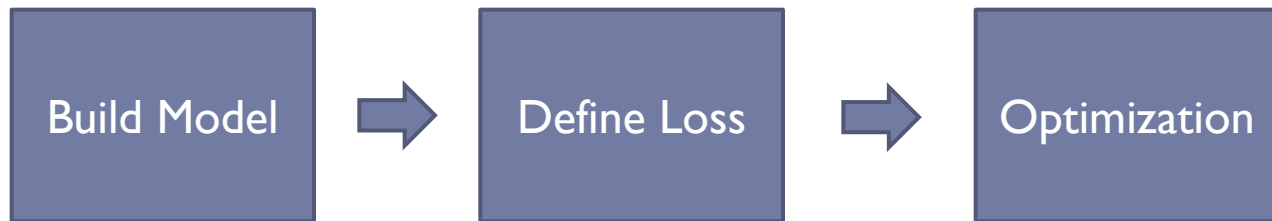# DNN Introduction

講者：Isaac

# Outline

▸ **Build DNN Model**

▸ **Define Loss Function**

   ▸ Mean square

   ▸ Cross entropy with softmax

▸ **Optimization**

   ▸ Gradient decent with moment

   ▸ Adagrad, RMS, and Adam

▸ **Predict/Validate Result**
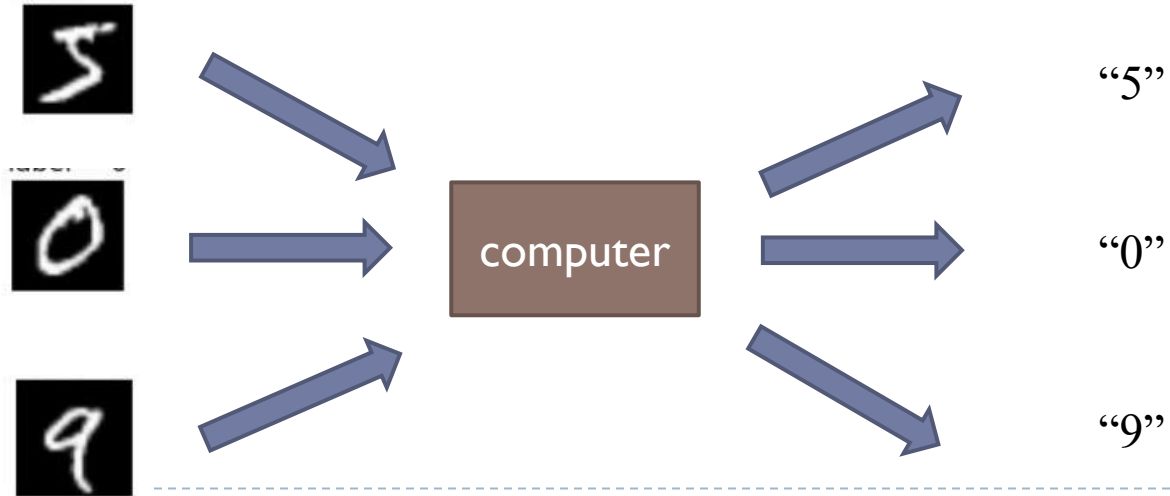
# Big Picture

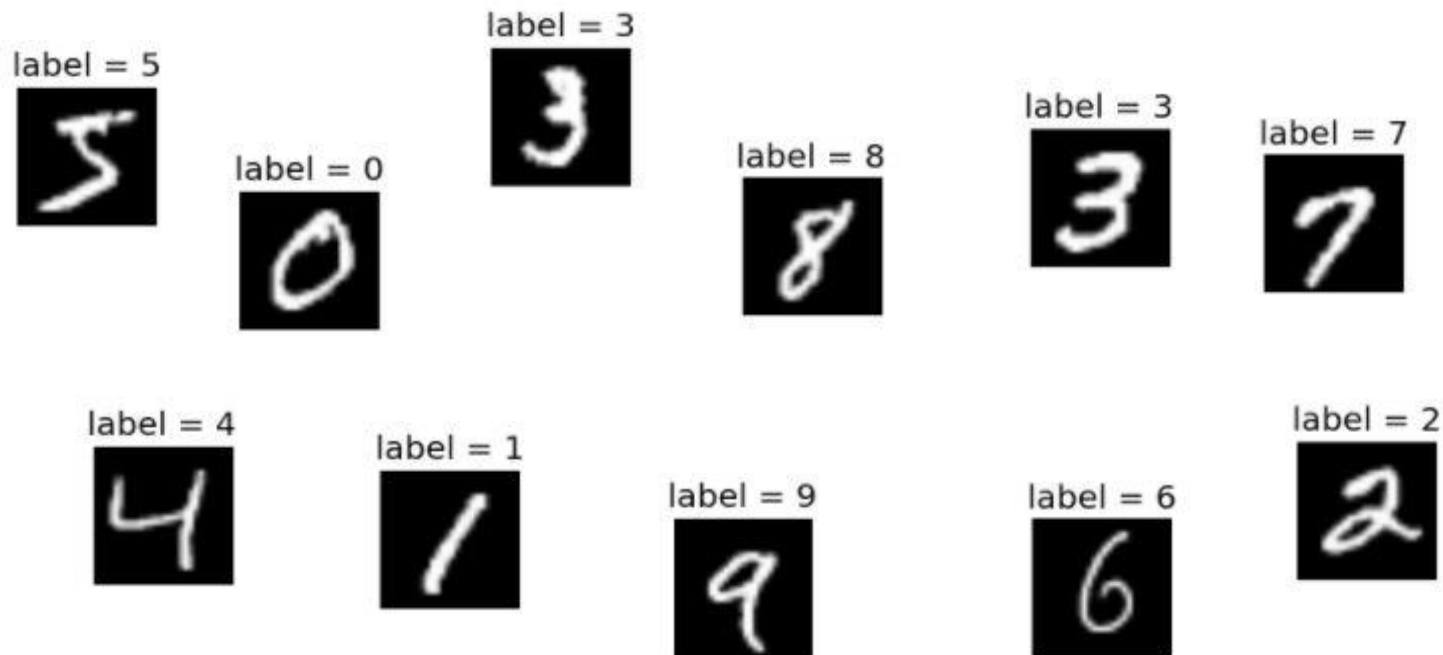Build Model → Define Loss → Optimization

Same logic when implementation in TensorFlow

Assume we want to build a handwritten system to recognize image from 0 to 9

Total 10000 Images with labels
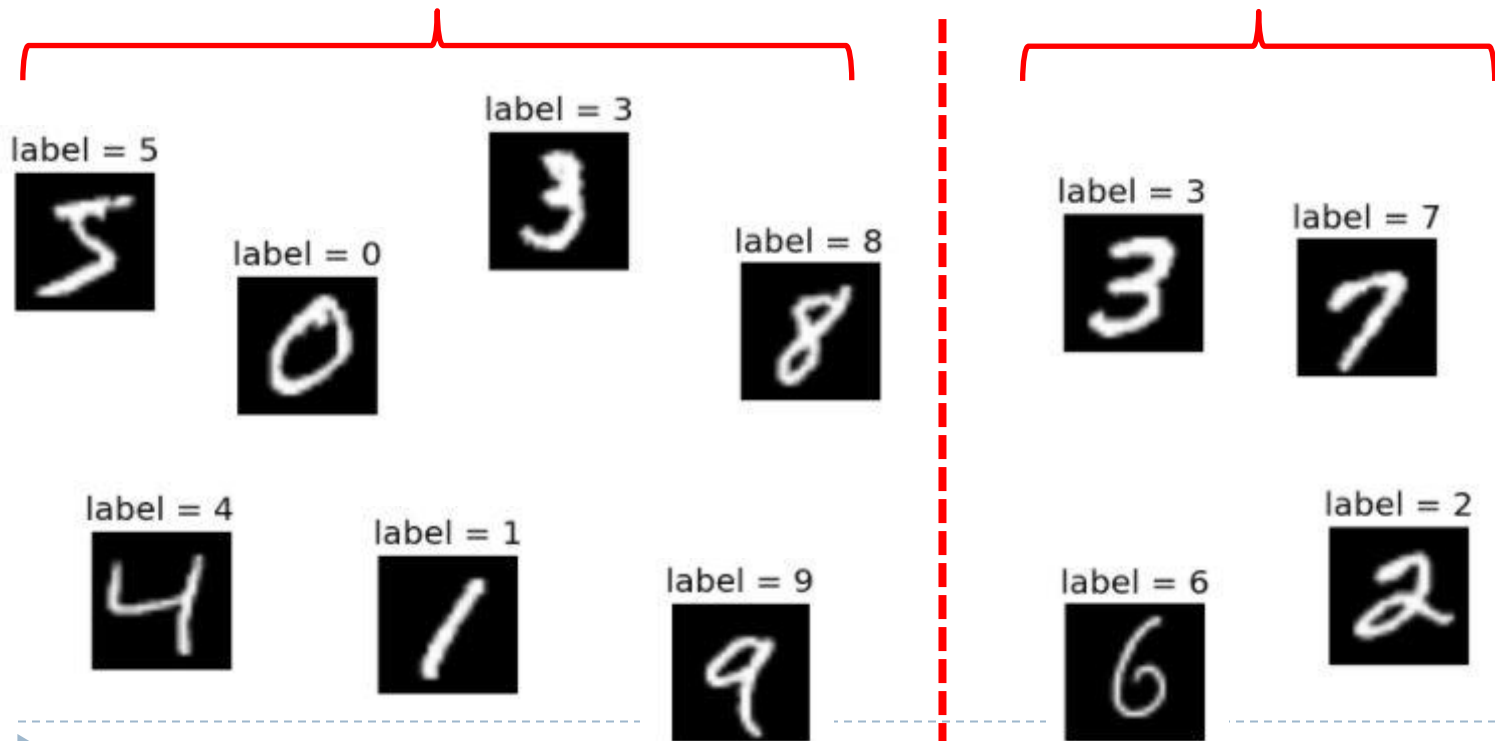
# Split the data



8000 images as training data

2000 images as testing data

label = 5

label = 3

label = 0

label = 8

label = 3

label = 7

label = 4

label = 1

label = 9

label = 6

label = 2

# Encode labels

8000 images as training data



One-hot encoding                    One-hot encoding

# What we want



$$f(x) \Rightarrow \begin{bmatrix} y_1 \\ y_2 \\ y_3 \\ y_4 \\ y_5 \\ y_6 \\ y_7 \\ y_8 \\ y_9 \\ y_{10} \end{bmatrix} \Longleftrightarrow \begin{bmatrix} 0 \\ 0 \\ 0 \\ 0 \\ 0 \\ 1 \\ 0 \\ 0 \\ 0 \\ 0 \end{bmatrix}$$

We want this two as close as possible

# What we want



2. Define Loss

$f(x)$

1. Build DNN Model

$$\begin{bmatrix} y_1 \\ y_2 \\ y_3 \\ y_4 \\ y_5 \\ y_6 \\ y_7 \\ y_8 \\ y_9 \\ y_{10} \end{bmatrix} \Longleftrightarrow \begin{bmatrix} 0 \\ 0 \\ 0 \\ 0 \\ 0 \\ 1 \\ 0 \\ 0 \\ 0 \\ 0 \end{bmatrix}$$

$$f(x) \Rightarrow \begin{bmatrix} y_1 \\ y_2 \\ y_3 \\ y_4 \\ y_5 \\ y_6 \\ y_7 \\ y_8 \\ y_9 \\ y_{10} \end{bmatrix} \Longleftrightarrow \begin{bmatrix} 0 \\ 0 \\ 0 \\ 0 \\ 0 \\ 1 \\ 0 \\ 0 \\ 0 \\ 0 \end{bmatrix}$$

# Finally, we get



$$f^*(x)$$

$$\begin{bmatrix} 0 \\ 0.2 \\ 0 \\ 0.1 \\ 0 \\ 0.8 \\ 0 \\ 0 \\ 0 \\ 0 \end{bmatrix}$$

We get it is '5'

# Build DNN Model

# Model Overview



$\sigma(x1 * w_{11} + x2 * w_{12} + x3 * w_{13} + b_1)$

$\sigma(x1 * w_{21} + x2 * w_{22} + x3 * w_{23} + b_2)$

$\sigma(x1 * w_{31} + x2 * w_{32} + x3 * w_{33} + b_3)$
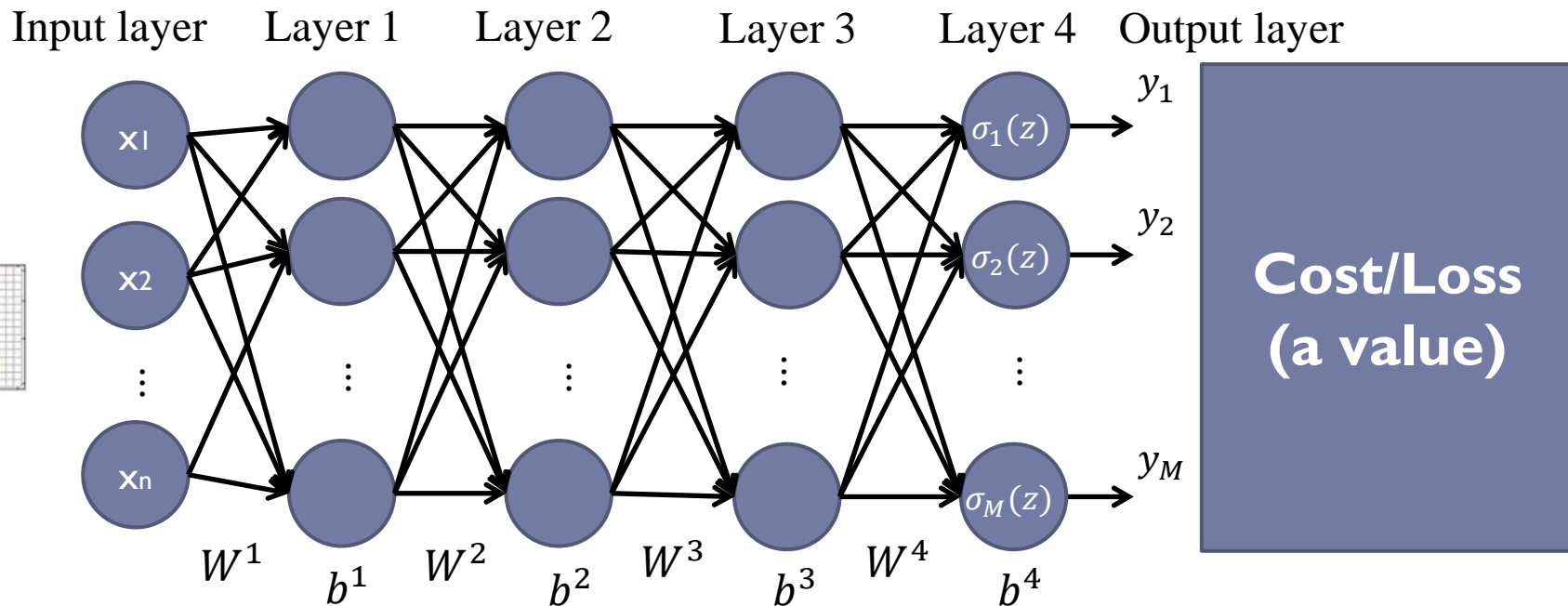
# Model Overview



$$\sigma(W * X + b)$$

$$X = \begin{bmatrix} x1 \\ x2 \\ x3 \end{bmatrix} \qquad W = \begin{bmatrix} w_{11} & w_{12} & w_{13} \\ w_{21} & w_{22} & w_{23} \\ w_{31} & w_{32} & w_{33} \end{bmatrix}$$

$$b = \begin{bmatrix} b1 \\ b2 \\ b3 \end{bmatrix}$$

▶ 14

# Model Overview

Input layer    Layer 1    Layer 2    Layer 3    Layer 4    Output layer

$x_1$

$x_2$

$x_n$

$\sigma_1(z)$    $y_1$

$\sigma_2(z)$    $y_2$

$\sigma_M(z)$    $y_M$

**Cost/Loss (a value)**

$W^1$    $b^1$    $W^2$    $b^2$    $W^3$    $b^3$    $W^4$    $b^4$

DNN also called "Multilayer perceptron"

# Model Overview



$$y = f(x) = \sigma(W^L \dots \sigma(W^2 \sigma(W^1 X + b^1) + b^2) \dots + b^L)$$

# Model Overview

Input layer     Layer 1     Layer 2             Layer L    Output layer
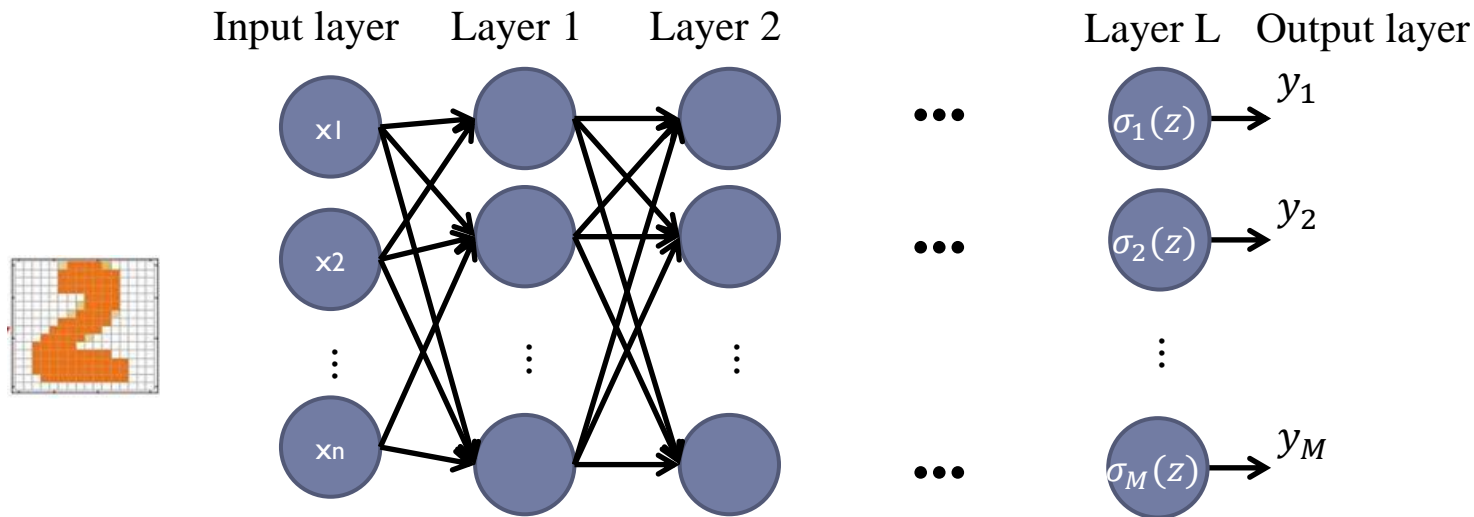
$$y = f(x) = \sigma(W^L \ldots \sigma(W^2 \sigma(W^1 X + b^1) + b^2) \ldots + b^L)$$

Weight (variable)        Bias (variable)

# Model Overview

Input layer   Layer 1   Layer 2          Layer L   Output layer



$$y = f(x) = \sigma(W^L \ldots \sigma(W^2 \sigma(W^1 X + b^1) + b^2) \ldots + b^L)$$

We want to find the best parameter set :

$$\theta = \{W^1, b^1, W^2, b^2, \ldots, W^L, b^L\}$$

# Parameters Initialization

▸ The most common way are

    ▸ Use random normal distribution with zero mean and small standard deviation

▸ In complex model, it is very important to find better way to initialize the parameters



$$\theta = \{W^1, b^1, W^2, b^2, \dots, W^L, b^L\}$$

# Activation Function

▸ Give NN nonlinearity property

▸ Can be regarded as "ON" (1) or "OFF" (0)

Input layer    Layer 1    Layer 2                 Layer L  Output layer



$$y = f(x) = \sigma(W^L \ldots \sigma(W^2 \sigma(W^1 X + b^1) + b^2) \ldots + b^L)$$

# Activation Function

▶ There are many kinds of activation

   ▶ relu, sigmoid, elu, etc……

▶ Usually, we use relu as first try on building neural network

$$\frac{1}{1 + e^{-x}}$$

**Sigmoid function**

$\sigma(z)$

$a$

$a = z$

$a = 0$

$z$

**RELU**

<span style="color:red">**We usually use this now !**</span>

# Example-initialization



0.1

-0.5

0.3

-0.2

0.4

0.1

-0.3

-0.2

0.5

Bias = 0.1

0.3

0.2

-0.2

Bias = 0.2

-0.1

0.1

0.2

Bias = 0.1

Bias = 0.1

Bias = 0.2

Bias = 0.1

\# of parameters = 21

# Example-feed data



0.6

0.2

0.1

-0.5

0.3

-0.2

0.4

0.1

-0.3

-0.2

0.5

Bias = 0.1

0.3

0.2

-0.2

Bias = 0.2

-0.1

0.1

0.2

Bias = 0.1

Bias = 0.1

Bias = 0.2

Bias = 0.1

# Example-forward pass



$$relu\left(\begin{bmatrix} 0.1 & -0.2 \\ -0.5 & 0.4 \\ 0.3 & 0.1 \end{bmatrix}\begin{bmatrix} 0.6 \\ 0.2 \end{bmatrix} + \begin{bmatrix} 0.1 \\ 0.2 \\ 0.1 \end{bmatrix}\right)$$

$$= relu\left(\begin{bmatrix} 0.12 \\ -0.02 \\ 0.3 \end{bmatrix}\right)$$

$$= \begin{bmatrix} 0.12 \\ 0 \\ 0.3 \end{bmatrix}$$

# Example-forward pass



$$relu\left(\begin{bmatrix} -0.3 & 0.3 & -0.1 \\ -0.2 & 0.2 & 0.1 \\ 0.5 & -0.2 & 0.2 \end{bmatrix}\begin{bmatrix} 0.12 \\ 0 \\ 0.3 \end{bmatrix} + \begin{bmatrix} 0.1 \\ 0.2 \\ 0.1 \end{bmatrix}\right)$$

$$= relu\left(\begin{bmatrix} -0.066 \\ 0.006 \\ 0.12 \end{bmatrix}\right)$$

$$= \begin{bmatrix} 0 \\ 0.006 \\ 0.12 \end{bmatrix}$$

# Define Loss Function

# Loss Function

‣ There are many kinds of loss function
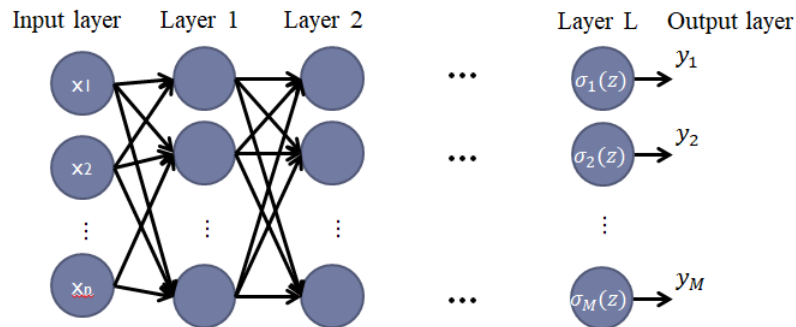
‣ Usually, it is a function that map multi-variables to a single value

‣ We will introduce two loss function in DNN

  ‣ Mean square

  ‣ Cross-entropy

# What is one-hot encode



$$\begin{bmatrix} 0.15 \\ 0.2 \\ 0.3 \\ \vdots \end{bmatrix}\quad\begin{matrix} \text{``0''} \\ \text{``1''} \\ \text{``2''} \end{matrix}$$

**Output of NN look like this**

$$\begin{bmatrix} 0 \\ 1 \\ 0 \\ \vdots \end{bmatrix}\quad\begin{matrix} \text{``0''} \\ \text{``1''} \\ \text{``2''} \end{matrix}$$

**What we want (one hot encode)**

# Mean Square

$$(y - \hat{y})^2$$

$$\begin{bmatrix} 0.15 \\ 0.2 \\ 0.3 \\ \vdots \end{bmatrix} \quad \begin{matrix} \text{``0''} \\ \text{``1''} \\ \text{``2''} \end{matrix} \qquad\qquad \begin{bmatrix} 0 \\ 1 \\ 0 \\ \vdots \end{bmatrix} \quad \begin{matrix} \text{``0''} \\ \text{``1''} \\ \text{``2''} \end{matrix}$$

# Cross-entropy with Softmax

▸ Information

  ▸ $log\left(\frac{1}{p_i}\right)$ where $p_i$ is probability of an event

| | |
|---|---|
| **Sun rises in the east tomorrow** | **It will rain tomorrow in Taiwan** |

Which is more informative?

# Cross-entropy with Softmax

▸ Entropy

   ▸ Expected value(mean) of information contained in each message

▸ Entropy can be seen as index of uncertainty

   ▸ Bigger mean more chaos

▸ Cross-entropy

   ▸ Measurement on the difference between two probability distribution

   ▸ Different distribution apply on entropy

   ▸ Cross-entropy is greater than entropy

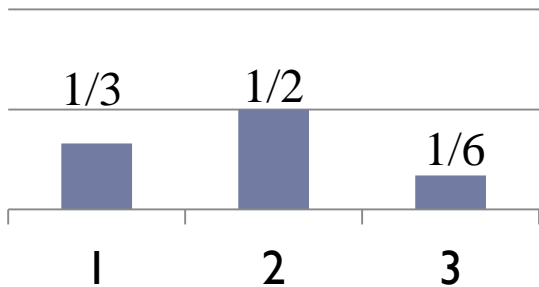$$H(y) = \sum_i y_i \log \frac{1}{y_i} = -\sum_i y_i \log y_i$$

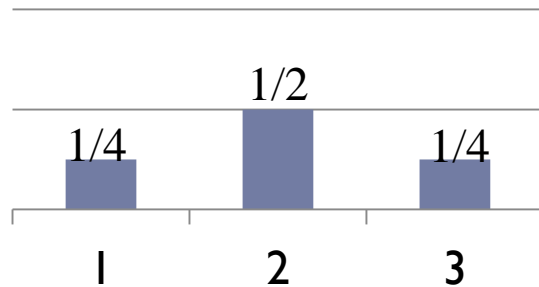Entropy

$$H(y) = -\sum_i y_i log \hat{y}_i$$

Cross-entropy

# Example

**Probability distribution 1**

$1/3$  $1/2$

$1/6$

1  2  3

**Probability distribution 2**

$1/2$

$1/4$  $1/4$

1  2  3

Entropy on distribution 1
$= 1/3 * \log(3) + 1/2 * \log(2) + 1/6 * \log(6)$

Entropy on distribution 2
$= 1/4 * \log(4) + 1/2 * \log(2) + 1/4 * \log(4)$

Cross-entropy on distribution 1 over distribution 2
$= 1/3 * \log(4) + 1/2 * \log(2) + 1/6 * \log(4)$

Cross-entropy on distribution 2 over distribution 1
$= 1/4 * \log(3) + 1/2 * \log(2) + 1/4 * \log(6)$

# Example

Entropy on distribution 1
$= 1/3 * \log(3) + 1/2 * \log(2) + 1/6 * \log(6)$
$= 0.439$

Entropy on distribution 2
$= 1/4 * \log(4) + 1/2 * \log(2) + 1/4 * \log(4)$
$= 0.452$

Cross-entropy on distribution 1 over distribution 2
$= 1/3 * \log(4) + 1/2 * \log(2) + 1/6 * \log(4) = 0.456$

Cross-entropy on distribution 2 over distribution 1
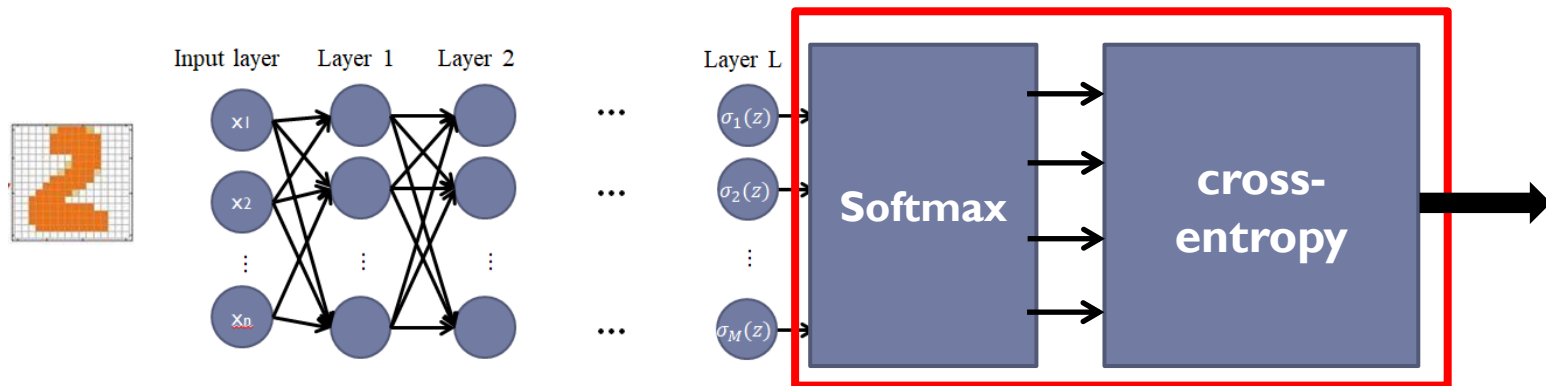$= 1/4 * \log(3) + 1/2 * \log(2) + 1/4 * \log(6) = 0.464$

- Cross-entropy is greater than entropy
  Cross-entropy on distribution 1 over 2 > Entropy on distribution 1
  Cross-entropy on distribution 2 over 1 > Entropy on distribution 2
- If two distribution become closer
  - Value of cross-entropy is closer to entropy

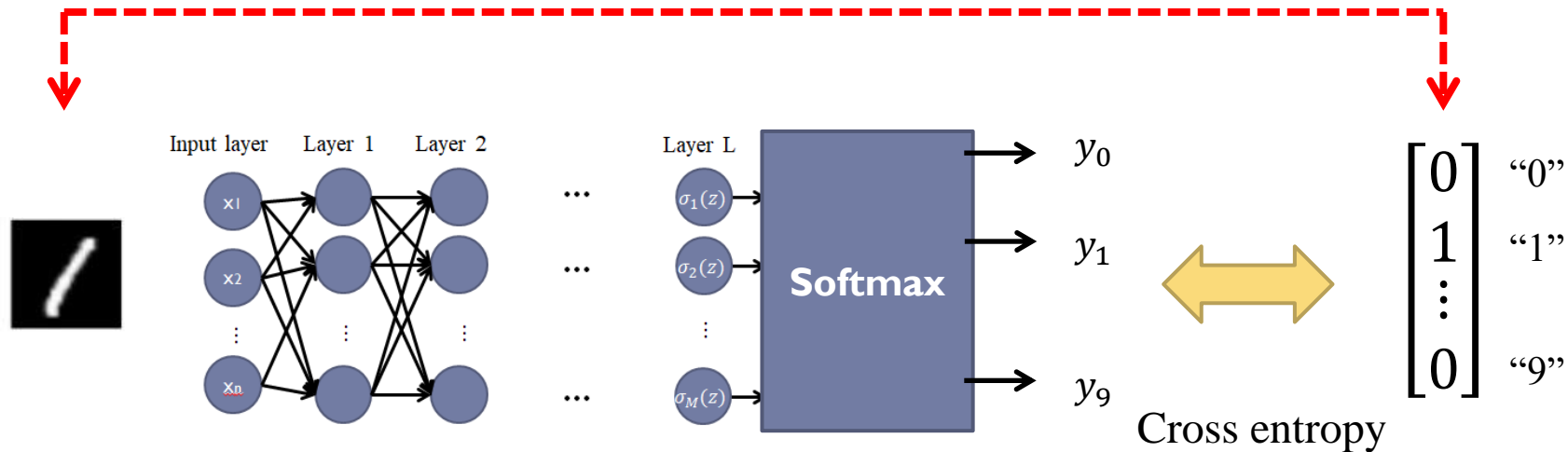# Cross-entropy with Softmax

▸ Cross-entropy usually come with softmax layer in NN

▸ Softmax function squash all of elements in vector to [0, 1]

$$\sigma(\mathbf{z})_j = \frac{e^{z_j}}{\sum_{k=1}^{K} e^{z_k}} \quad \text{for } j = 1, \ldots, K.$$

# Cross-entropy with Softmax



Input layer    Layer 1    Layer 2 ... Layer L

Softmax

$y_0$
$y_1$
$y_9$

$$\begin{bmatrix} 0 \\ 1 \\ \vdots \\ 0 \end{bmatrix}$$ "0" "1" "9"

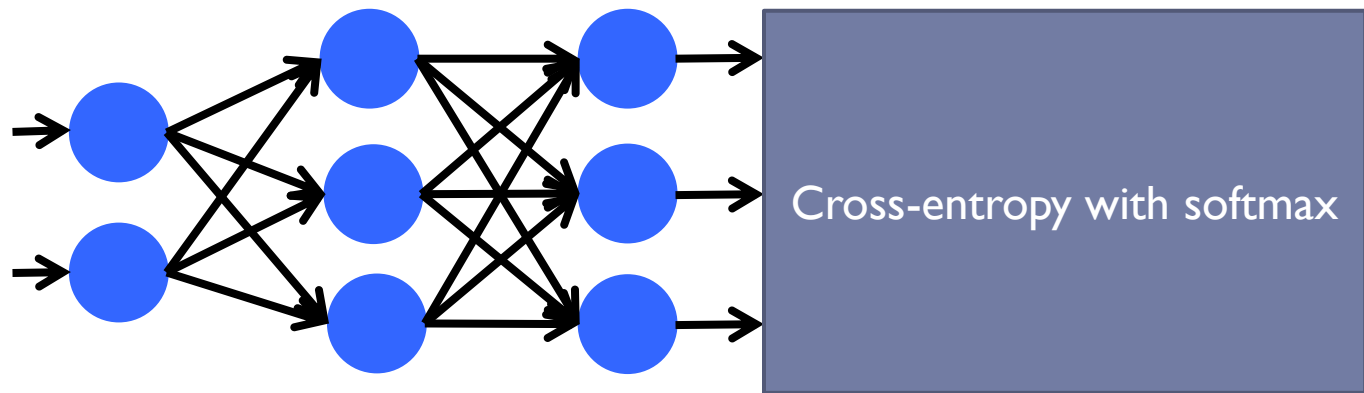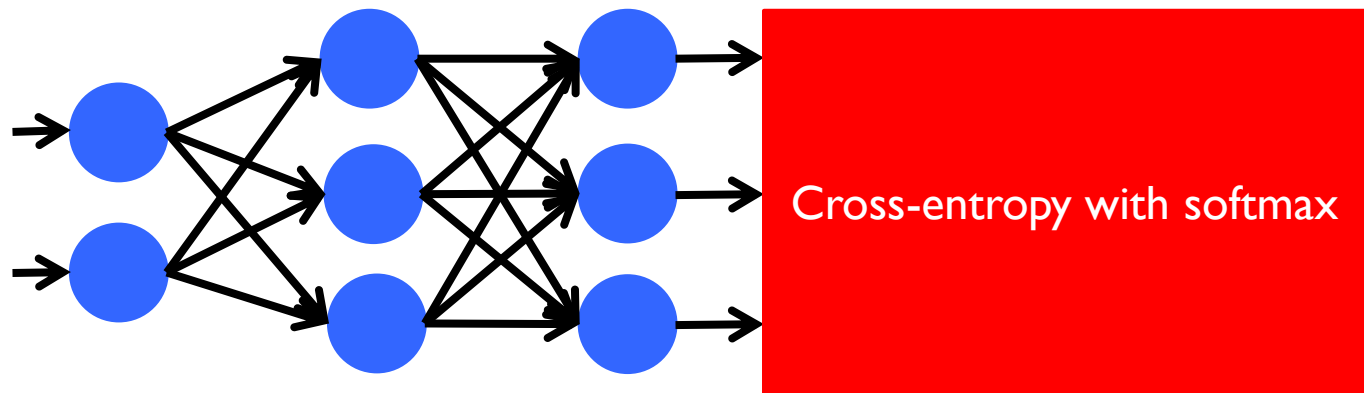Cross entropy

$$H(y) = -\sum_i y_i log\hat{y}_i$$

$$Cost = -(0 * \log(y_0) + 1 * \log(y_1) + 0 * \log(y_2) + \cdots + 0 * \log(y_9))$$

# Example-forward pass

Cross-entropy with softmax

# Example-forward pass



Cross-entropy with softmax
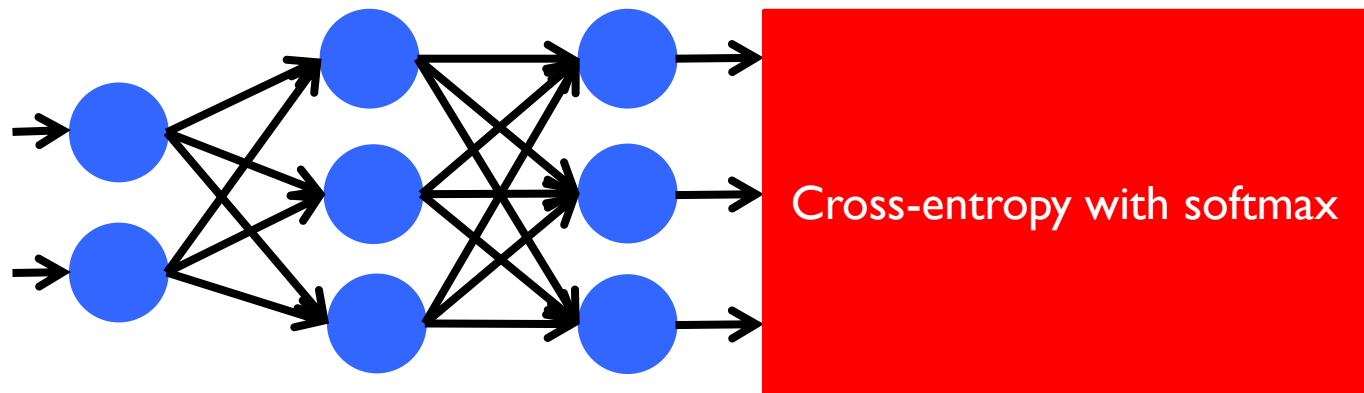
$$\text{Softmax}\left(\begin{bmatrix} 0 \\ 0.006 \\ 0.12 \end{bmatrix}\right) = \begin{bmatrix} 0.319 \\ 0.321 \\ 0.36 \end{bmatrix}$$

$$\sigma(\mathbf{z})_j = \frac{e^{z_j}}{\sum_{k=1}^{K} e^{z_k}} \quad \text{for } j = 1, \ldots, K.$$

# Example-forward pass



What we expect
(label)

$$\begin{bmatrix} 0 \\ 1 \\ 0 \end{bmatrix}$$

Cross-entropy with softmax

- 0 * ln(0.319) – 1 * ln(0.321) – 0 * ln(0.36)
= 1.1363

$$-\sum_{i=0}^{class\ \#} \hat{y}_i \ln(y_i)$$

This is large at the beginning. During training, this should decrease.

# How to Feed data

Training Data             Testing Data

| | |
|---|---|
| 80% | 20% |

**x:**      **y: "2" (label)**

# How to Feed data

| How to feed data | Mean Square | Cross-entropy |
|---|---|---|
| All data at a time | $\dfrac{1}{N}\sum_{i=1}^{N}(y_i - \hat{y}_i)^2$ | $\dfrac{1}{N}\sum_{j=1}^{N}(-\sum_{i=0}^{class\,\#}\hat{y}_i\ln(y_i))$ |
| One data at a time | $(y - \hat{y})$ | $-\sum_{i=0}^{class\,\#}\hat{y}_i\ln(y_i)$ |
| Batch of data a time (B < N) | $\dfrac{1}{B}\sum_{i=1}^{B}(y_i - \hat{y}_i)^2$ | $\dfrac{1}{B}\sum_{j=1}^{B}(-\sum_{i=0}^{class\,\#}\hat{y}_i\ln(y_i))$ |

# How to Feed data

**Pick one training data at a time**



Input layer    Layer 1    Layer 2     Layer L

$x1$   $x2$   $xn$   $\sigma_1(z)$   $\sigma_2(z)$   $\sigma_M(z)$

**Softmax**    **cross-entropy**

$$-\sum_{i=0}^{class\ \#} \hat{y}_i \ln(y_i)$$

# How to Feed data

**Pick all training data at a time**

All training data



$$\frac{1}{N}\sum_{j=1}^{N}\left(-\sum_{i=0}^{class\ \#}\hat{y}_i\ln(y_i)\right)$$

# How to Feed data

A batch of data
(you can define)

**Pick a batch of data at a time**



$$\frac{1}{B}\sum_{j=1}^{B}(-\sum_{i=0}^{class\ \#}\hat{y}_i\ln(y_i))$$

# Optimization

# Optimization

- We will introduce these Optimization methods
  - Gradient decent
  - Gradient decent with moment
  - Adagrad
  - RMSprop
  - Adam

# Gradient Descent

▶ An algorithm that find the minimum of a function



$Randomly\ select\ \theta_1\ as\ start\ point$

$$Compute\ \frac{dC(\theta_1)}{d\theta}$$

$$\theta_2 \leftarrow \theta_1 - \eta\frac{dC(\theta_1)}{d\theta}$$

$$Compute\ \frac{dC(\theta_2)}{d\theta}$$

$$\theta_3 \leftarrow \theta_2 - \boxed{\eta}\frac{dC(\theta_2)}{d\theta}$$

Learning rate

# Gradient Descent



$$\theta = \begin{bmatrix} x \\ y \end{bmatrix} \qquad \nabla C(\theta) = \begin{bmatrix} dz/dx \\ dz/dy \end{bmatrix}$$

*Randomly select $\theta_1$ as start point*

*Compute $\nabla C(\theta_1)$*

$\theta_2 \leftarrow \theta_1 - \eta \nabla C(\theta_1)$

*Compute $\nabla C(\theta_2)$*

$\theta_3 \leftarrow \theta_2 - \eta \nabla C(\theta_2)$

# Gradient Descent



cost

too big

too small

just make

# of parameters updates

# Gradient Descent



Cost

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

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

# Gradient Descent with Momentum

$Randomly\ select\ \theta_1\ as\ start\ point$

$Compute\ \nabla\theta_1$
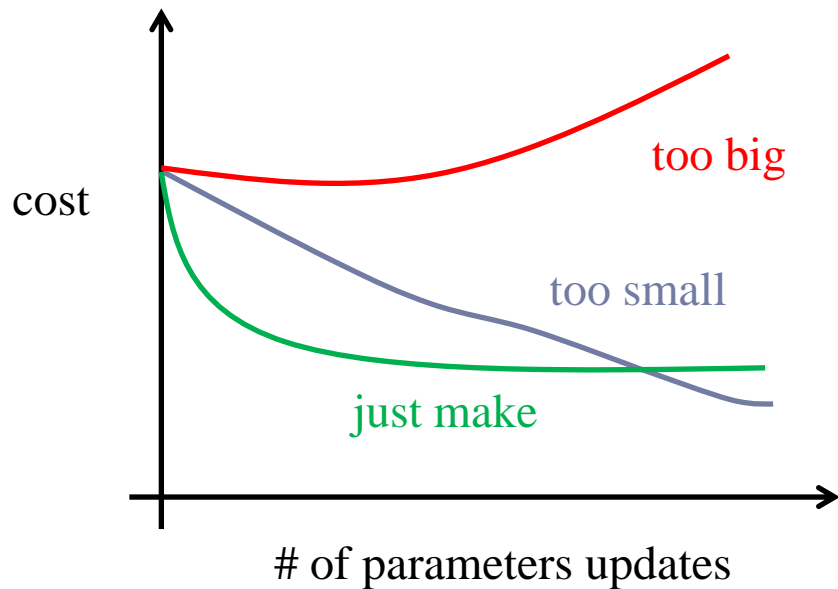
$\theta_2 \leftarrow \theta_1 - \eta\nabla\theta_1$

$Compute\ \nabla\theta_2$

$\theta_3 \leftarrow \theta_2 - \eta\nabla\theta_2$

$\vdots$

**Without momentum**

$Randomly\ select\ \theta_1\ as\ start\ point$
$,initialize\ v_1 = 0$

$Compute\ \nabla\theta_1, v_2 = \lambda v_1 - \eta\nabla\theta_1$

$\theta_2 \leftarrow \theta_1 + v_2$

$Compute\ \nabla\theta_2, v_3 = \lambda v_2 - \eta\nabla\theta_2$

$\theta_3 \leftarrow \theta_2 + v_3$

$\vdots$

**With momentum**

# Gradient Descent

- Drawback of gradient descent
  - Learning rate would not decay over time
  - Direction of learning rate is fixed
- Adaptive learning rate is needed

# Adagrad

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

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

$$\eta^t = \frac{\eta}{\sqrt{t+1}} \qquad g^t = \frac{\partial L(\theta^t)}{\partial w}$$

$$w^3 \leftarrow w^2 - \frac{\eta^2}{\sigma^2} g^2 \qquad \sigma^2 = \sqrt{\frac{1}{3}[(g^0)^2 + (g^1)^2 + (g^2)^2]}$$

$$\vdots$$

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

# RMSProp

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

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

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

$$\vdots$$

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

# 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)
$t \leftarrow 0$ (Initialize timestep) $\searrow$ for RMSprop
**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)
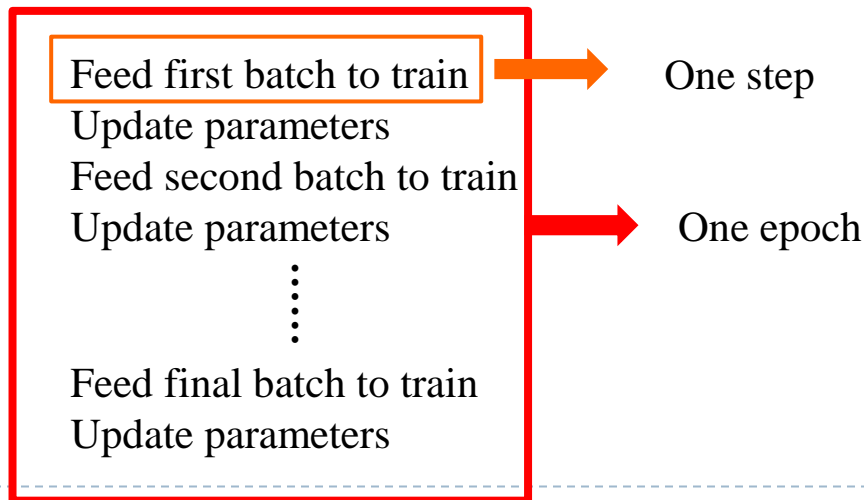  $\theta_t \leftarrow \theta_{t-1} - \alpha \cdot \widehat{m}_t/(\sqrt{\widehat{v}_t} + \epsilon)$ (Update parameters)
**end while**
**return** $\theta_t$ (Resulting parameters)

# What is one epoch/step

▸ **One epoch**

  ▸ One pass of all the training data

▸ **One step**

  ▸ Pass a batch of training data (batch size is user defined)

Feed first batch to train → One step

Update parameters
Feed second batch to train
Update parameters → One epoch

⋮

Feed final batch to train
Update parameters

# Backpropagation

▶ An efficient way to compute gradient given a set of variable

$$\frac{\partial C}{\partial w_{ij}^l} \quad \textcolor{red}{?}$$

# Backpropagation

$$y = f(x) = \sigma(W^L \dots \sigma(W^2 \sigma(W^1 X + b^1) + b^2) \dots + b^L)$$

$$z^1 = W^1 X + b^1$$
$$a^1 = \sigma(z^1)$$
$$z^2 = W^2 a^1 + b^2$$
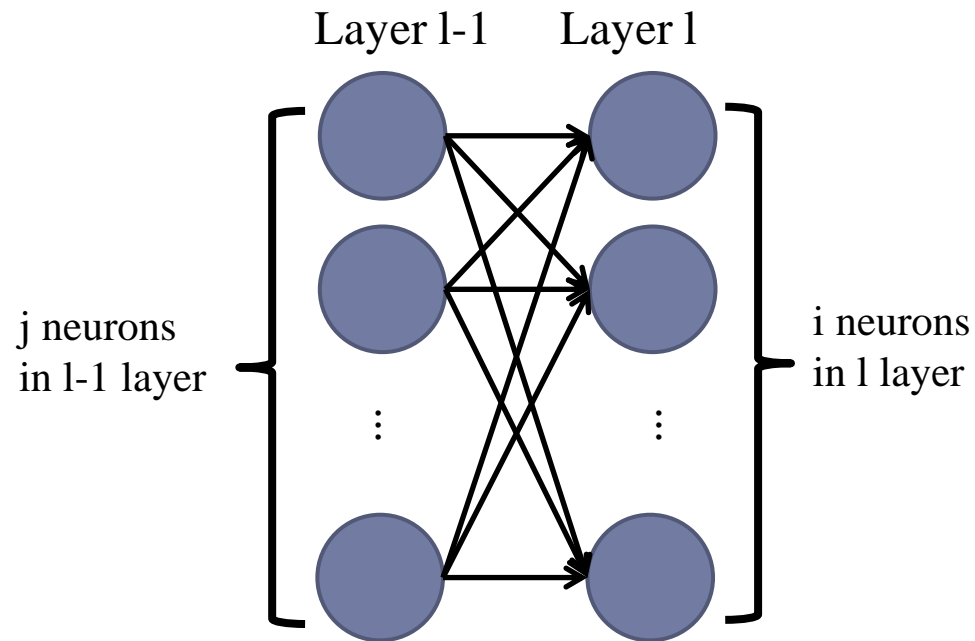$$a^2 = \sigma(z^2)$$
$$\vdots$$
$$z^l = W^l a^{l-1} + b^l$$
$$a^l = \sigma(z^l)$$

$$\frac{\partial C}{\partial w_{ij}^l} = \frac{\partial C}{\partial z_i^l} * \frac{\partial z_i^l}{\partial w_{ij}^l}$$

1. Calculate $\dfrac{\partial z_i^l}{\partial w_{ij}^l}$

2. Calculate $\dfrac{\partial C}{\partial z_i^l}$ (**error signal** $\delta_i^l$ )

# Backpropagation



Layer l-1    Layer l

j neurons
in l-1 layer

i neurons
in l layer

$$\frac{\partial z_i^l}{\partial w_{ij}^l}$$

$if\ l = 1\ (input\ layer\ \rightarrow\ layer\ 1)$

$$z_i^1 = \sum_j w_{ij}^1 x_j + b_i^1 \qquad \frac{\partial z_i^1}{\partial w_{ij}^1} = x_j$$

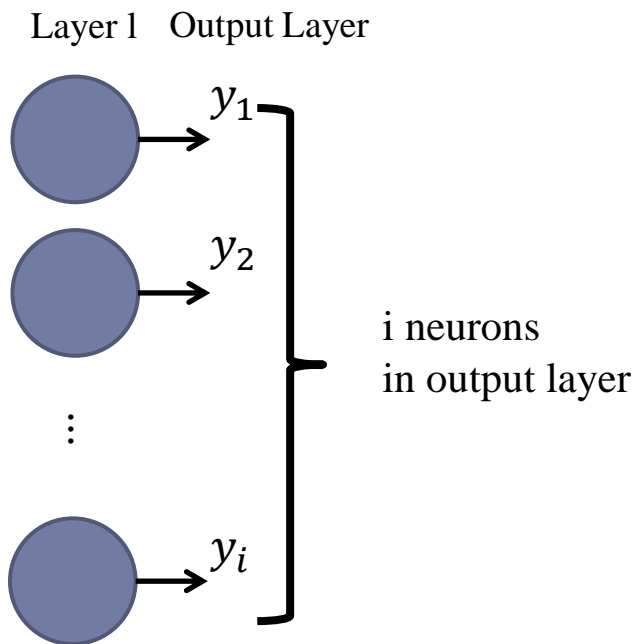$if\ l > 1\ (layer\ l - 1\ \rightarrow\ layer\ 1)$

$$z_i^l = \sum_j w_{ij}^l a_j^{l-1} + b_i^l \qquad \frac{\partial z_i^l}{\partial w_{ij}^l} = a_j^{l-1}$$

# Backpropagation

$$\frac{\partial C}{\partial z_i^l} \quad (\textbf{error signal} \quad \delta_i^l \ )$$

1. Compute $\delta^L$ (L layer error signal)
2. Compute relation between $\delta^l$ and $\delta^{l+1}$

# Backpropagation

Layer 1    Output Layer

$$\frac{\partial C}{\partial z_i^l} \quad (\textbf{error signal} \quad \delta_i^l \ )$$

$y_1$

$y_2$

i neurons
in output layer

$\vdots$

$y_i$

1. Compute $\delta^L$ (L layer error signal)

$$\delta_i^L = \frac{\partial C}{\partial z_i^L} = \frac{\partial y_i}{\partial z_i^L} * \frac{\partial C}{\partial y_i} = \sigma'(z_i^L) * \frac{\partial C}{\partial y_i}$$

# Backpropagation

Layer 1    Output Layer

$y_1$

$y_2$

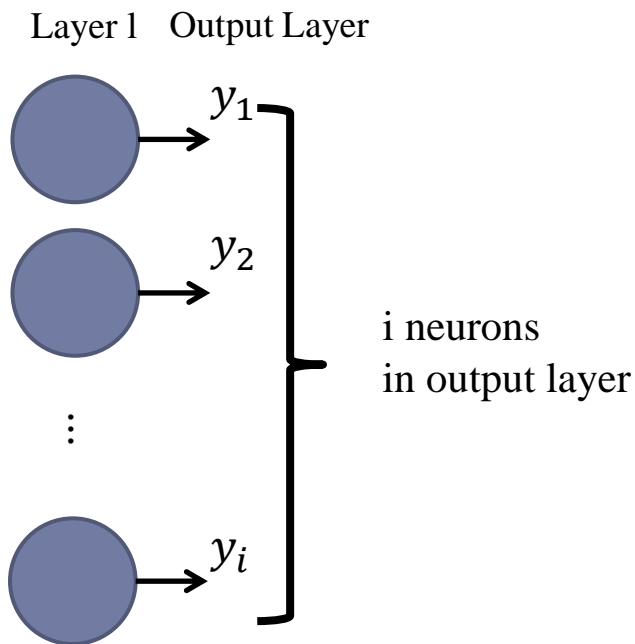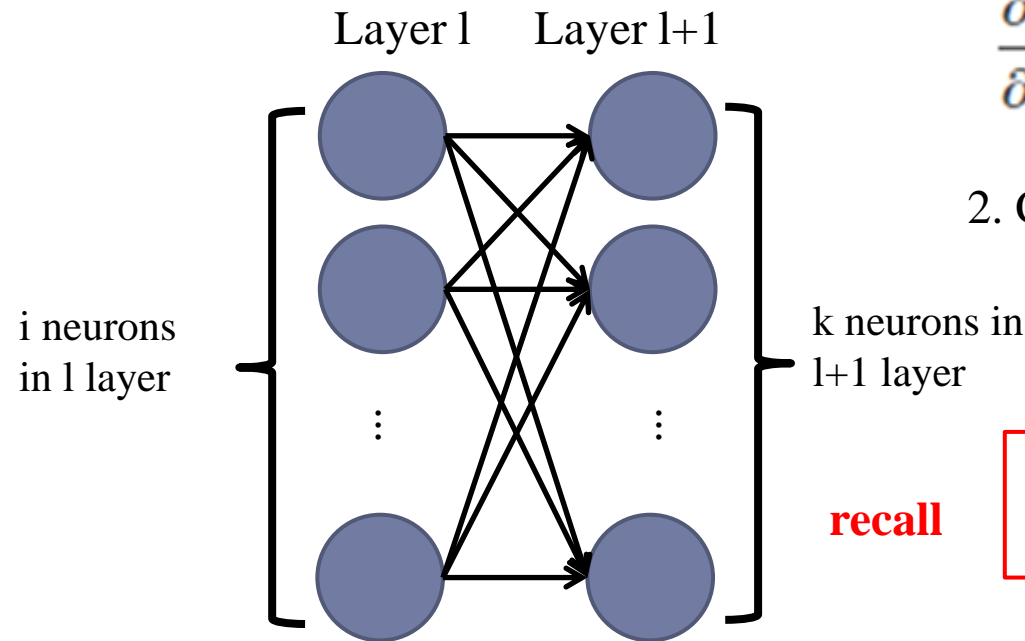$\vdots$

$y_i$

i neurons
in output layer

$$\frac{\partial C}{\partial z_i^l} \quad (\textbf{error signal} \quad \delta_i^l \; )$$

1. Compute $\delta^L$ (L layer error signal)

$$\delta^L = \sigma'(z^L) \bullet \nabla C(y)$$

$$\sigma'(z^L) = \begin{bmatrix} \sigma'(z_1^L) \\ \sigma'(z_2^L) \\ \vdots \\ \sigma'(z_n^L) \end{bmatrix} \qquad \nabla C'(y) = \begin{bmatrix} \frac{\partial C}{\partial y_1} \\ \frac{\partial C}{\partial y_2} \\ \vdots \\ \frac{\partial C}{\partial y_n} \end{bmatrix}$$

# Backpropagation

Layer l    Layer l+1



i neurons
in l layer

k neurons in
l+1 layer

$\dfrac{\partial C}{\partial z_i^l}$  (**error signal**  $\delta_i^l$ )

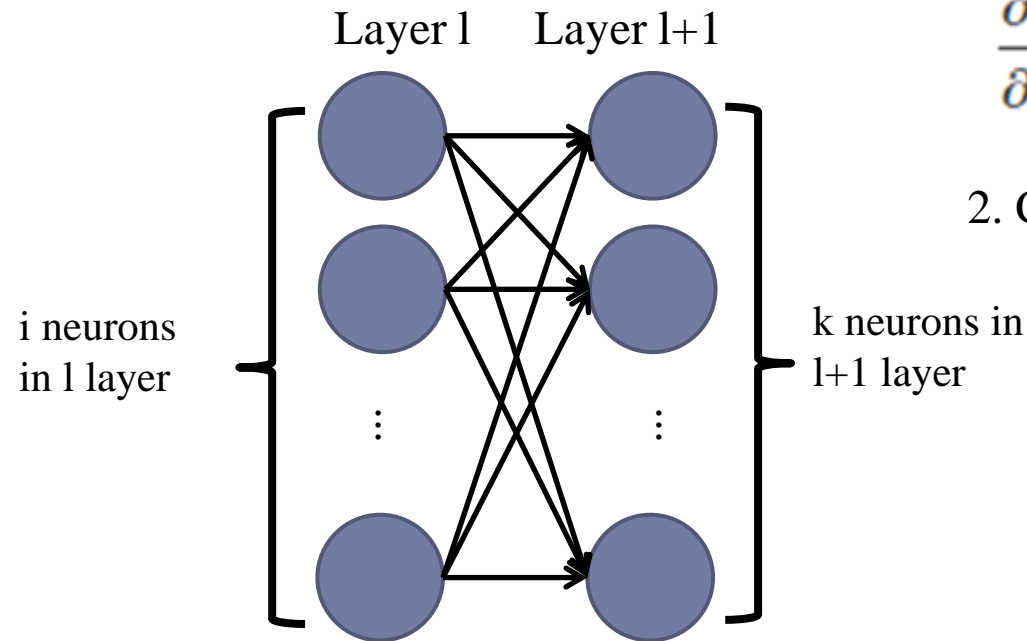2. Compute relation between $\delta^l$ and $\delta^{l+1}$

$$\delta_i^l = \frac{\partial C}{\partial z_i^l} = \frac{\partial a_i^l}{\partial z_i^l} \sum_k \frac{\partial z_k^{l+1}}{\partial a_i^l} \frac{\partial C}{\partial z_k^{l+1}}$$

**recall**

$$\frac{\partial a_i^l}{\partial z_i^l} = \sigma'(z_i^l), \; \frac{\partial z_k^{l+1}}{\partial a_i^l} = w_{ki}^{l+1}, \; \frac{\partial C}{\partial z_k^{l+1}} = \delta_k^{l+1}$$

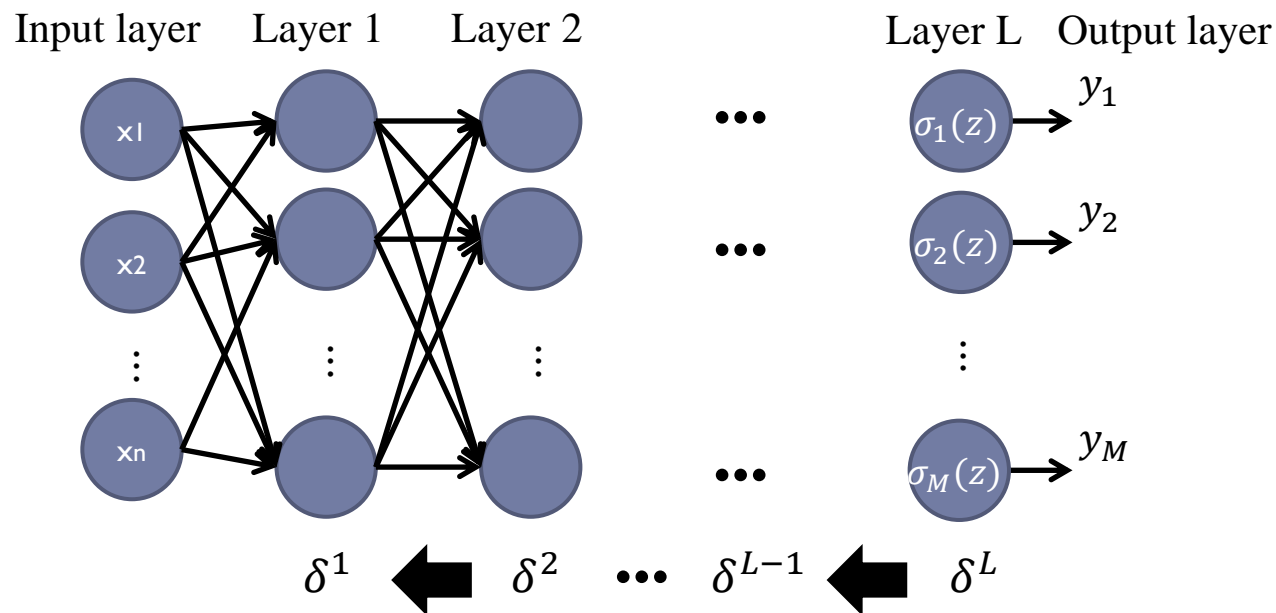$$\delta_i^l = \sigma'(z_i^l) \sum_k w_{ki}^{l+1} \delta_k^{l+1}$$

# Backpropagation

Layer l    Layer l+1

i neurons
in l layer

k neurons in
l+1 layer

$$\frac{\partial C}{\partial z_i^l} \quad (\textbf{error signal} \quad \delta_i^l \,)$$
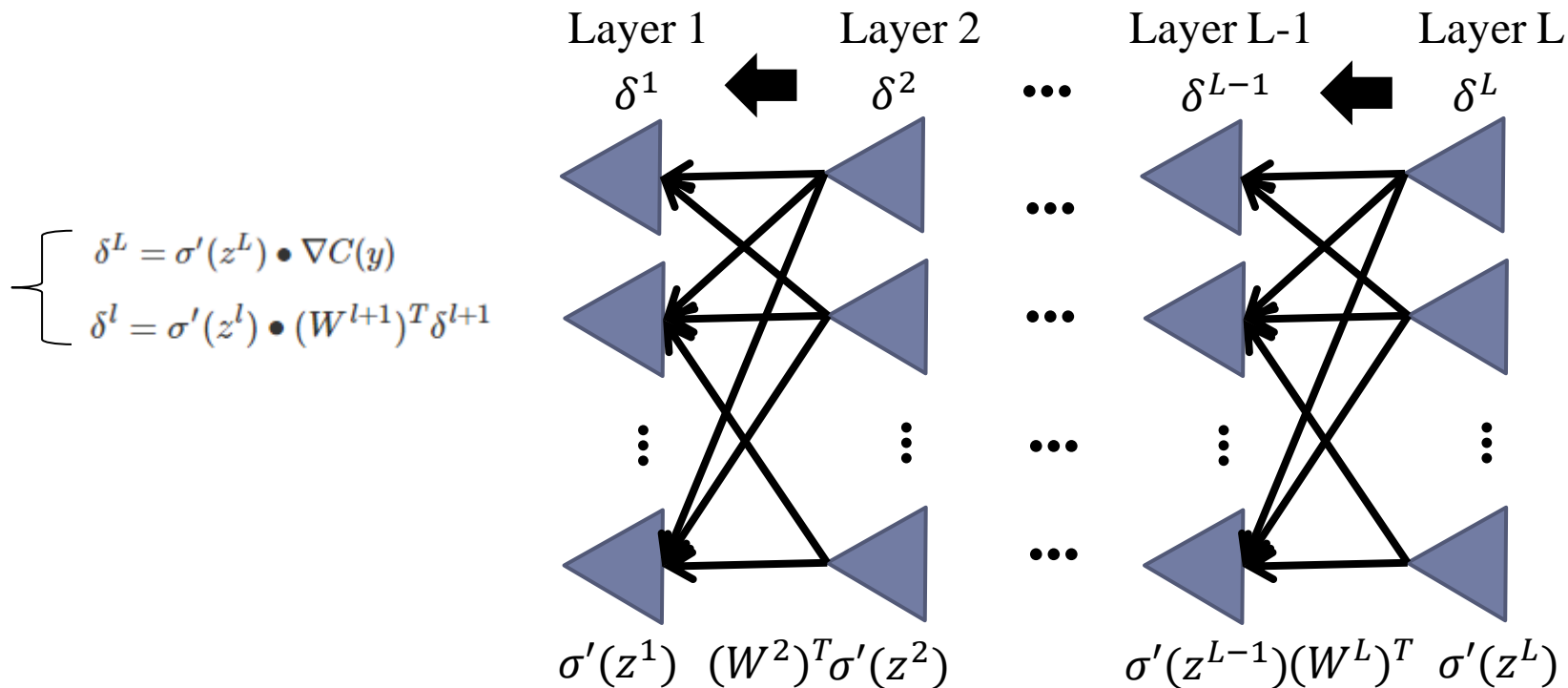
2. Compute relation between $\delta^l$ and $\delta^{l+1}$

$$\delta^l = \sigma'(z^l) \bullet (W^{l+1})^T \delta^{l+1}$$

# Backpropagation

$$\delta^L = \sigma'(z^L) \bullet \nabla C(y)$$

$$\delta^l = \sigma'(z^l) \bullet (W^{l+1})^T \delta^{l+1}$$

Input layer    Layer 1    Layer 2        Layer L    Output layer

x1

x2

$\vdots$

xn

$\sigma_1(z)$ → $y_1$

$\sigma_2(z)$ → $y_2$

$\vdots$

$\sigma_M(z)$ → $y_M$

•••         •••

$\delta^1$ ⬅ $\delta^2$ ••• $\delta^{L-1}$ ⬅ $\delta^L$

# Backpropagation

Layer 1　　　　Layer 2　　　　Layer L-1　　　Layer L

$\delta^1$　　　　　$\delta^2$　　　•••　　$\delta^{L-1}$　　　$\delta^L$

$$\delta^L = \sigma'(z^L) \bullet \nabla C(y)$$

$$\delta^l = \sigma'(z^l) \bullet (W^{l+1})^T \delta^{l+1}$$

$\sigma'(z^1)$　$(W^2)^T \sigma'(z^2)$　　　$\sigma'(z^{L-1})(W^L)^T$　$\sigma'(z^L)$

▶ 65

# Backpropagation

$$\frac{\partial C}{\partial w_{ij}^l} = \frac{\partial C}{\partial z_i^l} * \frac{\partial z_i^l}{\partial w_{ij}^l}$$

$$\delta^L = \sigma'(z^L) \bullet \nabla C(y)$$

$$\delta^{L-1} = \sigma'(z^{L-1}) \bullet (W^L)^T \delta^L$$

$$\vdots$$

$$\delta^{l-1} = \sigma'(z^{l-1}) \bullet (W^l)^T \delta^l$$

$$\vdots$$

$if\ l = 1\ (input\ layer\ \rightarrow layer\ 1)$

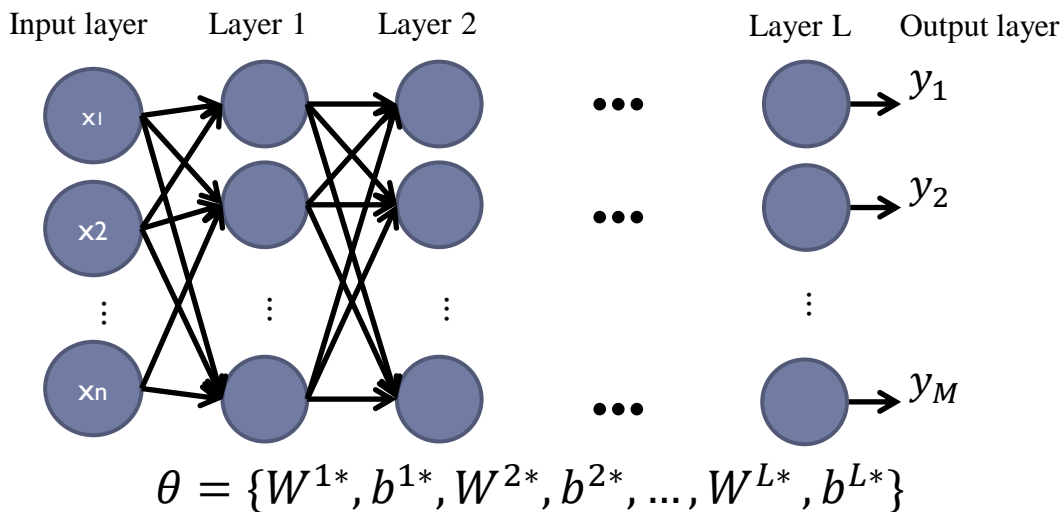$$z_i^1 = \sum_j w_{ij}^1 x_j^r + b_i^1 \qquad \frac{\partial z_i^1}{\partial w_{ij}^1} = x_j$$

$if\ l > 1 (layer\ l - 1\ \rightarrow layer\ 1)$

$$z_i^l = \sum_j w_{ij}^l a_j^{l-1} + b_i^l \qquad \frac{\partial z_i^l}{\partial w_{ij}^l} = a_j^{l-1}$$

# Predict/Validate Result

# Predict/Validate Result

▸ "*" mean a better set of parameters we have found

▸ Use test data to find model accuracy (supervised learning)

  ▸ # correct predict / # total data

| Input layer | Layer 1 | Layer 2 | | Layer L | Output layer |

$$\theta = \{W^{1*}, b^{1*}, W^{2*}, b^{2*}, \ldots, W^{L*}, b^{L*}\}$$

# Overfitting



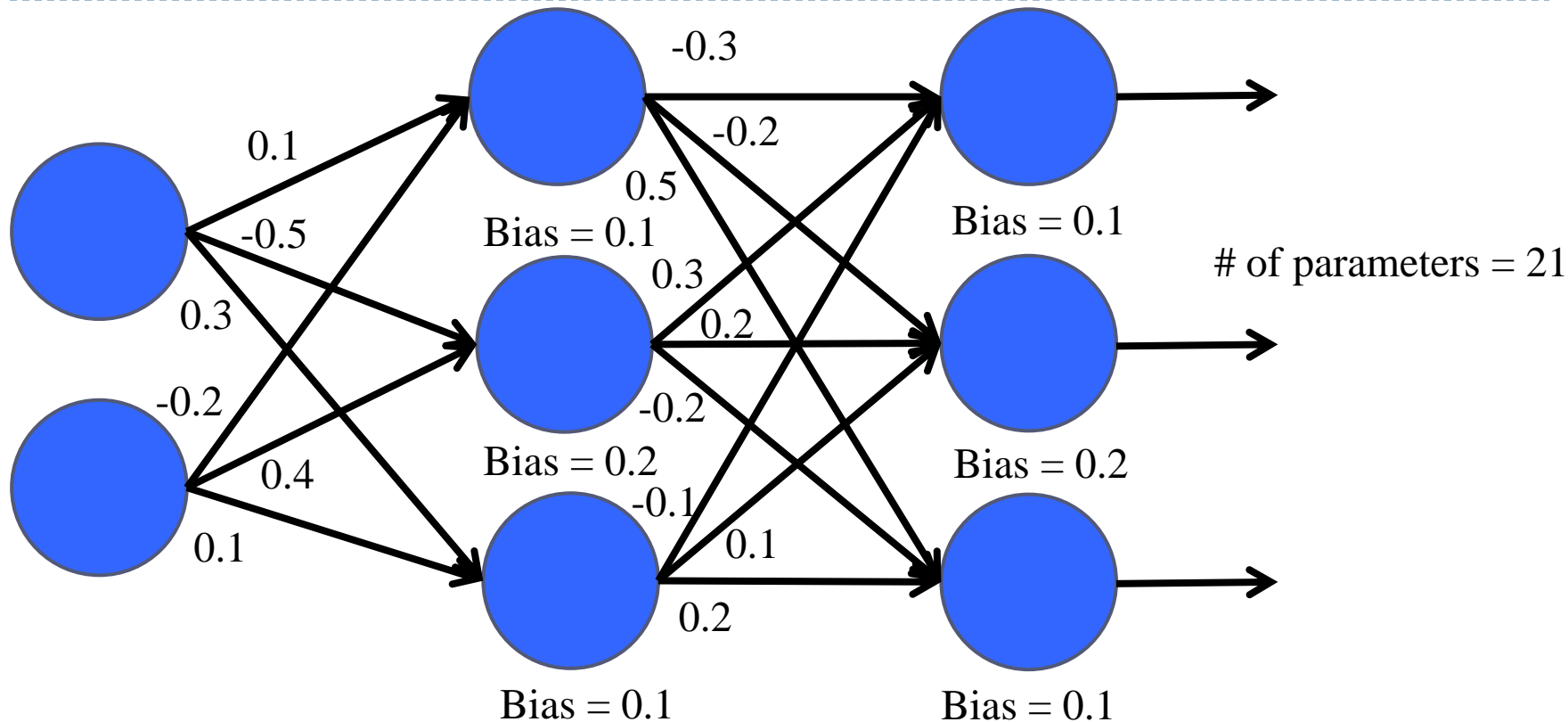Training error is going lower but testing error is not

# Overfitting

- **Prevent overfitting**
  - Use more data
  - Stronger regularization
  - Data augmentation
  - Reduce complexity of model

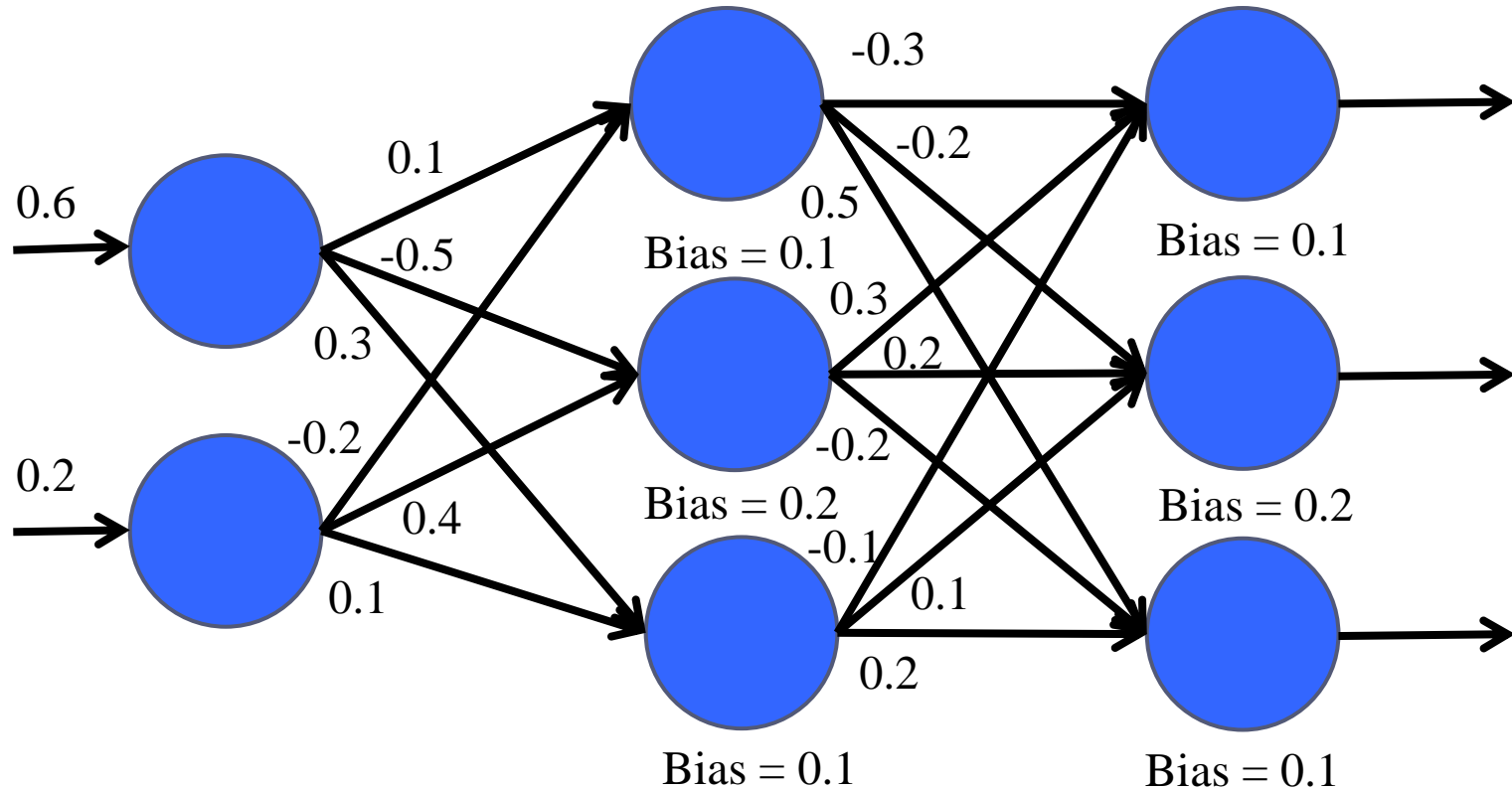# Quick Recap in this Lecture

1. Build DNN model
2. Define loss (mean square/cross entropy with softmax)
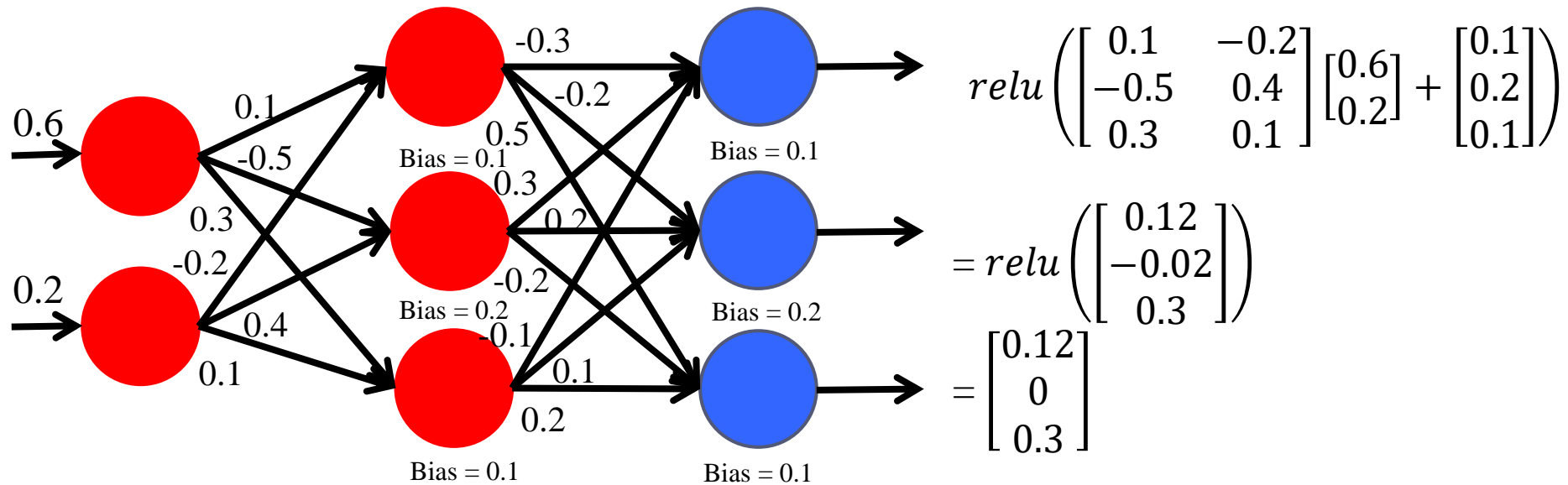3. Optimization (backpropagation)
4. Validate result
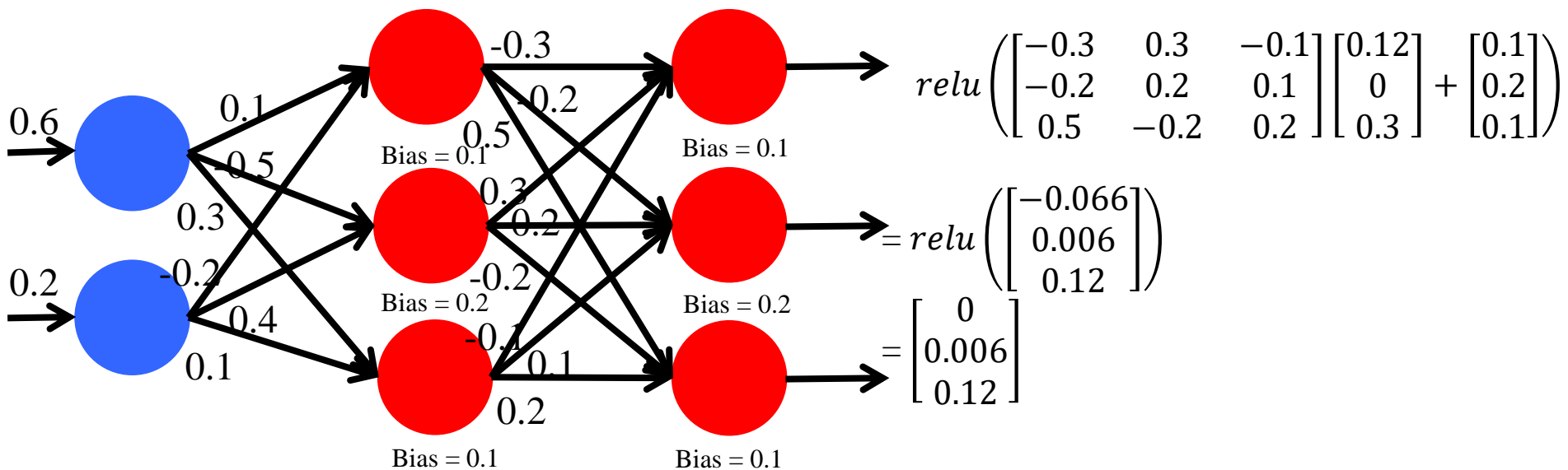
# Example-initialization



-0.3

0.1

-0.2

-0.5

0.5

Bias = 0.1

Bias = 0.1

0.3

0.3

# of parameters = 21

-0.2

0.2

0.4

-0.2

Bias = 0.2

Bias = 0.2

0.1

-0.1

0.1

0.2

Bias = 0.1

Bias = 0.1

▶ 72

# Example-feed data



0.6

0.2

0.1

-0.5

0.3

-0.2

0.4

0.1

-0.3

-0.2

0.5

0.3

0.2

-0.2

-0.1

0.1

0.2

Bias = 0.1

Bias = 0.2

Bias = 0.1

Bias = 0.1

Bias = 0.2

Bias = 0.1

# Example-forward pass



$$relu\left(\begin{bmatrix} 0.1 & -0.2 \\ -0.5 & 0.4 \\ 0.3 & 0.1 \end{bmatrix}\begin{bmatrix} 0.6 \\ 0.2 \end{bmatrix} + \begin{bmatrix} 0.1 \\ 0.2 \\ 0.1 \end{bmatrix}\right)$$

$$= relu\left(\begin{bmatrix} 0.12 \\ -0.02 \\ 0.3 \end{bmatrix}\right)$$

$$= \begin{bmatrix} 0.12 \\ 0 \\ 0.3 \end{bmatrix}$$

# Example-forward pass



$$relu\left(\begin{bmatrix} -0.3 & 0.3 & -0.1 \\ -0.2 & 0.2 & 0.1 \\ 0.5 & -0.2 & 0.2 \end{bmatrix}\begin{bmatrix} 0.12 \\ 0 \\ 0.3 \end{bmatrix} + \begin{bmatrix} 0.1 \\ 0.2 \\ 0.1 \end{bmatrix}\right)$$

$$= relu\left(\begin{bmatrix} -0.066 \\ 0.006 \\ 0.12 \end{bmatrix}\right)$$

$$= \begin{bmatrix} 0 \\ 0.006 \\ 0.12 \end{bmatrix}$$

# Example-forward pass



Cross-entropy with softmax

# Example-forward pass



Cross-entropy with softmax
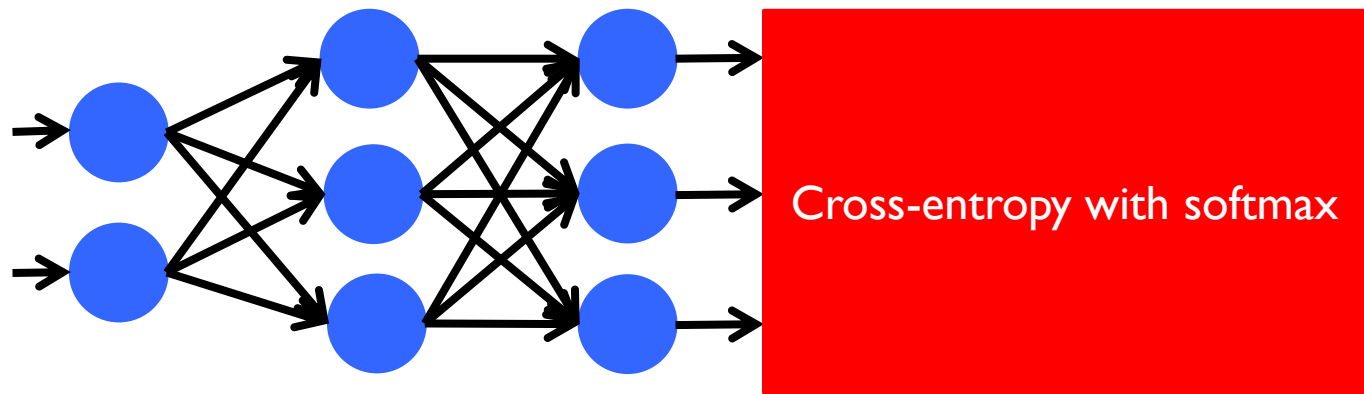
$$\text{Softmax}\left(\begin{bmatrix} 0 \\ 0.006 \\ 0.12 \end{bmatrix}\right) = \begin{bmatrix} 0.319 \\ 0.321 \\ 0.36 \end{bmatrix}$$

$$\sigma(\mathbf{z})_j = \frac{e^{z_j}}{\sum_{k=1}^{K} e^{z_k}} \quad \text{for } j = 1, \ldots, K.$$

# Example-forward pass



Cross-entropy with softmax

What we expect (label)

$$\begin{bmatrix} 0 \\ 1 \\ 0 \end{bmatrix}$$

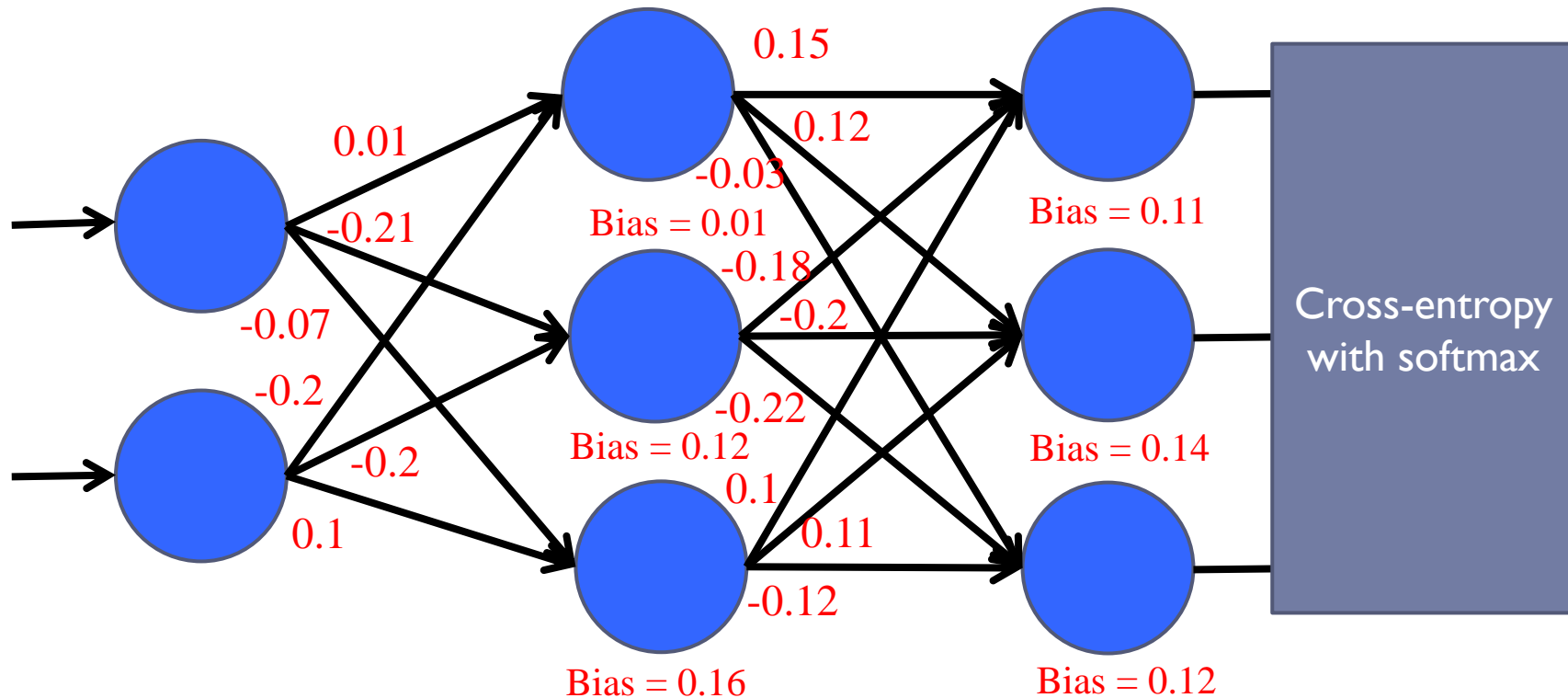$$- 0 * \ln(0.319) - 1 * \ln(0.321) - 0 * \ln(0.36)$$
$$= 1.1363$$

$$- \sum_{i=0}^{class \#} \hat{y}_i \ln(y_i)$$

This is large at the beginning. During training, this should decrease.
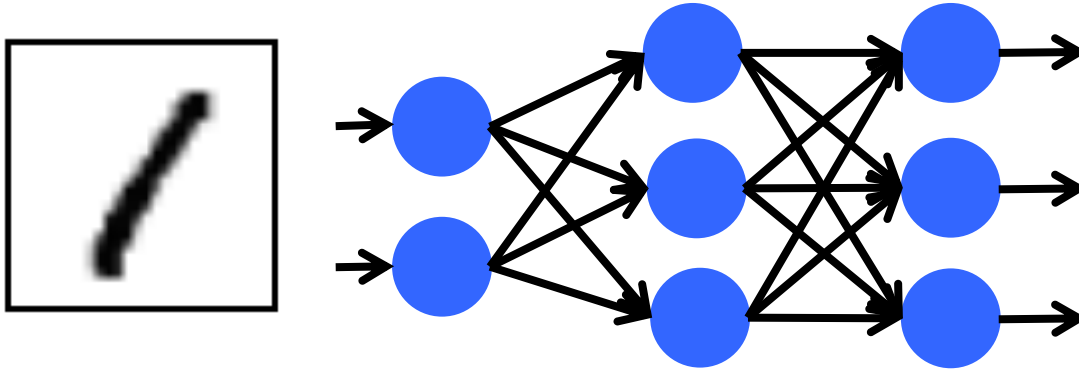
# Example-optimization

Use optimizer to optimize

and wait……

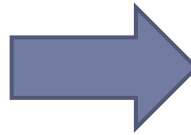# Example-after optimization

# Example-predict



Feed you test data
You may get something like this

$$\begin{bmatrix} 0.1 \\ 0.5 \\ 0.4 \end{bmatrix}$$

Model predict this is class #2
(0.5 is greater than others)