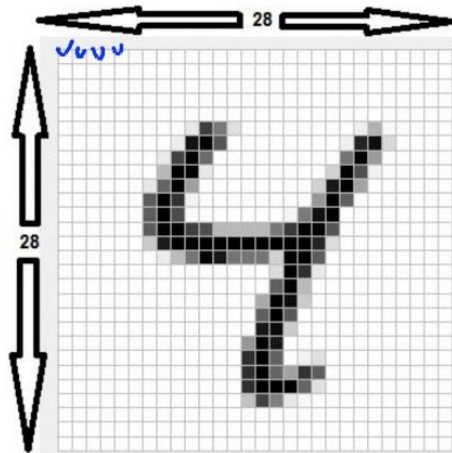# Ch. 2-4, 2-5

Pyo, Sanghun

# 2.5 Reviewing the first example

```python
from keras import models
from keras import layers
from keras.datasets import mnist
from keras.utils import to_categorical
from keras.utils import plot_model

(train_images, train_labels), (test_images, test_labels) = mnist.load_data()
```

MNIST data set

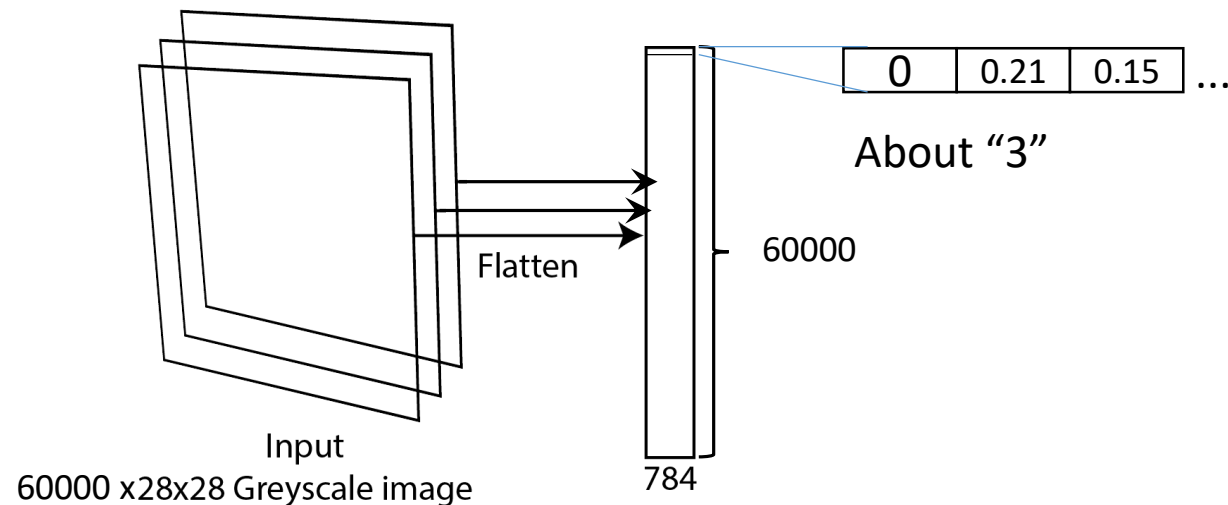

28X28X1 image file of handwritten digit number

Layer
- A key component of neural networks which is a kind of data processing filter.
- Extract more meaningful expressions from input data for a given problem.

# 2.5 Reviewing the first example

```
print(train_images.shape)        -> (60000, 28, 28)
print(len(train_labels))         -> 60000
print(train_labels)              -> [3, 4, 2, 1....]
print(test_images.shape)         -> (10000, 28, 28)
print(len(test_labels))          -> 10000
print(test_labels)


train_images = train_images.reshape((60000, 28 * 28))
train_images = train_images.astype('float32') / 255      # normalization
print("train_images_reshape?")
print(train_images.shape)
```
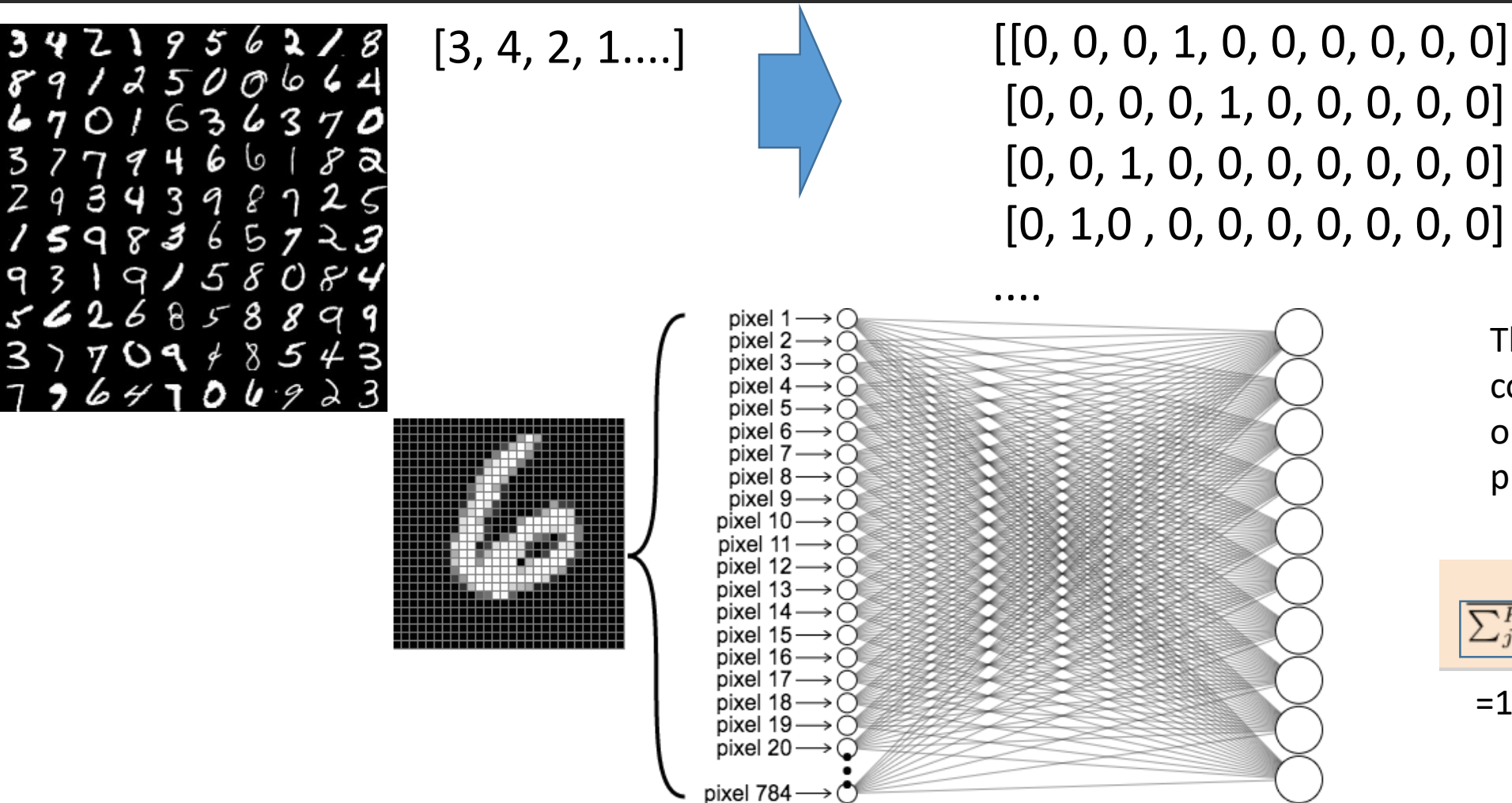


| 0 | 0.21 | 0.15 | ... |

About "3"

Flatten

60000

Input
60000 x28x28 Greyscale image

784

# 2.5 Reviewing the first example

```
train_labels = to_categorical(train_labels)
print("train_labels_categorical")
print(train_labels)
test_labels = to_categorical(test_labels)
```



[3, 4, 2, 1....]

[[0, 0, 0, 1, 0, 0, 0, 0, 0, 0]
[0, 0, 0, 0, 1, 0, 0, 0, 0, 0]
[0, 0, 1, 0, 0, 0, 0, 0, 0, 0]
[0, 1,0 , 0, 0, 0, 0, 0, 0, 0]

....

The "dense" layer, which is fully connected to next dense layer or "softmax" layer needs the properly encoded labels
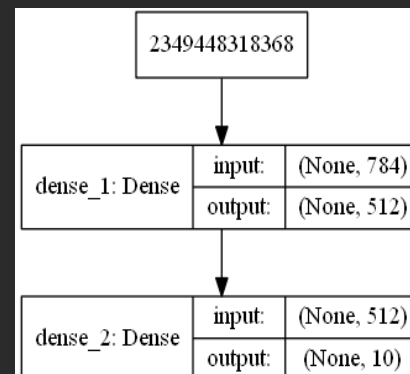
$$\frac{\exp(t_i)}{\sum_{j=1}^{K} \exp(t_j)} = \text{softmax}(t_i)$$
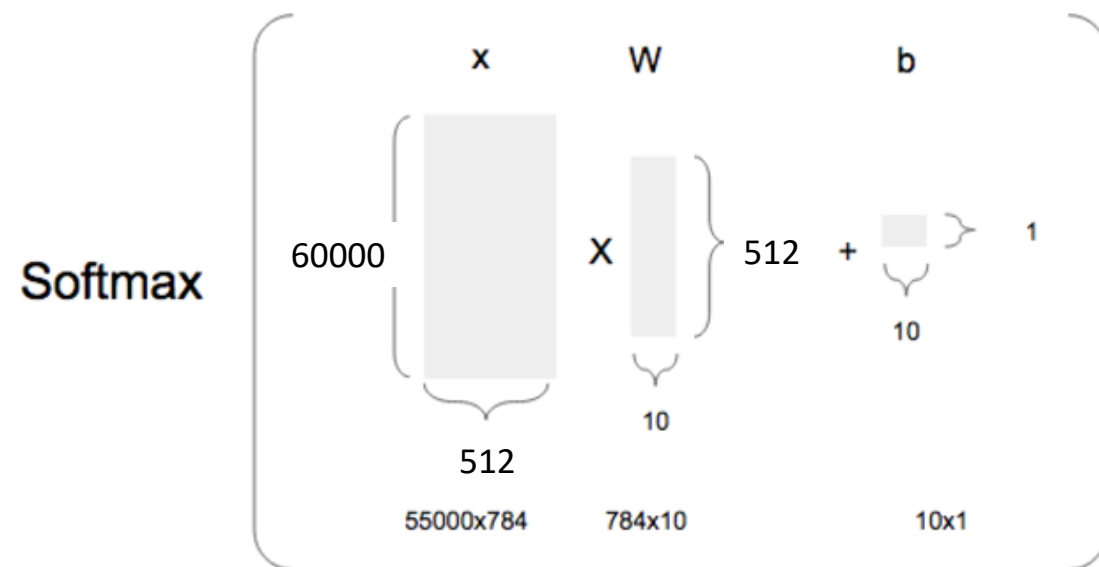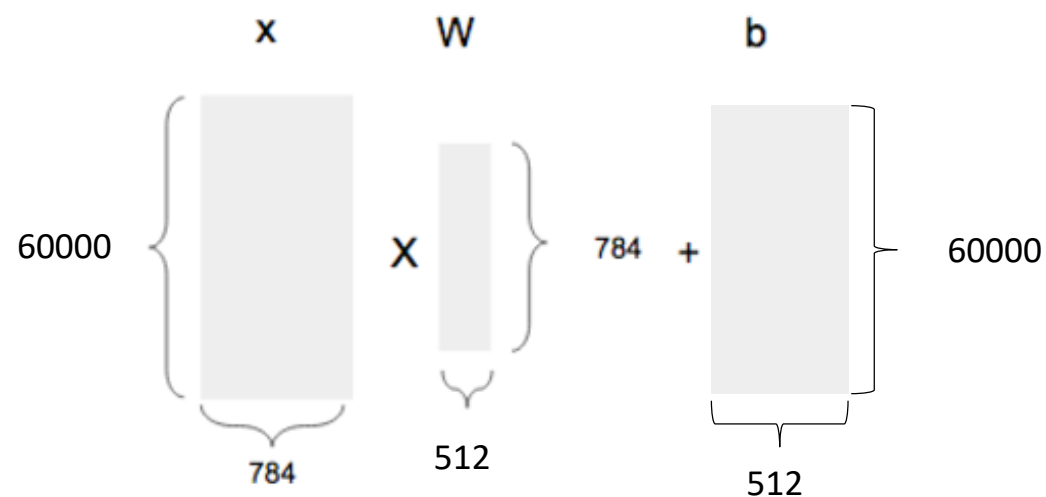
=1

# 2.5 Reviewing the first example

```python
network = models.Sequential()
network.add(layers.Dense(512, activation='relu', input_shape=(28 * 28,)))
network.add(layers.Dense(10, activation='softmax'))

network.summary()
plot_model(network, to_file='model.png',show_shapes=True,
show_layer_names=True)
```



A user just need to set the shape of output. In Keras, the weigh and bias is automatically calculated.
The problem of this example is the classification problem. Thus, the later including "Softmax" needs to match up with the encoded labels.

# 2.4 Neural Network Engine: Gradient-based Optimization
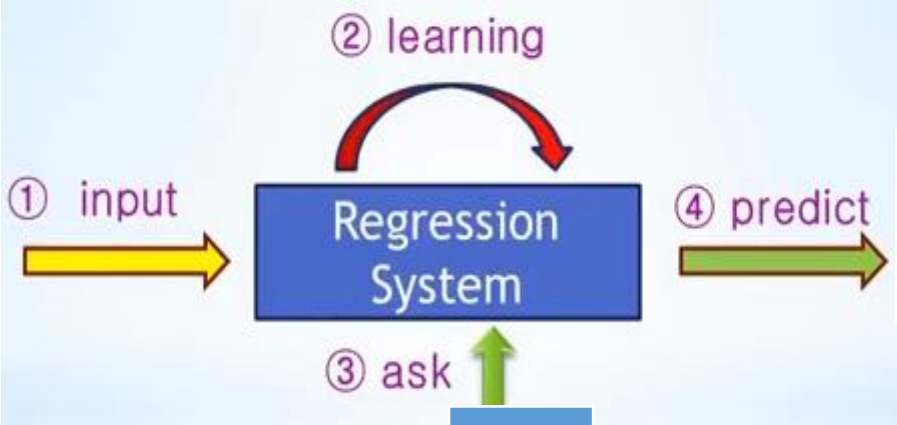
```python
network.compile(optimizer='rmsprop',
                loss='categorical_crossentropy',
                metrics=['accuracy'])
network.fit(train_images, train_labels, epochs=5, batch_size=128)

test_loss, test_acc = network.evaluate(test_images, test_labels)
print('test_acc:', test_acc)
```

For example, if the batch size is 128, you will see only 128 data samples, get a gradient, and update the weight and bias. You can update weights much faster than gradient decent using the whole data. This method of obtaining weight is called Stochastic Gradient Descent method.
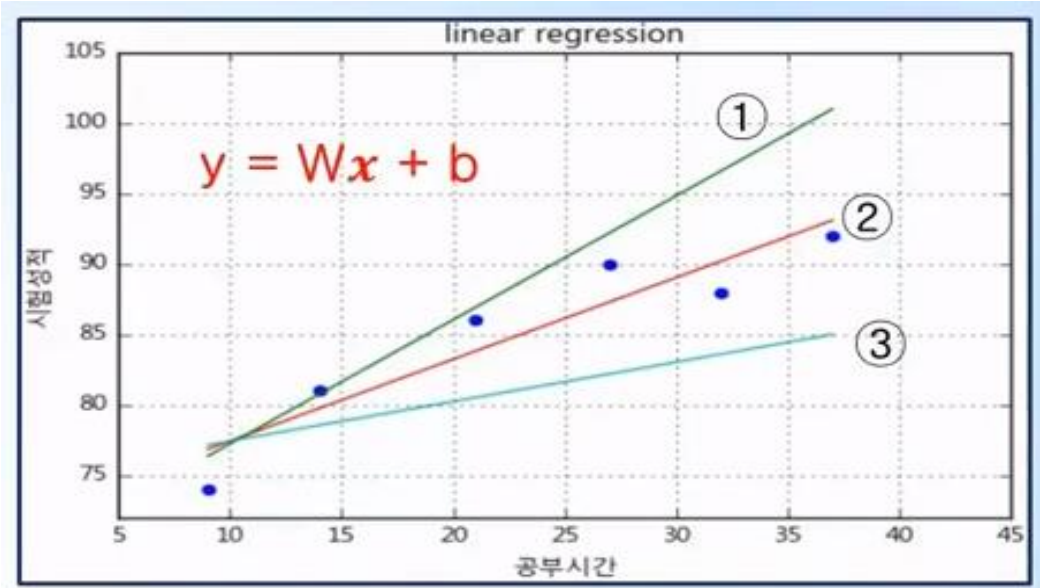
Gradient Descent for linear regression

| X | t |
|---|---|
| 9 | 74 |
| 14 | 81 |
| 21 | 86 |
| 27 | 90 |

② learning

① input

Regression System

④ predict

| X | t |
|---|---|
| 55 | ? |
| 60 | ? |

③ ask

| X |
|---|
| 55 |
| 60 |

$$y = Wx + b$$

linear regression

공부시간

# 2.4 Neural Network Engine: Gradient-based Optimization
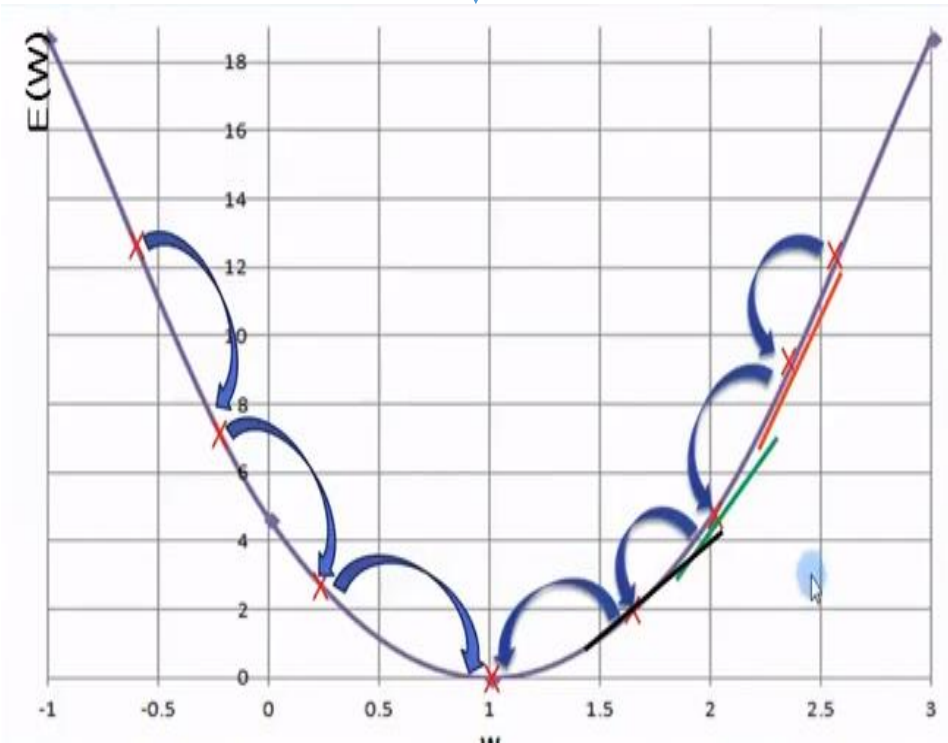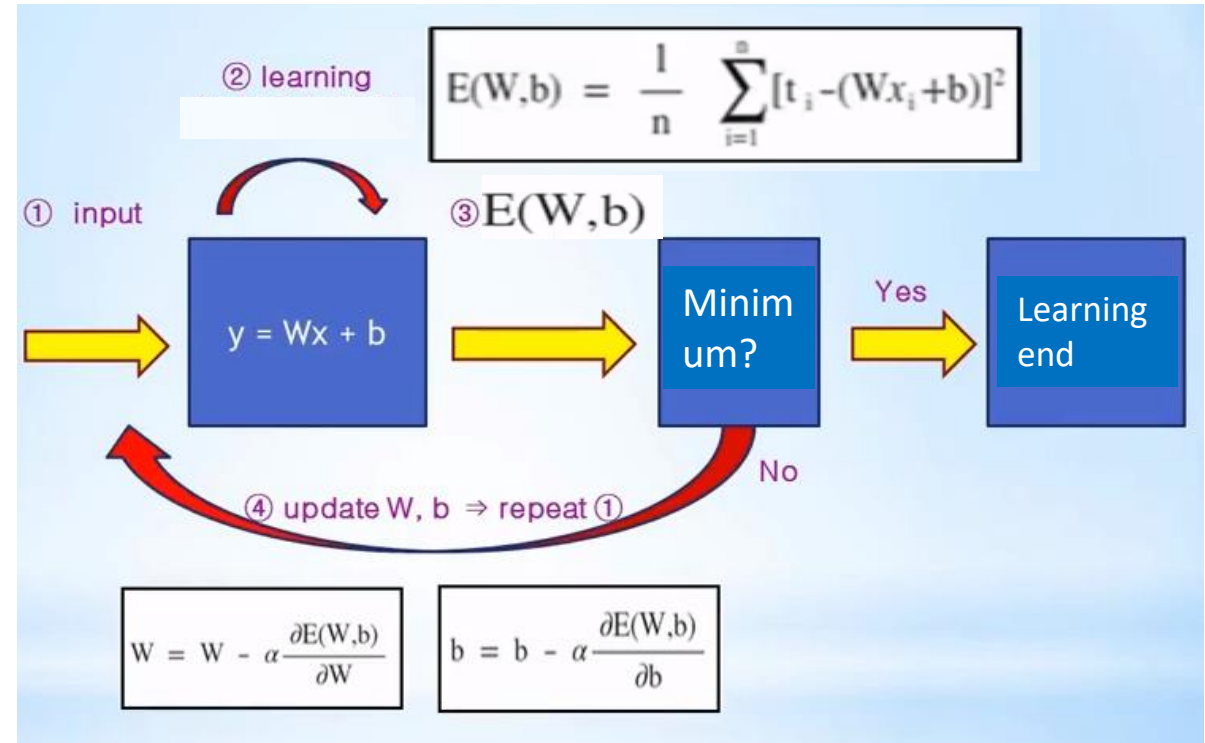
Gradient Descent for linear regression

$$y = Wx + b$$

$$\text{loss function} = E(W,b) = \frac{1}{n} \sum_{i=1}^{n} [t_i - y_i]^2 = \frac{1}{n} \sum_{i=1}^{n} [t_i - (Wx_i + b)]^2$$
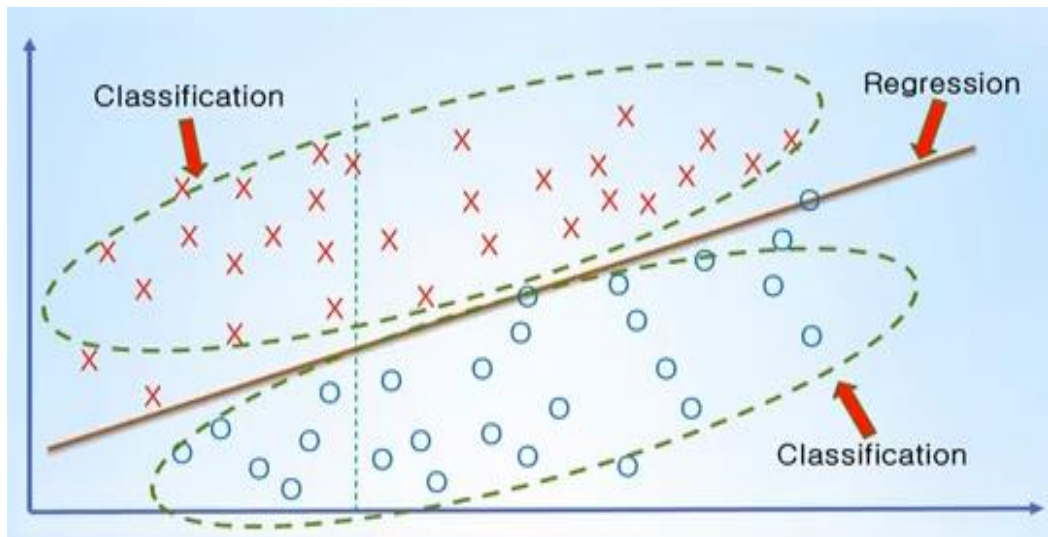
Linear regression is changed to convex function optimization problem

$$W = W - \alpha \frac{\partial E(W,b)}{\partial W} \qquad b = b - \alpha \frac{\partial E(W,b)}{\partial b}$$

Learning rate



② learning

$$E(W,b) = \frac{1}{n} \sum_{i=1}^{n} [t_i - (Wx_i + b)]^2$$

① input      ③ E(W,b)

y = Wx + b → Minimum? → (Yes) → Learning end

No

④ update W, b ⇒ repeat ①

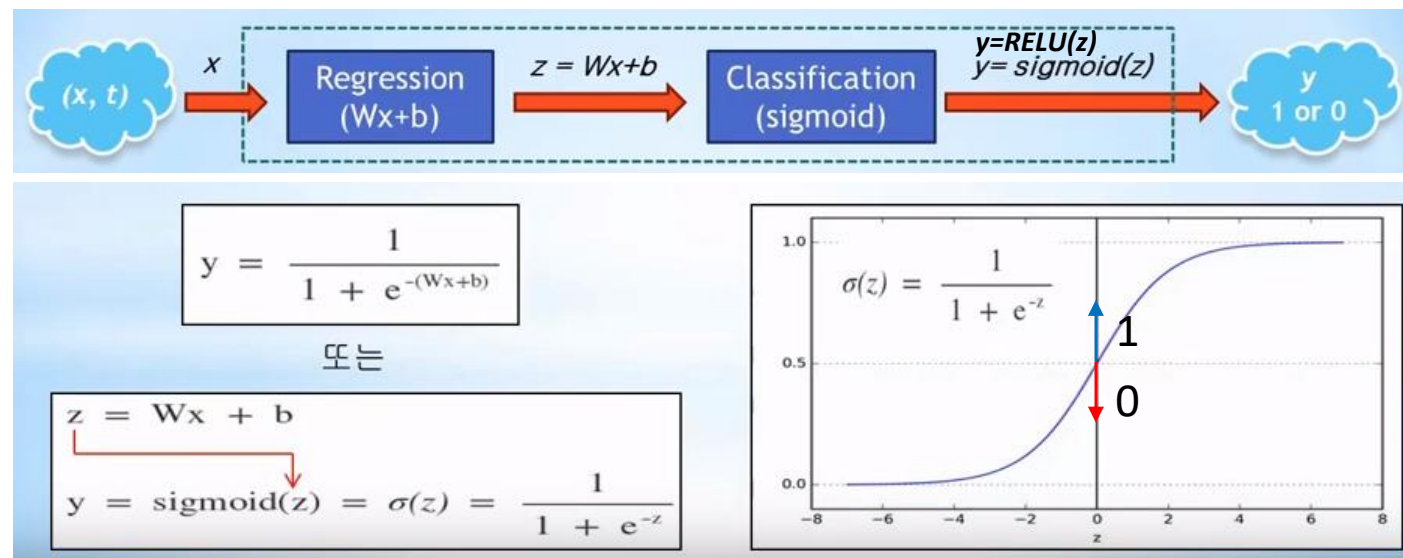$$W = W - \alpha \frac{\partial E(W,b)}{\partial W} \qquad b = b - \alpha \frac{\partial E(W,b)}{\partial b}$$

The concept of gradient descent for classification



1) Find the optimal fitted line representing the training data characteristics and distribution
2) Classify data as "up" or "down" based on its best fitted line



$$y = \frac{1}{1 + e^{-(Wx+b)}}$$

$$z = Wx + b$$

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

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

The final output (y) of classification problem is logically 1 or 0 by the sigmoid or RELU function. Thus, a loss function is needed that is different from linear regression with continuous values.

Cross-entropy

$$y = \frac{1}{1 + e^{-(Wx+b)}} \quad , \quad t_i = 0 \ or \ 1$$

The result of linear regression

$$E(W,b) = -\sum_{i=1}^{n} \{t_i \log y_i + (1-t_i)\log(1-y_i)\}$$
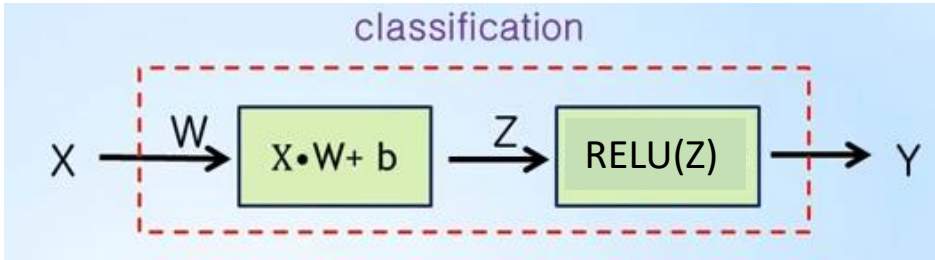
$$W = W - \alpha \frac{\partial E(W,b)}{\partial W} \qquad b = b - \alpha \frac{\partial E(W,b)}{\partial b}$$

# The concept of gradient descent for classification for multi variable case

```
network.compile(optimizer='rmsprop',
            loss='categorical_crossentropy',
            metrics=['accuracy'])
network.fit(train_images, train_labels, epochs=5, batch_size=128)
```

classification

$$X \xrightarrow{W} \boxed{X \bullet W + b} \xrightarrow{Z} \boxed{RELU(Z)} \rightarrow Y$$

| X1 | X2 | T |
|----|----|---|
| 2  | 4  | 0 |
| 4  | 11 | 0 |
| 6  | 6  | 0 |
| 8  | 5  | 0 |
| 10 | 7  | 1 |
| 12 | 16 | 1 |
| 14 | 8  | 1 |
| 16 | 3  | 1 |
| 18 | 7  | 1 |

$$\begin{pmatrix} 2 & 4 \\ 4 & 11 \\ 6 & 6 \\ 8 & 5 \\ 10 & 7 \\ 12 & 16 \\ 14 & 8 \\ 16 & 3 \\ 18 & 7 \end{pmatrix} \bullet W + b \implies \begin{pmatrix} z_1 \\ z_2 \\ z_3 \\ z_4 \\ z_5 \\ z_6 \\ z_7 \\ z_8 \\ z_9 \end{pmatrix} \xrightarrow[\text{calculation}]{\text{RELU}} \begin{pmatrix} y_1 \\ y_2 \\ y_3 \\ y_4 \\ y_5 \\ y_6 \\ y_7 \\ y_8 \\ y_9 \end{pmatrix} \xleftrightarrow[\text{son}]{\text{Compari}} \begin{pmatrix} 0 \\ 0 \\ 0 \\ 0 \\ 1 \\ 1 \\ 1 \\ 1 \\ 1 \end{pmatrix}$$

Dot product

Input  X

regression Output  Z

classification Result  Y

Ground truth  T

$$( 9 \times 2 ) \bullet ( 2 \times 1 ) = ( 9 \times 1 ) \qquad ( 9 \times 1 ) \qquad ( 9 \times 1 )$$

The concept of gradient descent for classification for multi variable case
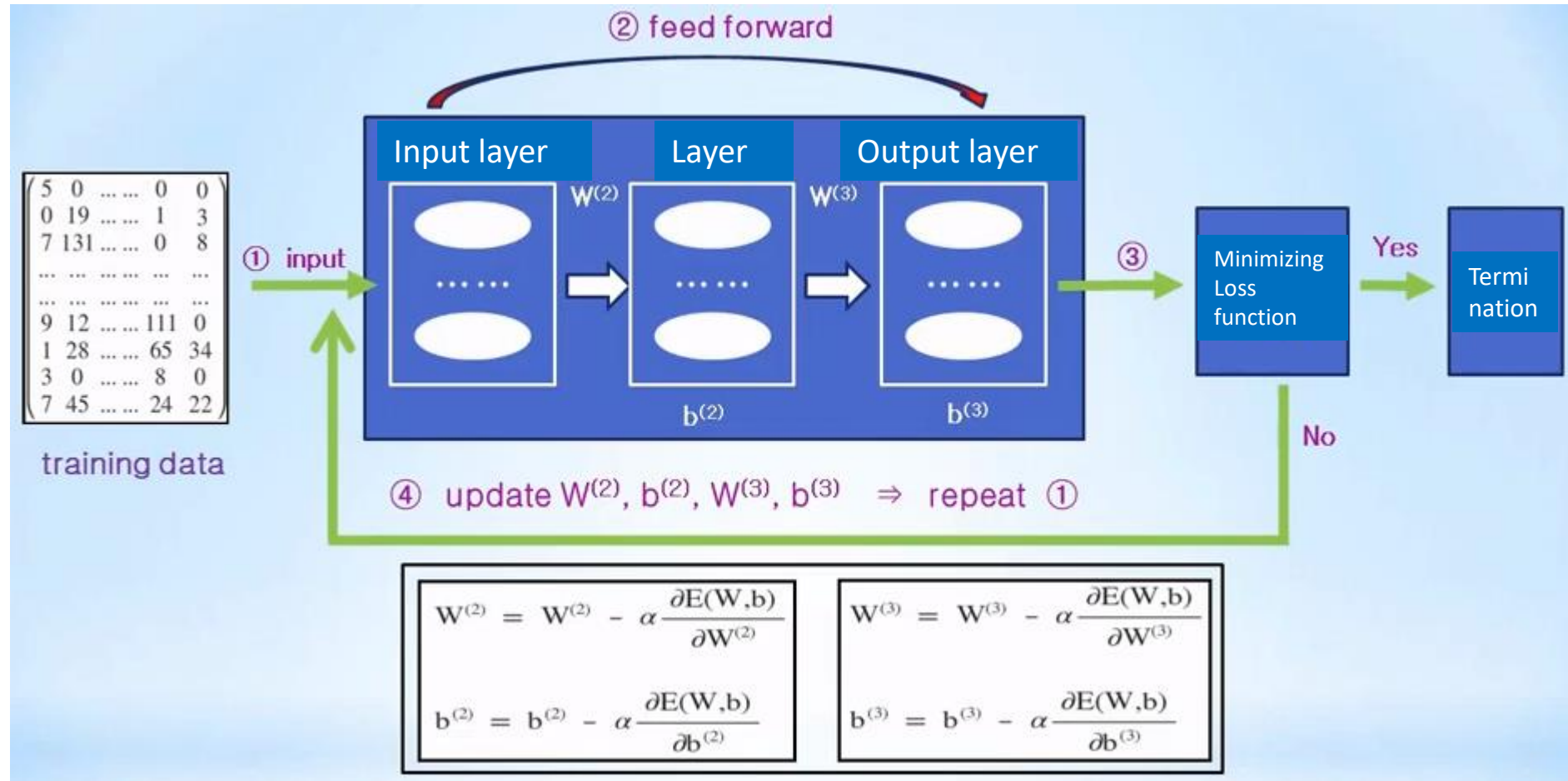
```
network.compile(optimizer='rmsprop',
                loss='categorical_crossentropy',
                metrics=['accuracy'])
network.fit(train_images, train_labels, epochs=5, batch_size=128)
```
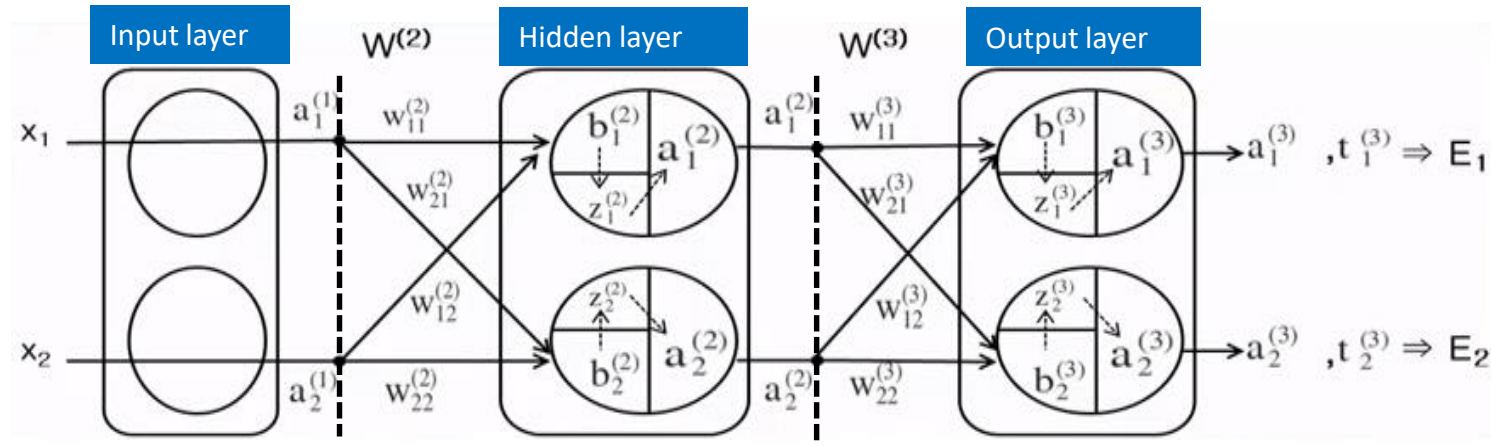
If the this code is successfully excuted..

1batch

| 0 | 0 | 0.01 |

About "6"

60000 (1epoch)

784

# Back propagation – conceptual access with sigmoid

Problem of optimization process only using calculus

② feed forward

Input layer | Layer | Output layer

$W^{(2)}$  $W^{(3)}$

① input

③ Minimizing Loss function

Yes

Termi nation

No

$\begin{pmatrix} 5 & 0 & \cdots\cdots & 0 & 0 \\ 0 & 19 & \cdots\cdots & 1 & 3 \\ 7 & 131 & \cdots\cdots & 0 & 8 \\ \cdots & \cdots & \cdots\cdots & \cdots & \cdots \\ \cdots & \cdots & \cdots\cdots & \cdots & \cdots \\ 9 & 12 & \cdots\cdots & 111 & 0 \\ 1 & 28 & \cdots\cdots & 65 & 34 \\ 3 & 0 & \cdots\cdots & 8 & 0 \\ 7 & 45 & \cdots\cdots & 24 & 22 \end{pmatrix}$

training data

$b^{(2)}$  $b^{(3)}$

④ update $W^{(2)}$, $b^{(2)}$, $W^{(3)}$, $b^{(3)}$  ⇒ repeat ①

$$W^{(2)} = W^{(2)} - \alpha \frac{\partial E(W,b)}{\partial W^{(2)}}$$

$$b^{(2)} = b^{(2)} - \alpha \frac{\partial E(W,b)}{\partial b^{(2)}}$$

$$W^{(3)} = W^{(3)} - \alpha \frac{\partial E(W,b)}{\partial W^{(3)}}$$

$$b^{(3)} = b^{(3)} - \alpha \frac{\partial E(W,b)}{\partial b^{(3)}}$$

Excessive numerical differential operation causes increased computation cost!!
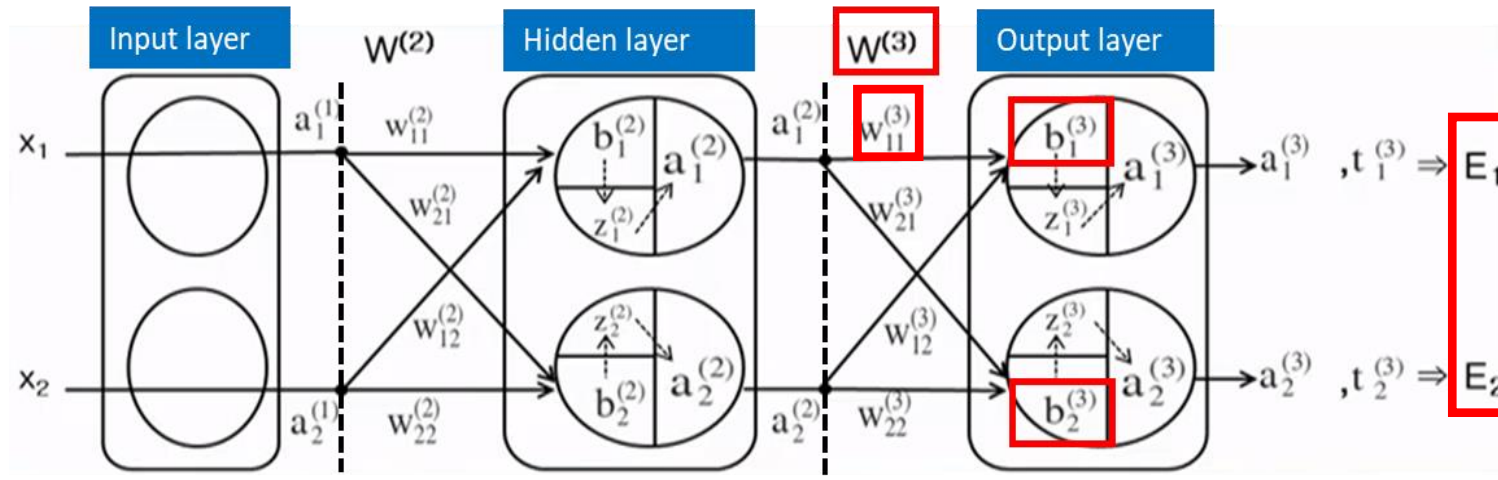
# Back propagation – conceptual access with sigmoid



Neural network diagram with Input layer, $W^{(2)}$, Hidden layer, $W^{(3)}$, Output layer.

Inputs: $x_1$, $x_2$

$a_1^{(1)}$, $w_{11}^{(2)}$, $w_{21}^{(2)}$, $w_{12}^{(2)}$, $w_{22}^{(2)}$, $a_2^{(1)}$

Hidden: $b_1^{(2)}$, $a_1^{(2)}$, $z_1^{(2)}$, $z_2^{(2)}$, $b_2^{(2)}$, $a_2^{(2)}$

$a_1^{(2)}$, $w_{11}^{(3)}$, $w_{21}^{(3)}$, $w_{12}^{(3)}$, $w_{22}^{(3)}$, $a_2^{(2)}$

Output: $b_1^{(3)}$, $a_1^{(3)}$, $z_1^{(3)}$, $z_2^{(3)}$, $a_2^{(3)}$, $b_2^{(3)}$

$\rightarrow a_1^{(3)}$ , $t_1^{(3)} \Rightarrow E_1$

$\rightarrow a_2^{(3)}$ , $t_2^{(3)} \Rightarrow E_2$

$$\frac{\partial \text{sigmoid}(z)}{\partial z} = \frac{\partial}{\partial z}\left(\frac{1}{1+e^{-z}}\right)$$

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

$$= \frac{e^{-z}}{(1+e^{-z})^2}$$

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

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

$$= \frac{1}{(1+e^{-z})} \times \left(1 - \frac{1}{(1+e^{-z})}\right)$$

$$= \text{sigmoid}(z) \times (1 - \text{sigmoid}(z))$$

|  | Linear regression value (z) | Output (a) |
|---|---|---|
| Input | 입력 층에는 가중치가 없기 때문에 선형회귀 값은 적용하지 않음 | $a_1^{(1)} = x_1$ $\quad$ $a_2^{(1)} = x_2$ |
| Hidden | $z_1^{(2)} = a_1^{(1)}w_{11}^{(2)} + a_2^{(1)}w_{12}^{(2)} + b_1^{(2)}$ | $a_1^{(2)} = \text{sigmoid}(z_1^{(2)})$ |
|  | $z_2^{(2)} = a_1^{(1)}w_{21}^{(2)} + a_2^{(1)}w_{22}^{(2)} + b_2^{(2)}$ | $a_2^{(2)} = \text{sigmoid}(z_2^{(2)})$ |
| output | $z_1^{(3)} = a_1^{(2)}w_{11}^{(3)} + a_2^{(2)}w_{12}^{(3)} + b_1^{(3)}$ | $a_1^{(3)} = \text{sigmoid}(z_1^{(3)})$ |
|  | $z_2^{(3)} = a_1^{(2)}w_{21}^{(3)} + a_2^{(2)}w_{22}^{(3)} + b_2^{(3)}$ | $a_2^{(3)} = \text{sigmoid}(z_2^{(3)})$ |

# Back propagation – conceptual access with sigmoid : example (1)
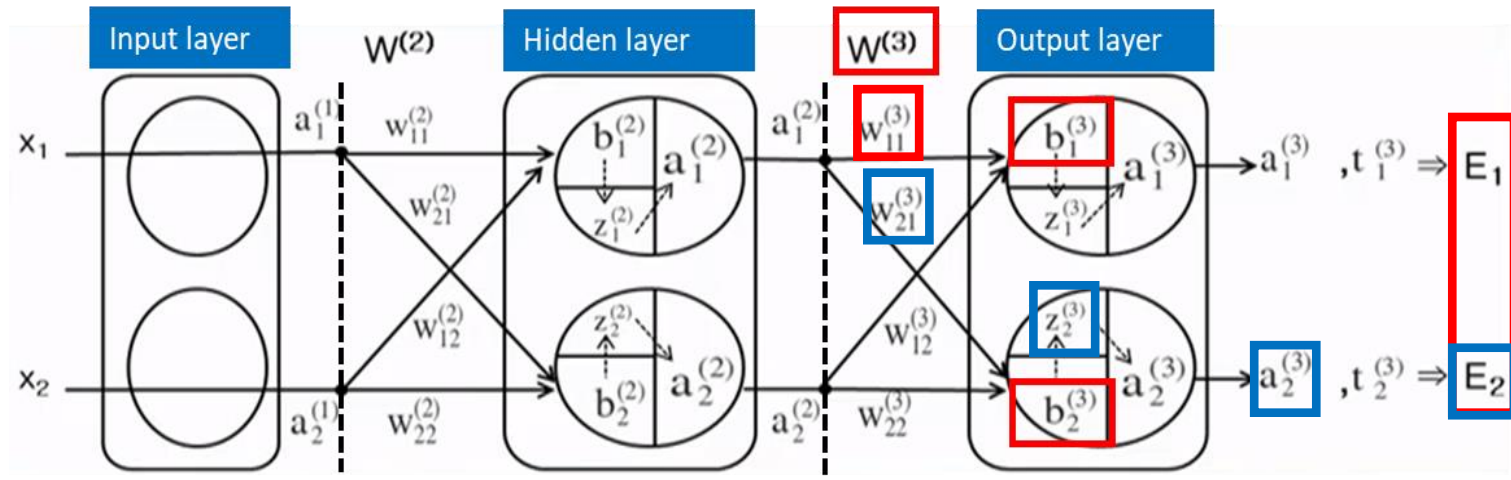


$$\frac{\partial E}{\partial w_{11}^{(3)}} = \frac{\partial E_1}{\partial w_{11}^{(3)}} + \frac{\partial E_2}{\partial w_{11}^{(3)}} \nearrow^{0}$$

$$= \frac{\partial E_1}{\partial a_1^{(3)}} \times \frac{\partial a_1^{(3)}}{\partial z_1^{(3)}} \times \frac{\partial z_1^{(3)}}{\partial w_{11}^{(3)}}$$

$$= \frac{\partial \left\{ \frac{1}{2}(t_1^{(3)} - a_1^{(3)})^2 \right\}}{\partial a_1^{(3)}} \times \frac{\partial \, sigmoid(z_1^{(3)})}{\partial z_1^{(3)}} \times \frac{\partial (a_1^{(2)} w_{11}^{(3)} + a_2^{(2)} w_{12}^{(3)} + b_1^{(3)})}{\partial w_{11}^{(3)}}$$

$$= (a_1^{(3)} - t_1^{(3)}) \times sigmoid(z_1^{(3)}) \times (1 - sigmoid(z_1^{(3)})) \times a_1^{(2)}$$

$$= (a_1^{(3)} - t_1^{(3)}) \times a_1^{(3)} \times (1 - a_1^{(3)}) \times a_1^{(2)}$$

where

$$E = E_1 + E_2$$

$$E_1 = \frac{1}{2}(t_1^{(3)} - a_1^{(3)})^2$$

$$a_1^{(3)} = sigmoid(z_1^{(3)})$$

$$z_1^{(3)} = a_1^{(2)} w_{11}^{(3)} + a_2^{(2)} w_{12}^{(3)} + b_1^{(3)}$$

Calculus becomes algebraic equation!

# Back propagation – conceptual access with sigmoid : example (2)



$$\frac{\partial E}{\partial w_{21}^{(3)}} = \frac{\partial E_1}{\partial w_{21}^{(3)}}^{0} + \frac{\partial E_2}{\partial w_{21}^{(3)}}$$

$$= \frac{\partial E_2}{\partial a_2^{(3)}} \times \frac{\partial a_2^{(3)}}{\partial z_2^{(3)}} \times \frac{\partial z_2^{(3)}}{\partial w_{21}^{(3)}}$$

$$= \frac{\partial\left\{\frac{1}{2}(t_2^{(3)}-a_2^{(3)})^2\right\}}{\partial a_2^{(3)}} \times \frac{\partial \text{sigmoid}(z_2^{(3)})}{\partial z_2^{(3)}} \times \frac{\partial(a_1^{(2)}w_{21}^{(3)}+a_2^{(2)}w_{22}^{(3)}+b_2^{(3)})}{\partial w_{21}^{(3)}}$$

$$= (a_2^{(3)} - t_2^{(3)}) \times \text{sigmoid}(z_2^{(3)}) \times (1-\text{sigmoid}(z_2^{(3)})) \times a_1^{(2)}$$

$$= (a_2^{(3)} - t_2^{(3)}) \times a_2^{(3)} \times (1 - a_2^{(3)}) \times a_1^{(2)}$$

where

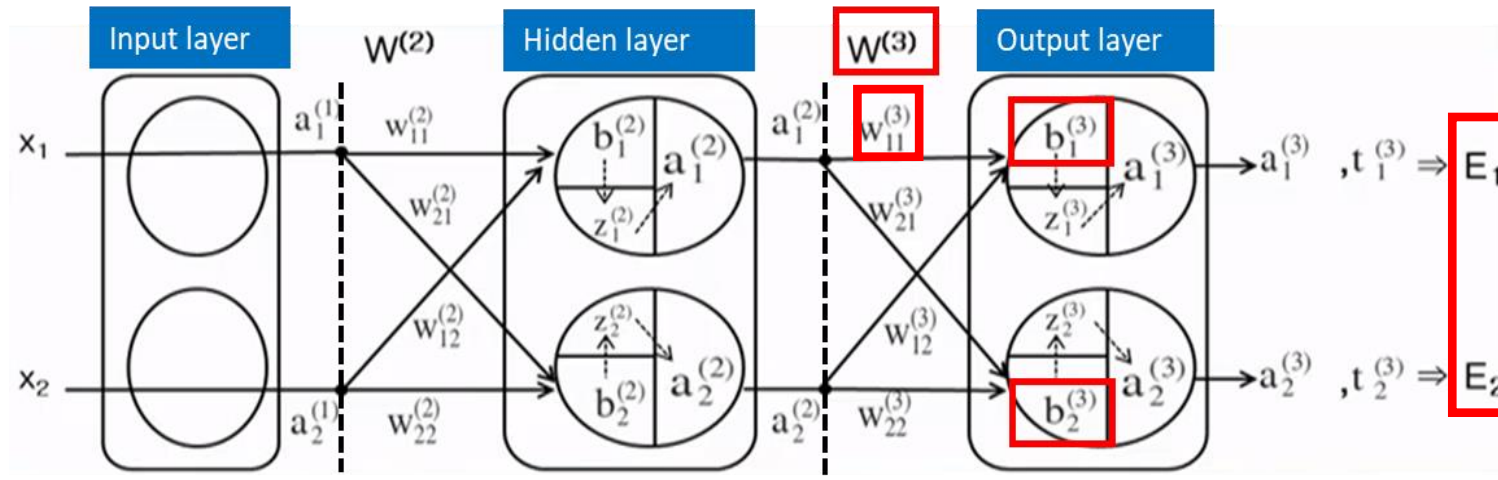$$E = E_1 + E_2$$

$$E_1 = \frac{1}{2}(t_1^{(3)}-a_1^{(3)})^2$$

$$a_1^{(3)} = \text{sigmoid}(z_1^{(3)})$$

$$z_1^{(3)} = a_1^{(2)}w_{11}^{(3)} + a_2^{(2)}w_{12}^{(3)} + b_1^{(3)}$$

Similarly, the others of weight in output layer (3) can be computed such as the back propagation of $w_{11}$ and $w_{12}$ in output layer.

# Back propagation – conceptual access with sigmoid : example (3)



$$\frac{\partial E}{\partial b_1^{(3)}} = \frac{\partial E_1}{\partial b_1^{(3)}} + \frac{\partial E_2}{\partial b_1^{(3)}}{}^{0}$$

$$= \frac{\partial E_1}{\partial a_1^{(3)}} \times \frac{\partial a_1^{(3)}}{\partial z_1^{(3)}} \times \frac{\partial z_1^{(3)}}{\partial b_1^{(3)}}$$

$$= \frac{\partial\left\{\frac{1}{2}(t_1^{(3)}-a_1^{(3)})^2\right\}}{\partial a_1^{(3)}} \times \frac{\partial \text{sigmoid}(z_1^{(3)})}{\partial z_1^{(3)}} \times \frac{\partial(a_1^{(2)}w_{11}^{(3)}+a_2^{(2)}w_{12}^{(3)}+b_1^{(3)})}{\partial b_1^{(3)}}$$

$$= (a_1^{(3)} - t_1^{(3)}) \times \text{sigmoid}(z_1^{(3)}) \times (1-\text{sigmoid}(z_1^{(3)})) \times 1$$

$$= (a_1^{(3)} - t_1^{(3)}) \times a_1^{(3)} \times (1 - a_1^{(3)}) \times 1$$

where

$$E = E_1 + E_2$$

$$E_1 = \frac{1}{2}(t_1^{(3)}-a_1^{(3)})^2$$

$$a_1^{(3)} = \text{sigmoid}(z_1^{(3)})$$

$$z_1^{(3)} = a_1^{(2)}w_{11}^{(3)} + a_2^{(2)}w_{12}^{(3)} + b_1^{(3)}$$

Let's think about the back propagation between hidden layer to E