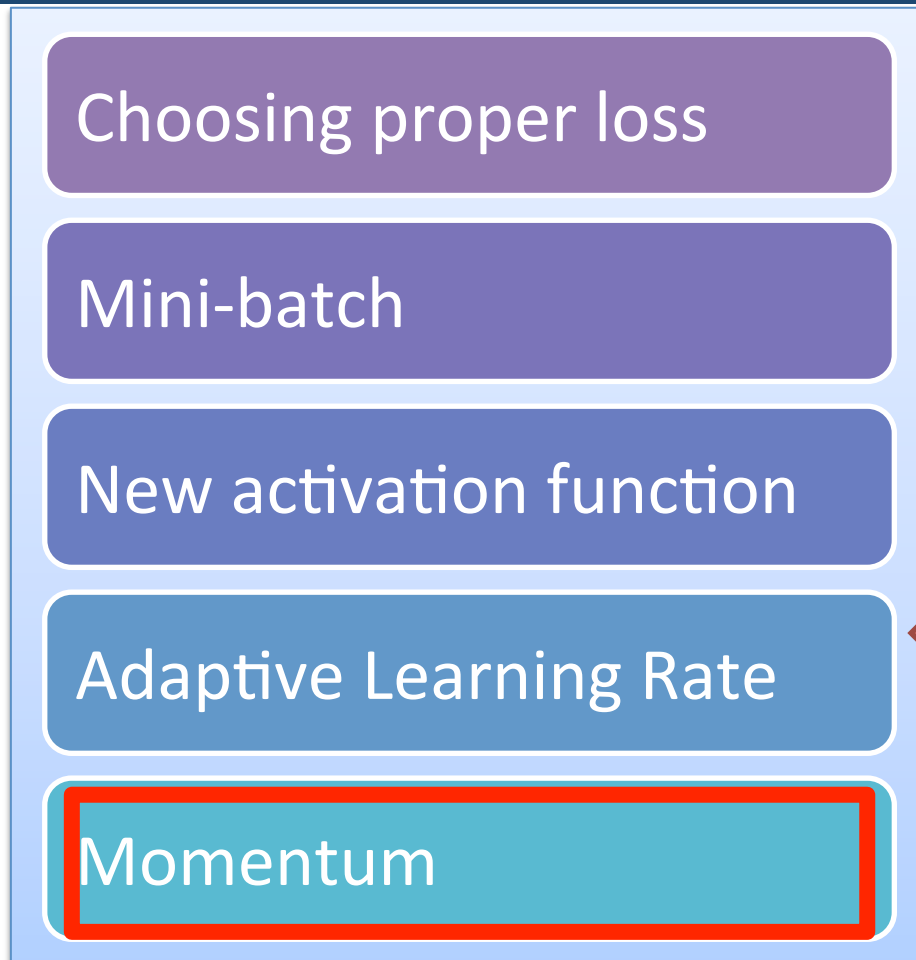Smaller Learning Rate

Smaller Derivatives

Larger Learning Rate

2. Smaller derivatives, larger learning rate, and vice versa

Why?

# Not the whole story ......

- Adagrad [John Duchi, JMLR'11]
- RMSprop
  - https://www.youtube.com/watch?v=O3sxAc4hxZU
- Adadelta [Matthew D. Zeiler, arXiv'12]
- "No more pesky learning rates" [Tom Schaul, arXiv'12]
- AdaSecant [Caglar Gulcehre, arXiv'14]
- Adam [Diederik P. Kingma, ICLR'15]
- Nadam
  - http://cs229.stanford.edu/proj2015/054_report.pdf

# *Recipe of Deep Learning*

Choosing proper loss

Mini-batch

New activation function

Adaptive Learning Rate

Momentum

Good Results on Testing Data?

YES

Good Results on Training Data?

YES

# Hard to find optimal network parameters



Total Loss

Very slow at the **plateau**

Stuck at saddle point

Stuck at local minima

$\partial L / \partial w \approx 0$

$\partial L / \partial w = 0$

$\partial L / \partial w = 0$

The value of a network parameter w

# In physical world ……

- Momentum

How about put this phenomenon in gradient descent?

# Momentum

Still not guarantee reaching global minima, but give some hope ……

cost

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

➡ Negative of $\partial L/\partial w$

⇢ Momentum

➡ Real Movement

$\partial L/\partial w = 0$

```python
model.compile(loss='categorical_crossentropy',
              optimizer=SGD(lr=0.1),
              metrics=['accuracy'])
```

```python
model.compile(loss='categorical_crossentropy',
              optimizer=Adam(),
              metrics=['accuracy'])
```

**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 1st moment vector)
  $v_0 \leftarrow 0$ (Initialize 2nd moment vector)
  $t \leftarrow 0$ (Initialize timestep)

```python
loss = nn.CrossEntropyLoss()
optimizer = torch.optim.Adam(model.parameters(), lr=0.1)
```

      at timestep $t$)
      ent estimate)
      moment estimate)
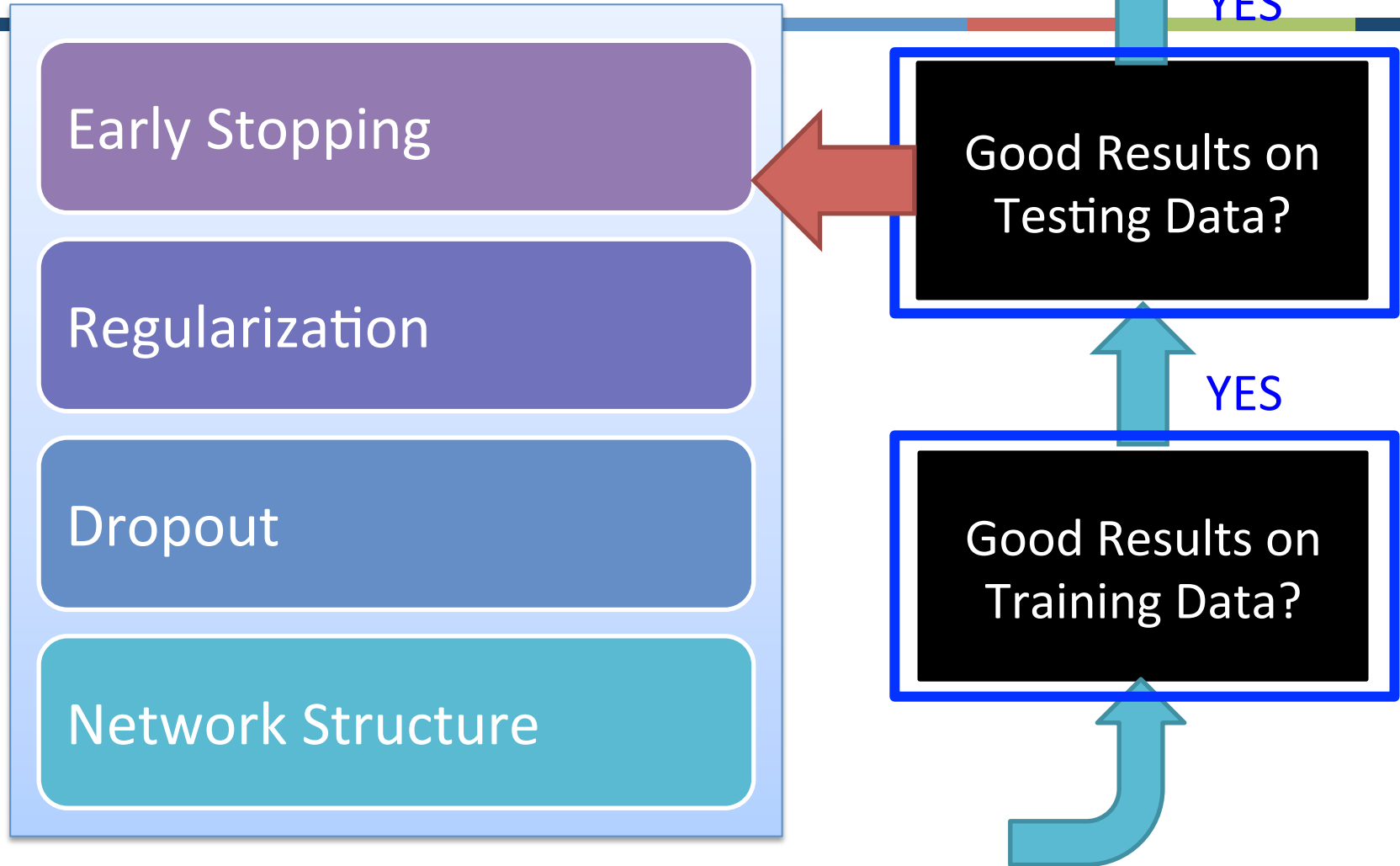  $\widehat{m}_t \leftarrow m_t/(1 - \beta_1^t)$ (Compute bias-corrected first moment estimate)
  $\widehat{v}_t \leftarrow v_t/(1 - \beta_2^t)$ (Compute bias-corrected second raw moment estimate)
  $\theta_t \leftarrow \theta_{t-1} - \alpha \cdot \widehat{m}_t/(\sqrt{\widehat{v}_t} + \epsilon)$ (Update parameters)
**end while**
**return** $\theta_t$ (Resulting parameters)

# Recipe of Deep Learning

- Early Stopping
- Regularization
- Dropout
- Network Structure

Good Results on Testing Data?

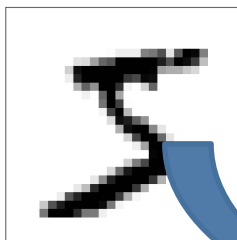Good Results on Training Data?

YES
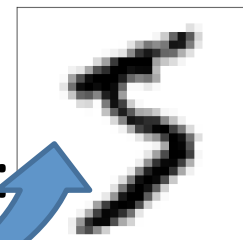
YES

# Panacea for Overfitting

- Have more training data

- **_Create_** more training data (?)
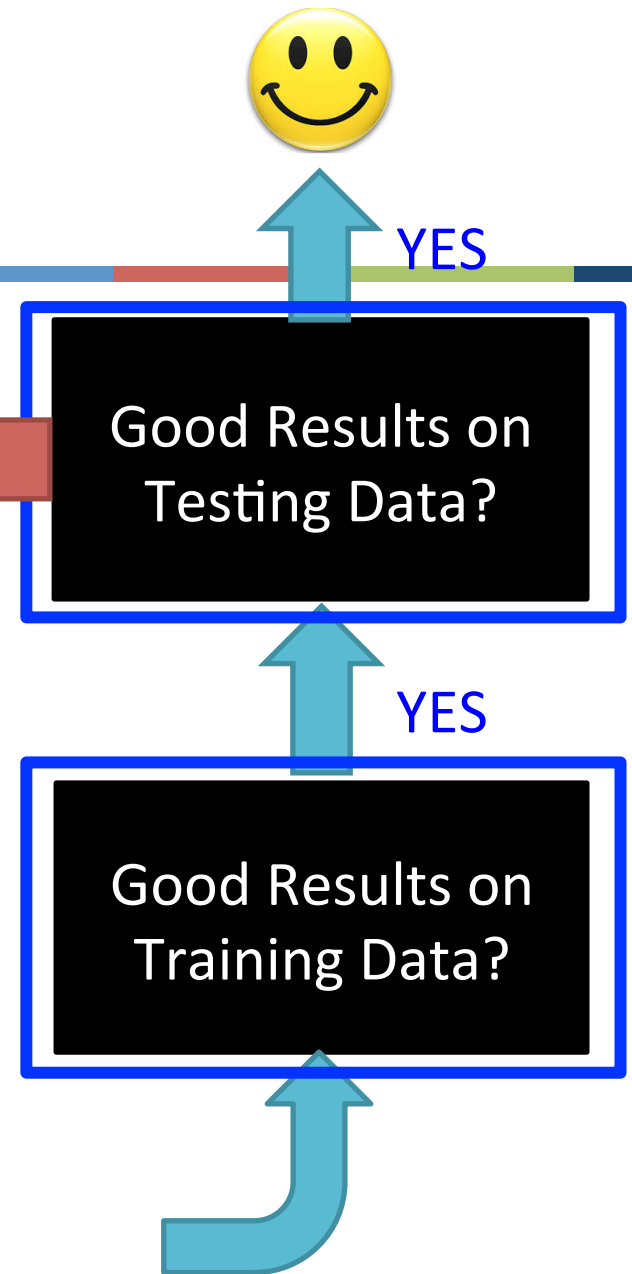
Handwriting recognition:

Original
Training Data:

Created
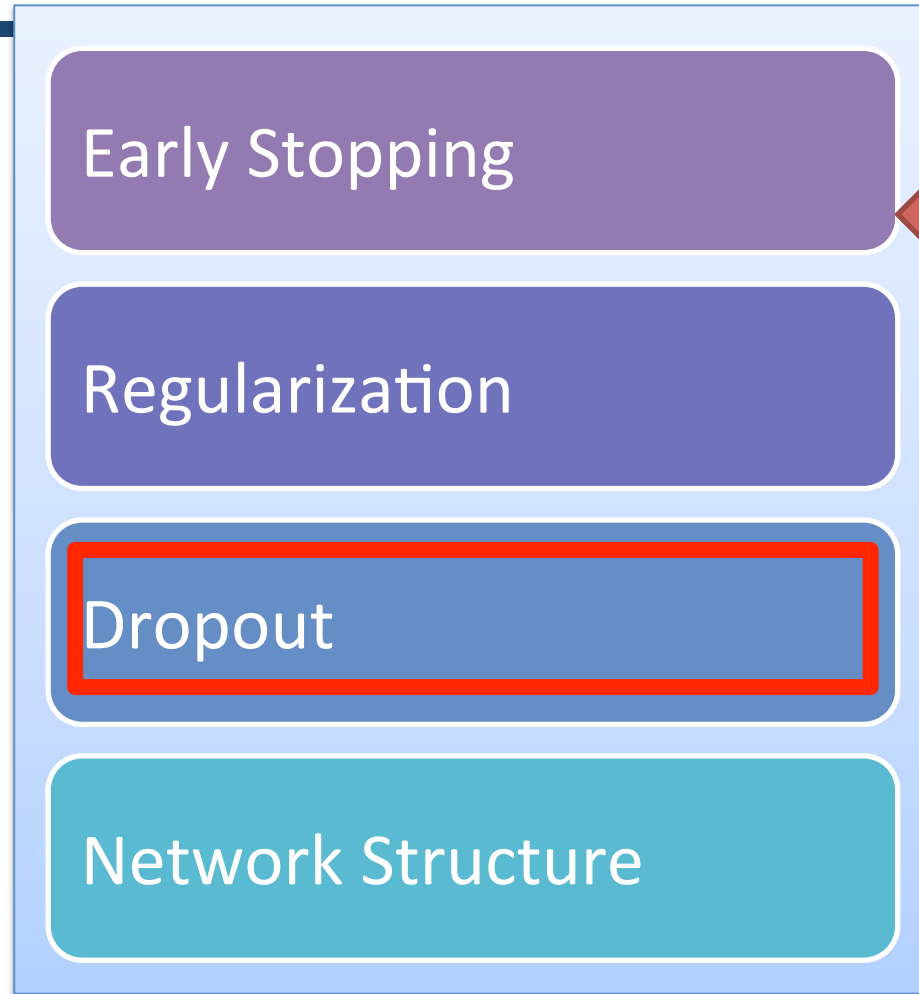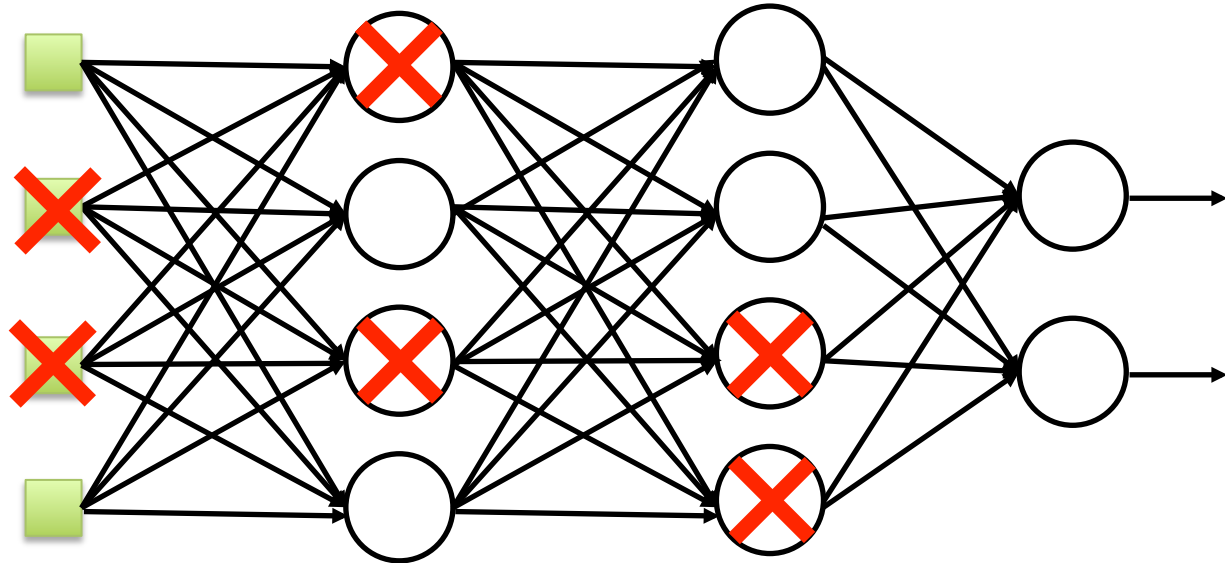Training Data:

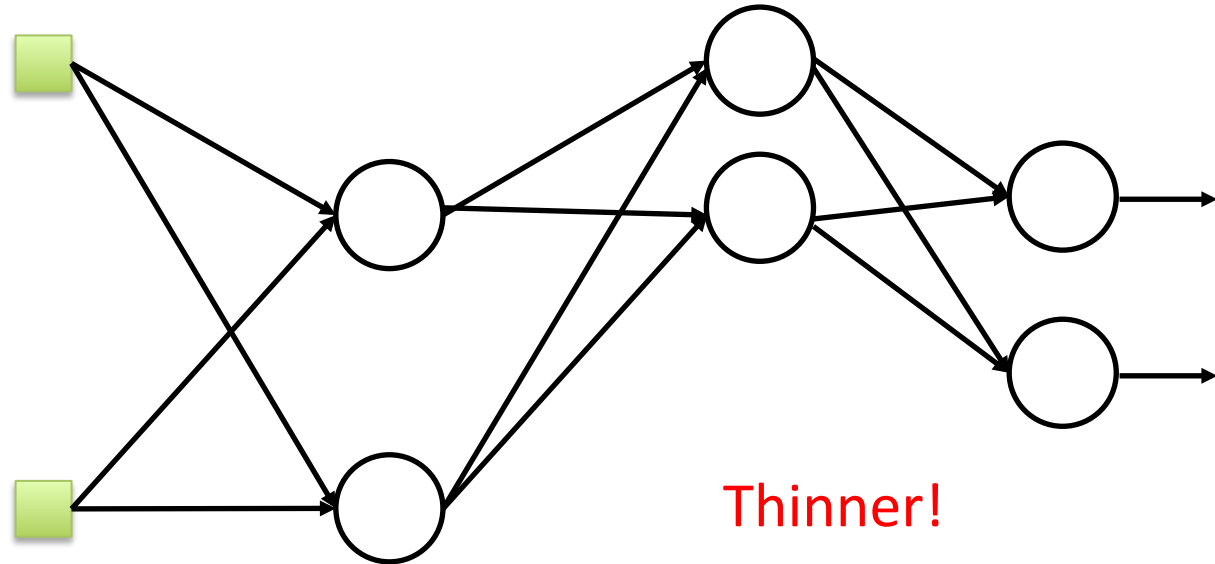Shift 15°

# *Recipe of Deep Learning*

# Dropout



**Training:**

> **Each time before updating the parameters**
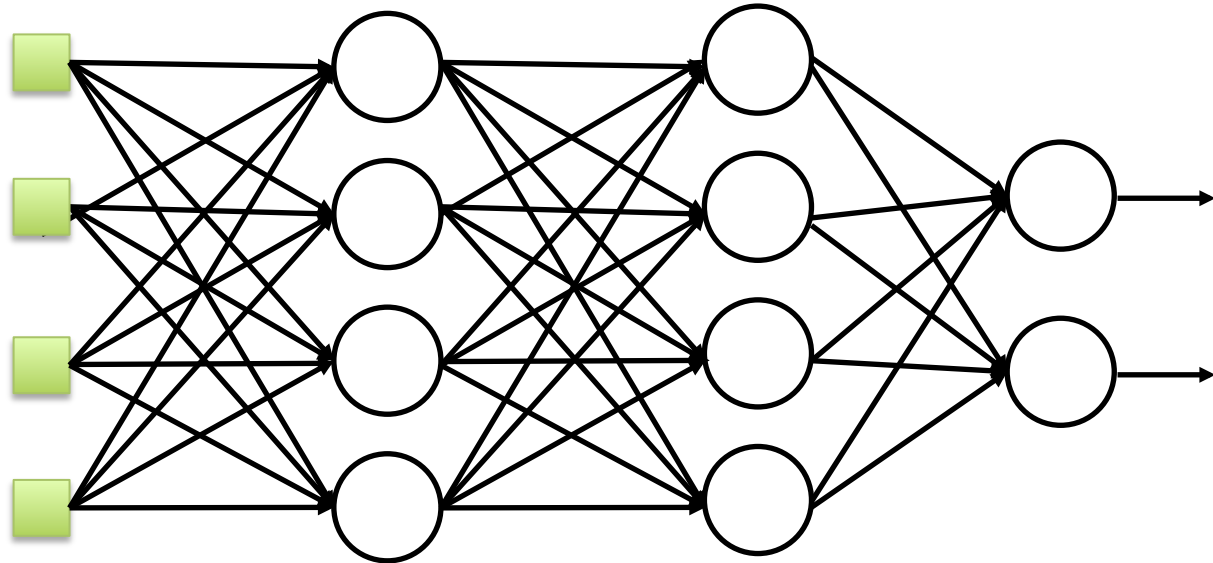>   - Each neuron has a probability of $p$ to dropout

# Dropout

Thinner!

➢ **Each time before updating the parameters**

- Each neuron has a probability of $p$ to dropout

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

- Using the new network for training

For each mini-batch, we resample the dropout neurons
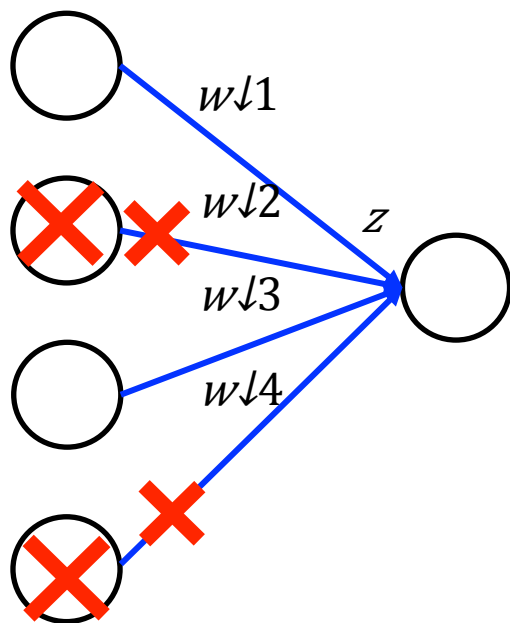
# Dropout

**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.

# Dropout - Intuitive Reason

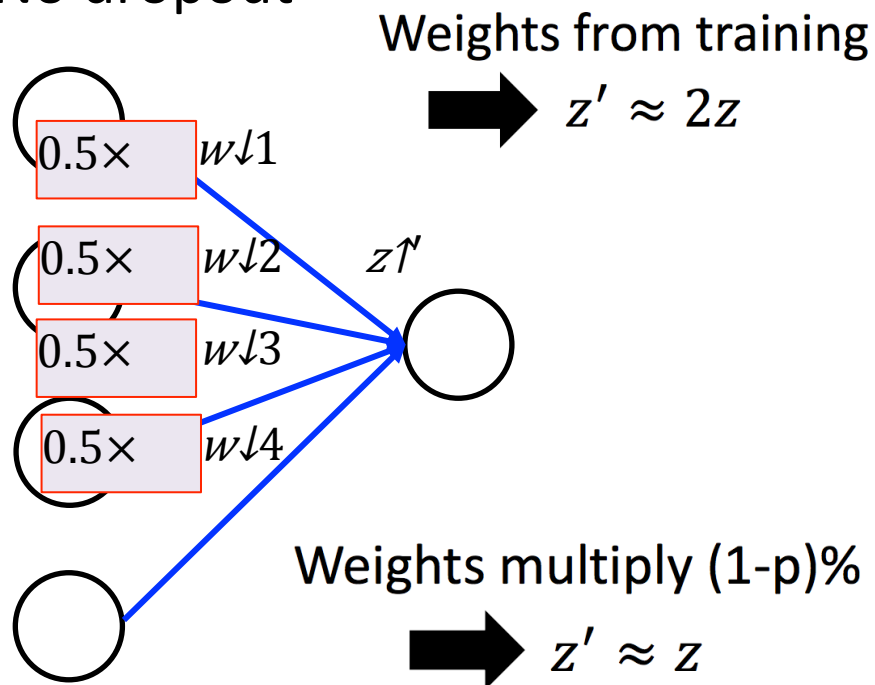- Why the weights should multiply (1-p) when testing?

**_Training of Dropout_**

Assume dropout rate is 0.5
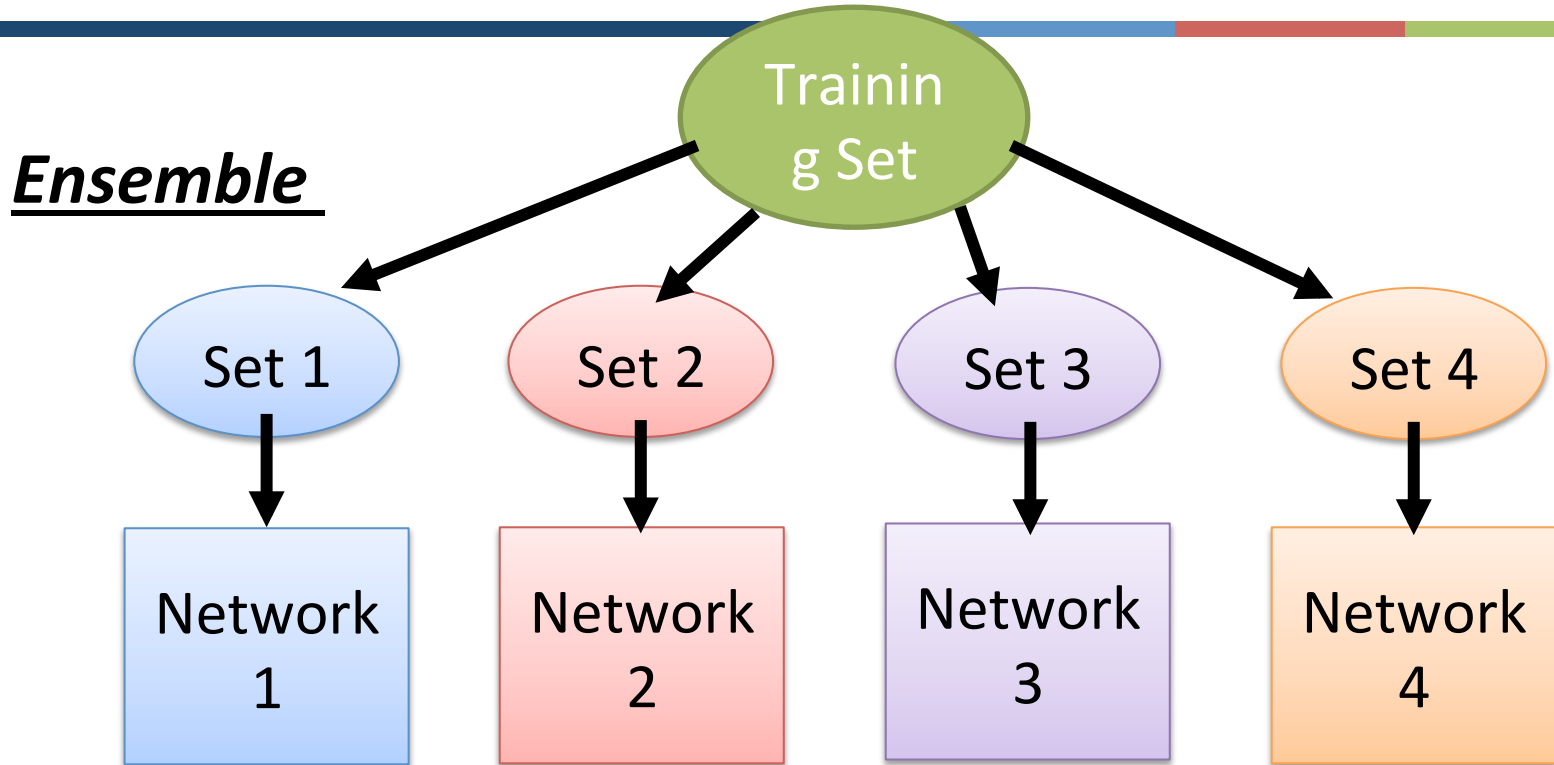


**_Testing of Dropout_**
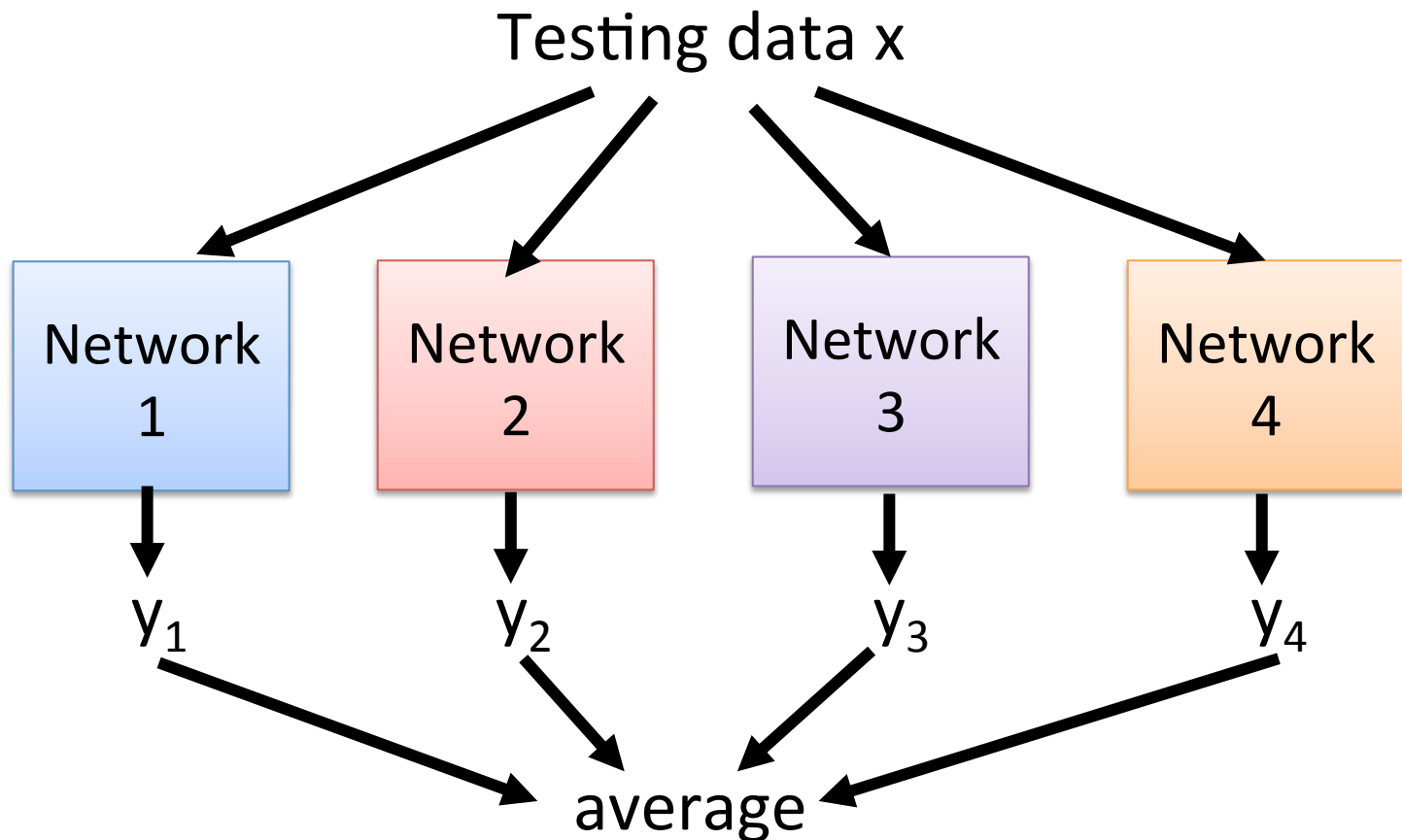
No dropout

Weights from training $\Rightarrow z' \approx 2z$

Weights multiply (1-p)% $\Rightarrow z' \approx z$

# Dropout is a kind of ensemble.



**_Ensemble_**

Train a bunch of networks with different structures

# Dropout is a kind of ensemble.

**_Ensemble_**

Testing data x

Network 1 → $y_1$

Network 2 → $y_2$

Network 3 → $y_3$

Network 4 → $y_4$

$y_1$, $y_2$, $y_3$, $y_4$ → average

# Dropout is a kind of ensemble.



**Training of Dropout**

M neurons

$2^M$ possible networks
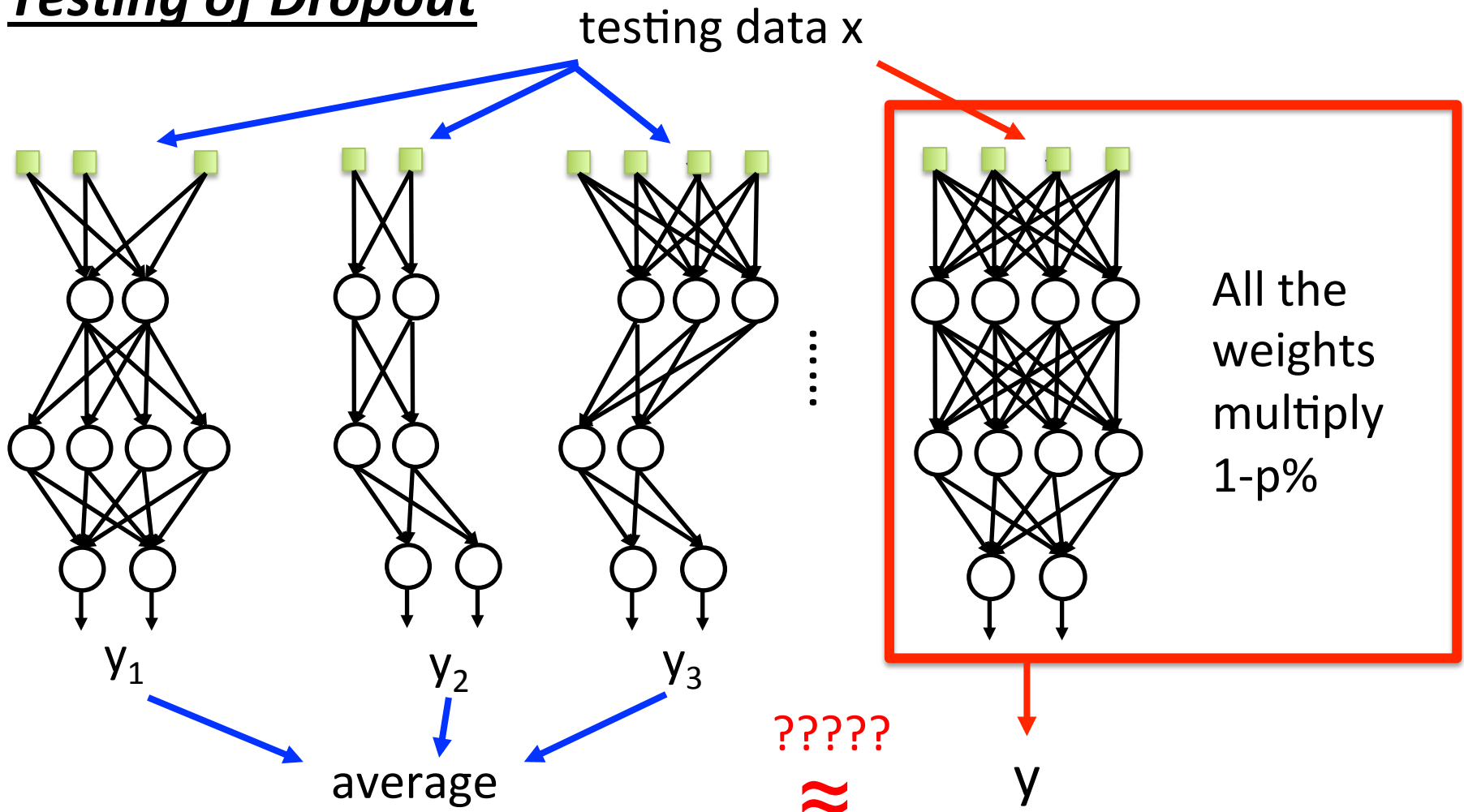
➢Using one mini-batch to train one network

➢Some parameters in the network are shared

# Dropout is a kind of ensemble.



**_Testing of Dropout_**

testing data x

$y_1$

$y_2$

$y_3$

All the weights multiply 1-p%

average

?????

$\approx$
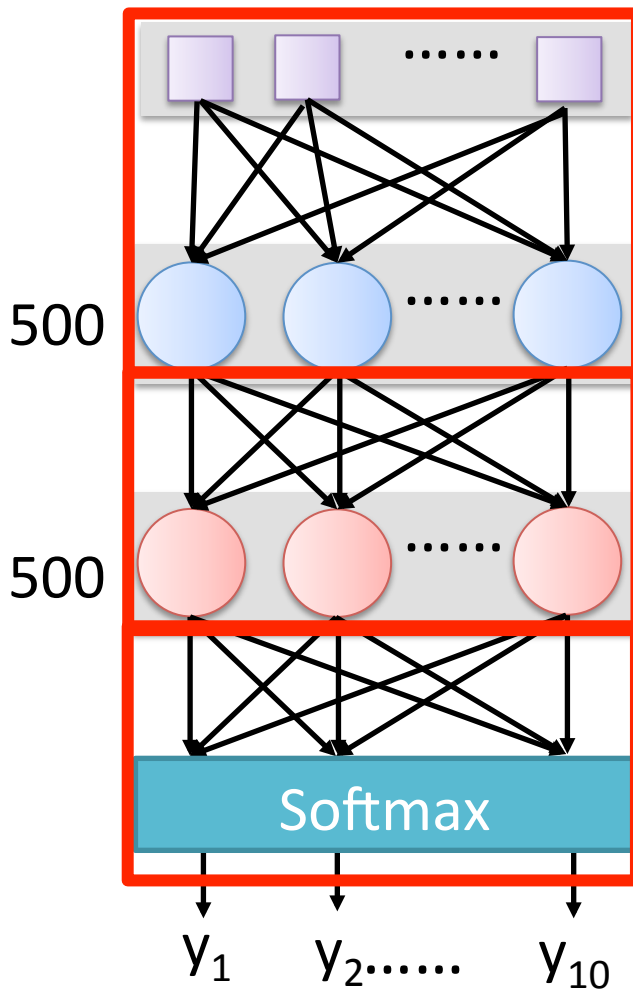
y

# More about dropout

- More reference for dropout [Nitish Srivastava, JMLR'14] [Pierre Baldi, NIPS'13][Geoffrey E. Hinton, arXiv'12]

- Dropout works better with Maxout [Ian J. Goodfellow, ICML'13]

- Dropconnect [Li Wan, *ICML'13*]

  - Dropout delete neurons

  - Dropconnect deletes the connection between neurons

- Annealed dropout [S.J. Rennie, SLT'14]

  - Dropout rate decreases by epochs

- Standout [J. Ba, NISP'13]

  - Each neural has different dropout rate

# Demo



```
model = Sequential()
```

```
model.add( Dense( input_dim=28*28,
                  output_dim=500 ))
model.add( Activation('sigmoid') )
```

model.add( dropout(0.8) )

```
model.add( Dense( output_dim=500 ) )
model.add( Activation('sigmoid') )
```

model.add( dropout(0.7) )

```
model.add( Dense(output_dim=10 ) )
model.add( Activation('softmax') )
```
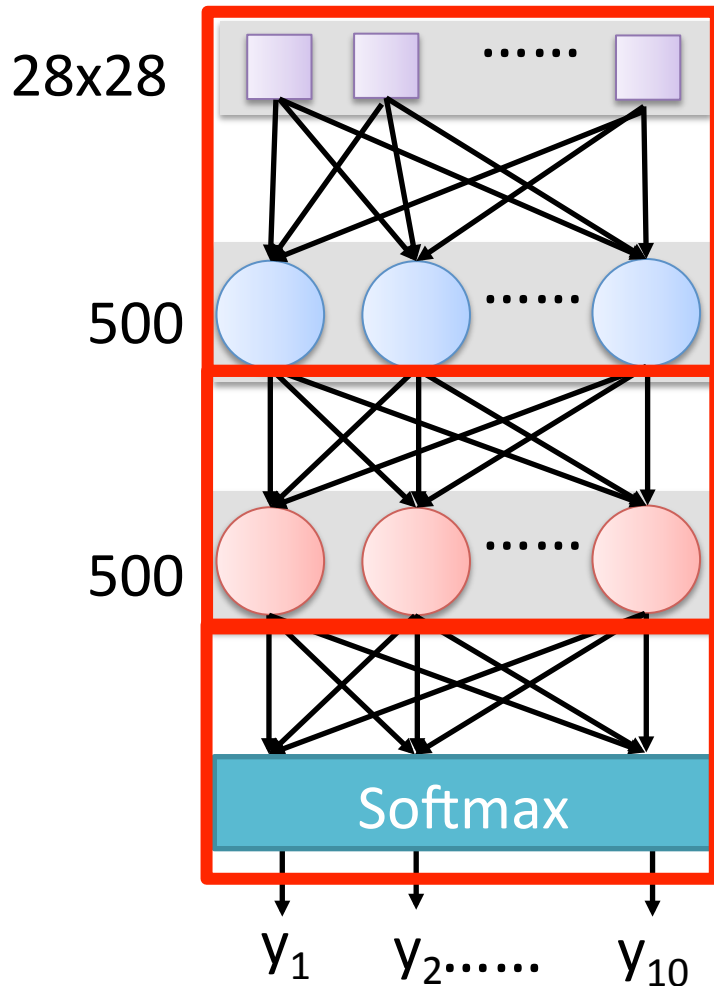
# PyTorch

28x28

500

500

Softmax
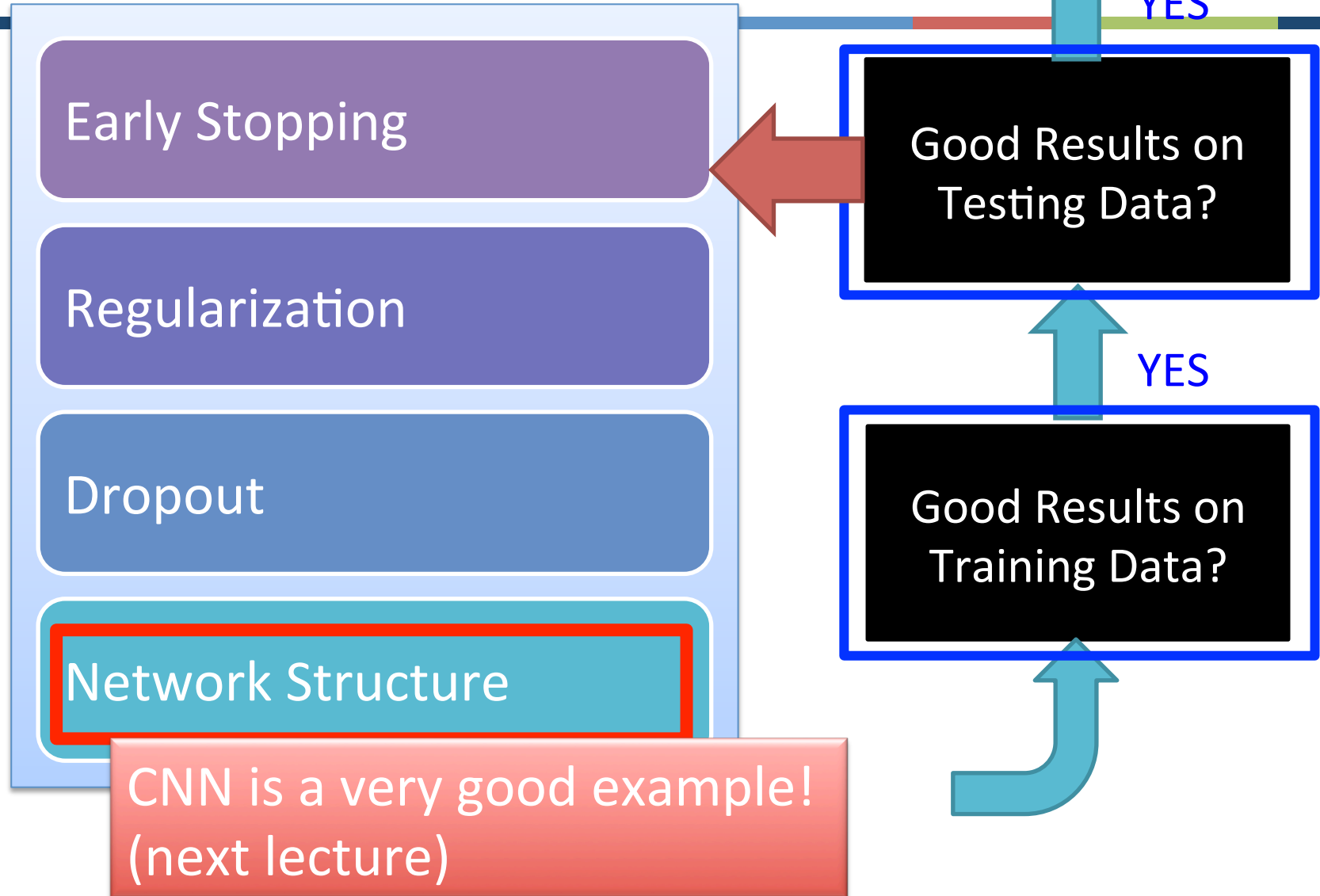
$y_1$  $y_2$......  $y_{10}$

```python
import torch.nn as nn
import torch.nn.functional as F

class MyNetwork(nn.Module):
    def __init__(self):
        super(MyNetwork, self).__init__()

        self.fc1 = nn.Linear(28 * 28, 500)
        self.fc2 = nn.Linear(500, 500)
        self.fc3 = nn.Linear(500, 10)
        self.do1 = nn.Dropout(0.8)
        self.do2 = nn.Dropout(0.7)

    def forward(self, x):
        x = F.sigmoid(self.fc1(x))
        x = self.do1(x)
        x = F.sigmoid(self.fc2(x))
        x = self.do2(x)
        return F.log_softmax(self.fc3(x))
```

# Concluding Remarks

# *Recipe of Deep Learning*

Step 1: define a set of function

Step 2: goodness of function

Step 3: pick the best function

Neural Network

Good Results on Testing Data?

Good Results on Training Data?

YES

NO

YES

NO