

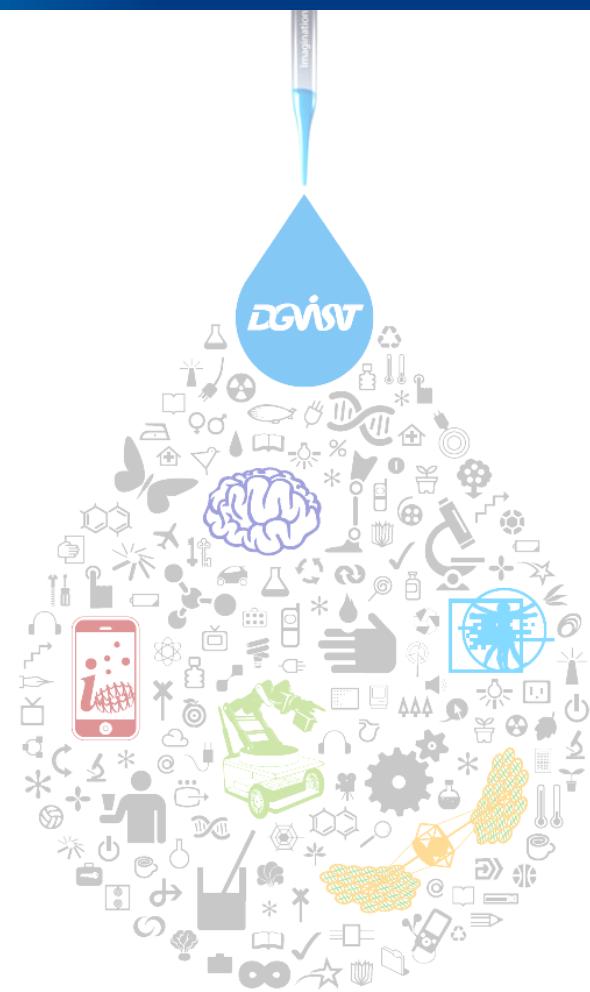


# Deep Learning Seminar

## CH 6. Deep Feedforward Networks

2017-07-26

Eunjeong Yi



# Chapter 6. Deep Feedforward Networks

**6.1 Example: Learning XOR**

**6.2 Gradient-Based Learning**

**6.3 Hidden Units**

**6.4 Architecture Design**

**6.5 Back-Propagation and Other Differentiation**

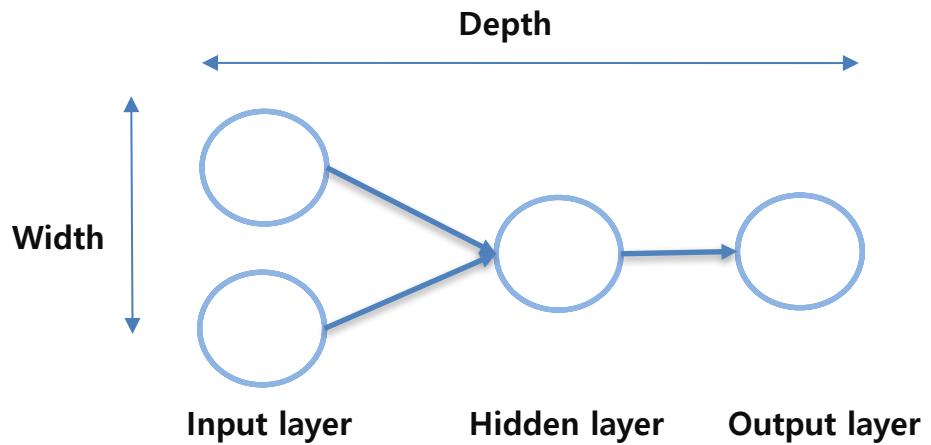
**6.6 Historical Notes**

# Deep feedforward network

- **No feedback connection**

- **Structure of model**

- Input layer, hidden layer, output layer
- The depth, width of model
- Cost function



# Example: Learning XOR

## ■ XOR function (“exclusive or”)

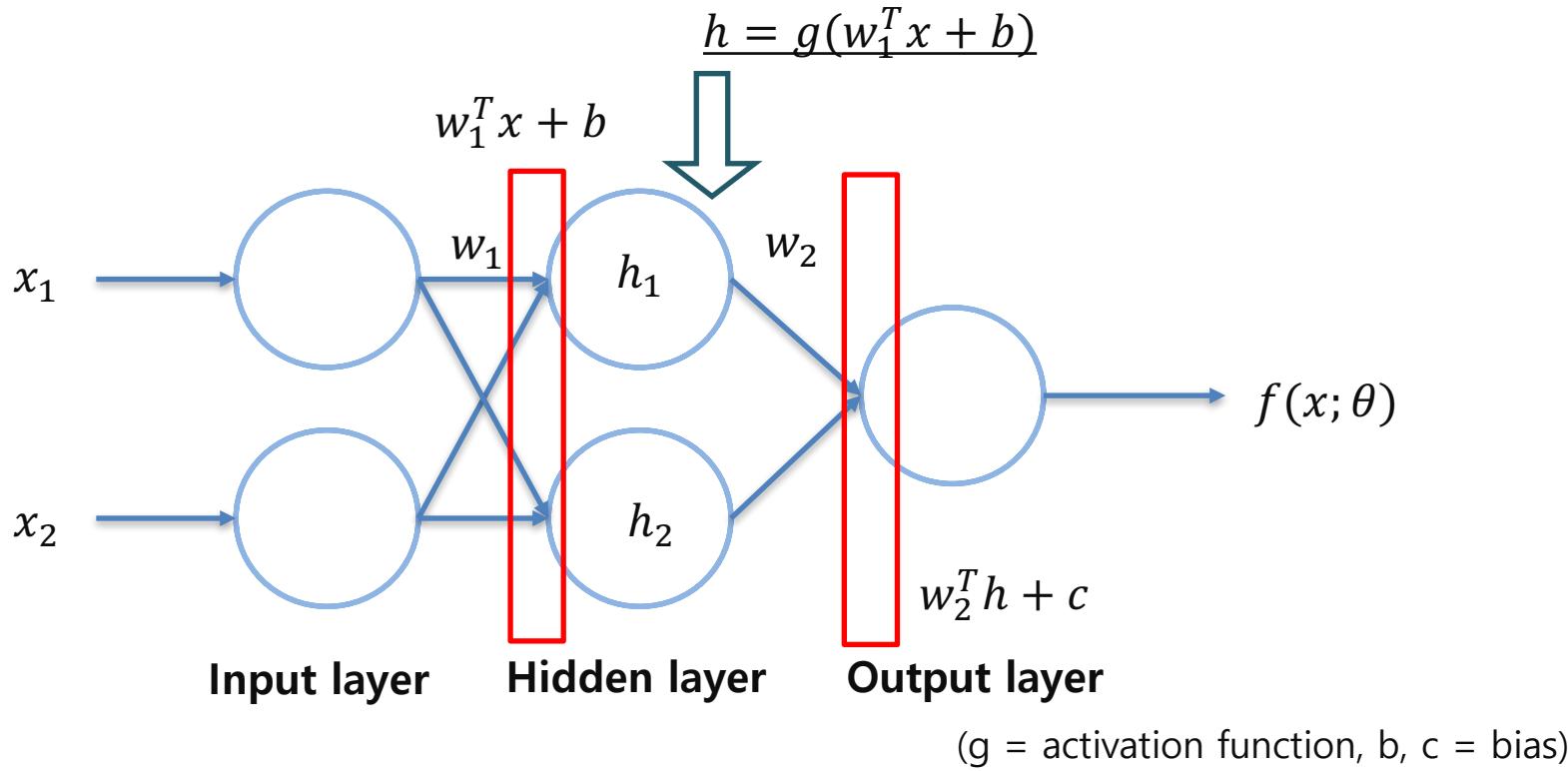
$x_1$	$x_2$	$x_1 \text{ } XOR \text{ } x_2$
0	0	0
0	1	1
1	0	1
1	1	0

- $\mathbb{X} = \{[0, 0]^T, [0, 1]^T, [1, 0]^T, [1, 1]^T\}$
- Mean squared error (MSE) loss function

$$J(\boldsymbol{\theta}) = \frac{1}{4} \sum_{\mathbf{x} \in \mathbb{X}} (f^*(\mathbf{x}) - f(\mathbf{x}; \boldsymbol{\theta}))^2$$

- $f^*(\mathbf{x})$ : correct answer
- $f(\mathbf{x}; \boldsymbol{\theta})$ : calculated result by neural network

# Example: Learning XOR



$$f(x; \theta) = f(x; w_1, b, w_2, c) = w_2^T g(w_1^T x + b) + c$$

# Gradient-based Learning

## ■ Loss function of neural network is non-convex

- Trained by using iterative, gradient-based optimizers

## ■ Cost function

- Learning Conditional Distribution
- Learning Conditional Statistics

## ■ Output layers

- Linear Units for Gaussian Output Distributions
- Sigmoid Units
- Softmax Units

# Learning Conditional Distribution

## ■ Negative log-likelihood as cross-entropy between training data and model distribution

$$J(\theta) = H(\hat{p}_{data}, p_{model}(y|x)) = - \sum_{x,y \sim \hat{p}_{data}} \hat{p}_{data} \log p_{model}(y|x)$$

$$J(\theta) = -\mathbb{E}_{x,y \sim \hat{p}_{data}} \log p_{model}(y|x)$$

- $\hat{p}_{data}$  : *training data – generating distribution*
- $p_{model}(y|x)$  : *probability distribution estimating  $\hat{p}_{data}$*
- $H(\hat{p}_{data}, p_{model}(y|x))$  : *cross entropy between  $\hat{p}_{data}$  and  $p_{model}$*

## ■ **log** function undoes the **exp** of some output units

# Learning Conditional Statistics

## ■ Mean Square Error (MSE)

$$f^* = \arg \min_f \mathbb{E}_{x,y \sim p_{data}} \|y - f(x)\|^2$$

## ■ Mean Absolute Error (MAE)

$$f^* = \arg \min_f \mathbb{E}_{x,y \sim p_{data}} \|y - f(x)\|_1$$

## ■ MSE, MAE often lead to poor results when used with gradient-based learning

# Linear Units for Gaussian Output Distributions

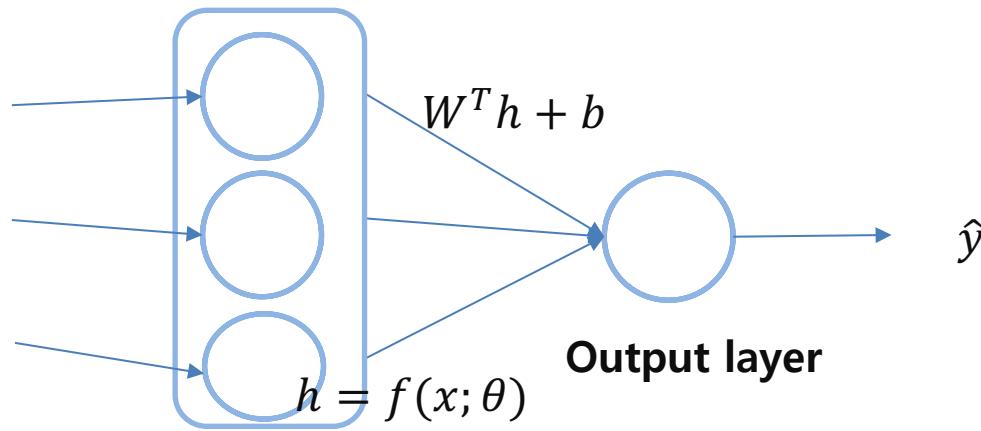
- Feature  $h$ , produce a vector  $\hat{y}$

$$\hat{y} = W^T h + b$$

- To produce the mean of a conditional Gaussian distribution

$$p(y|x) = \mathcal{N}(y; \hat{y}, I)$$

- No saturation → little difficult to gradient-based optimization

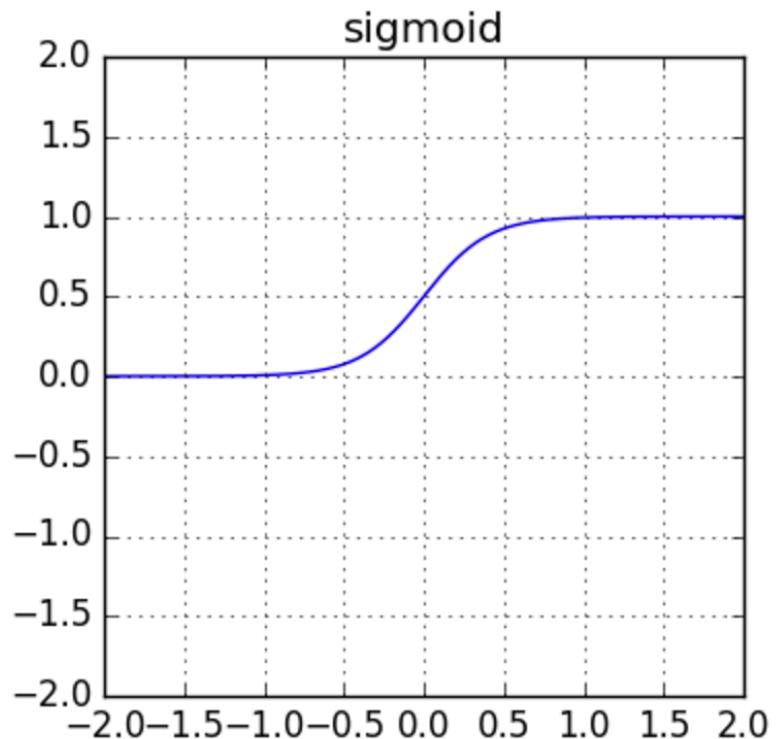


# Sigmoid Units

- **Binary classification**
- **Output:**  $\hat{y} = \sigma(w^T h + b)$

➤  $\sigma(z) = \frac{1}{1+e^{-z}}$

- **Saturate to 0 and 1**



© Copyright 2015-2017, CodeReclaimers

<http://neat-python.readthedocs.io/en/latest/activation.html>

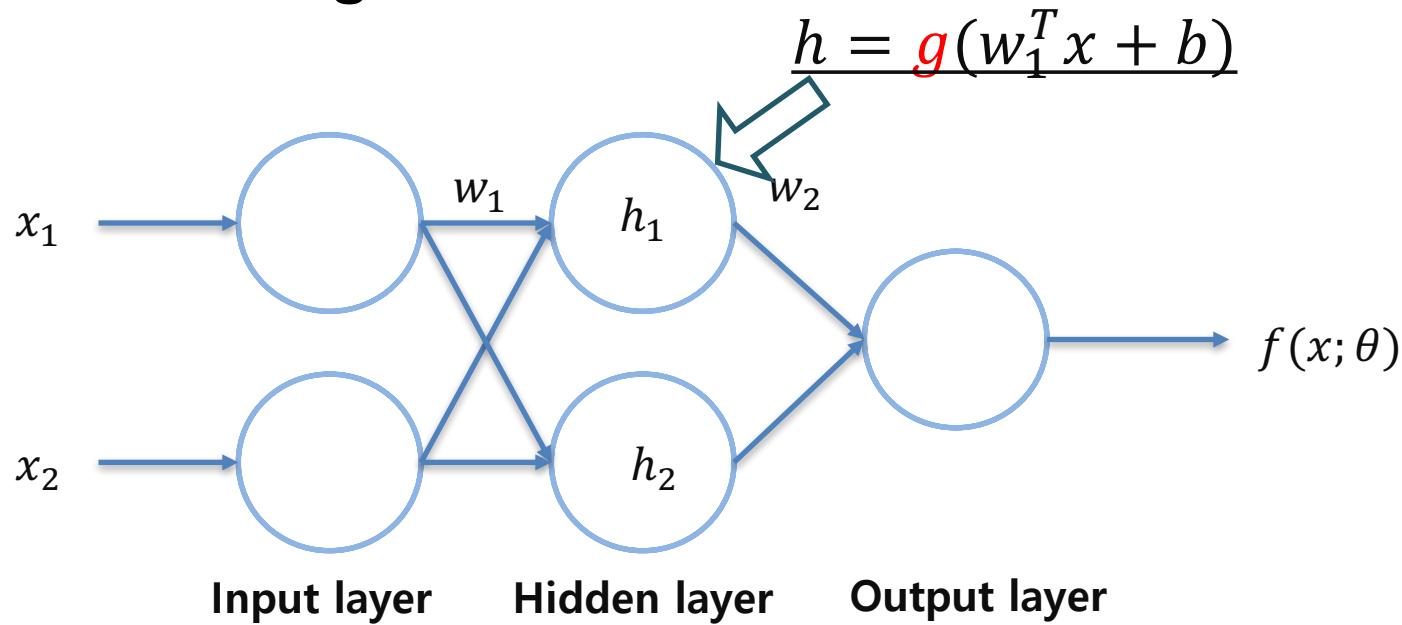
# Softmax Units

- Multiclass classification
- To generalize to the case of a discrete variable with n values, vector  $\hat{y}$ , with  $\hat{y}_i = P(y = i|x)$

$$\mathbf{z} = \mathbf{W}^T \mathbf{h} + \mathbf{b} \quad (z_i = \log \tilde{P}(y = i|x))$$
$$softmax(\mathbf{z})_i = \frac{exp(z_i)}{\sum_j exp(z_j)}$$

# Hidden Units

- How to choose the type of hidden unit to use in the hidden layers of the model
- Input:  $z = W^T x + b$
- Activation function  $g(z)$
- Ex) Rectified Linear Units (ReLU), Logistic Sigmoid and Hyperbolic Tangent

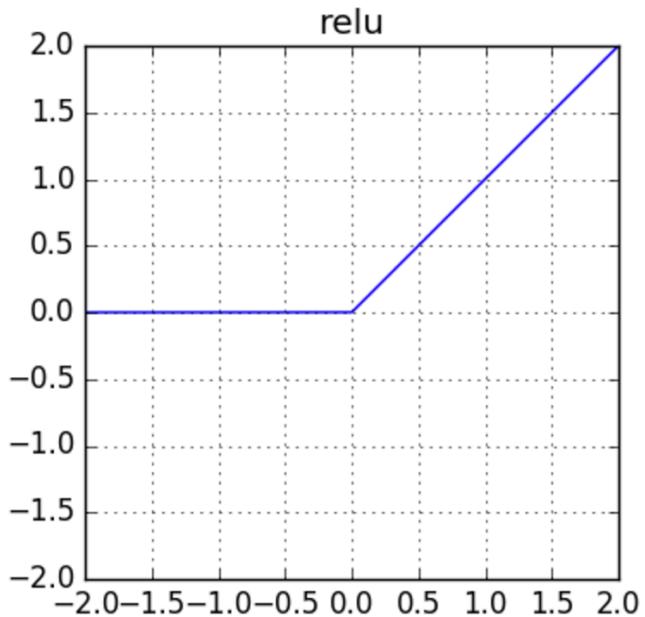


# Rectified Linear Units (ReLU)

- Activation function  $g(z)$

$$g(z) = \max(0, z)$$

- If model's behavior is closer to linear, models are easier to optimize



# Logistic Sigmoid and Hyperbolic Tangent

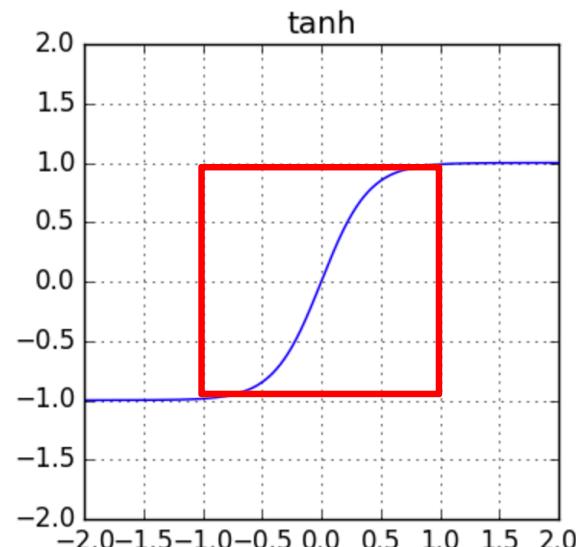
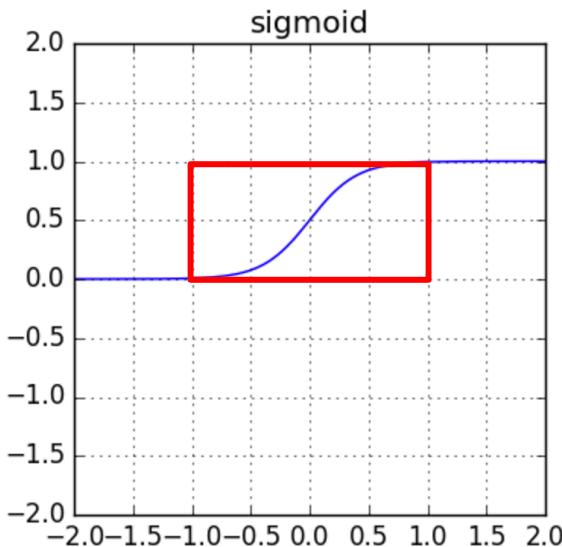
## ■ Activation function

$$g(z) = \sigma(z)$$

or

$$g(z) = \tanh(z) = 2\sigma(2z) - 1$$

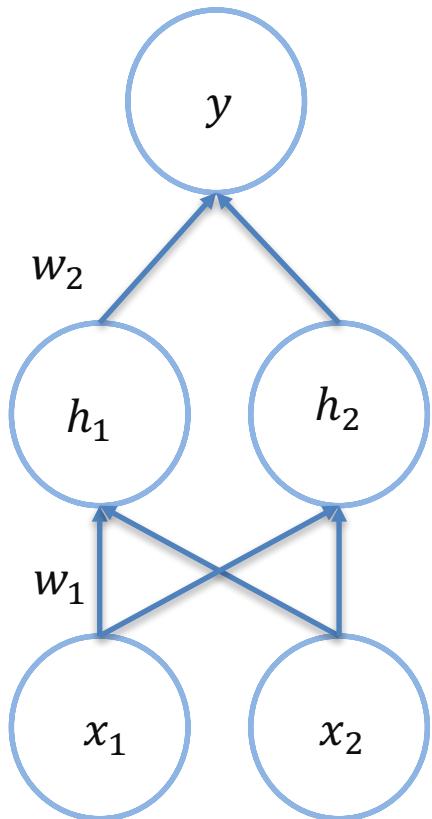
## ■ Saturation of sigmoidal units make gradient-based learning difficult



© Copyright 2015-2017, CodeReclaimers  
<http://neat-python.readthedocs.io/en/latest/activation.html>

# Back-Propagation

- Method to calculate the gradient of the loss function with respect to the weights in an artificial neural network



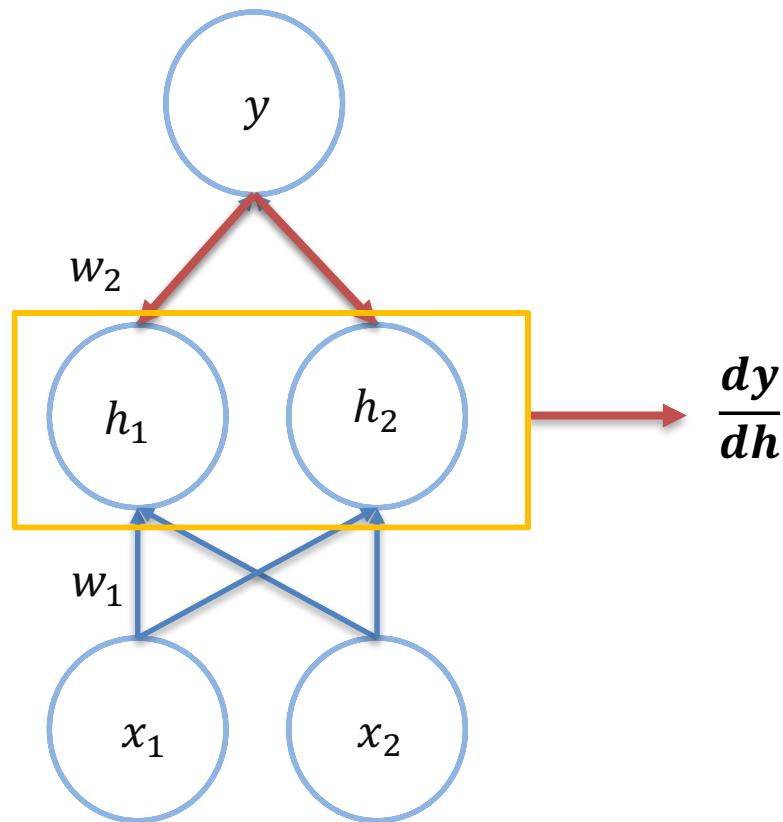
Forward propagation result

$$y = w_2^T h + b = w_2^T g(w_1^T x + c) + b$$

( $g$  = activation function,  $b$ ,  $c$  = bias)

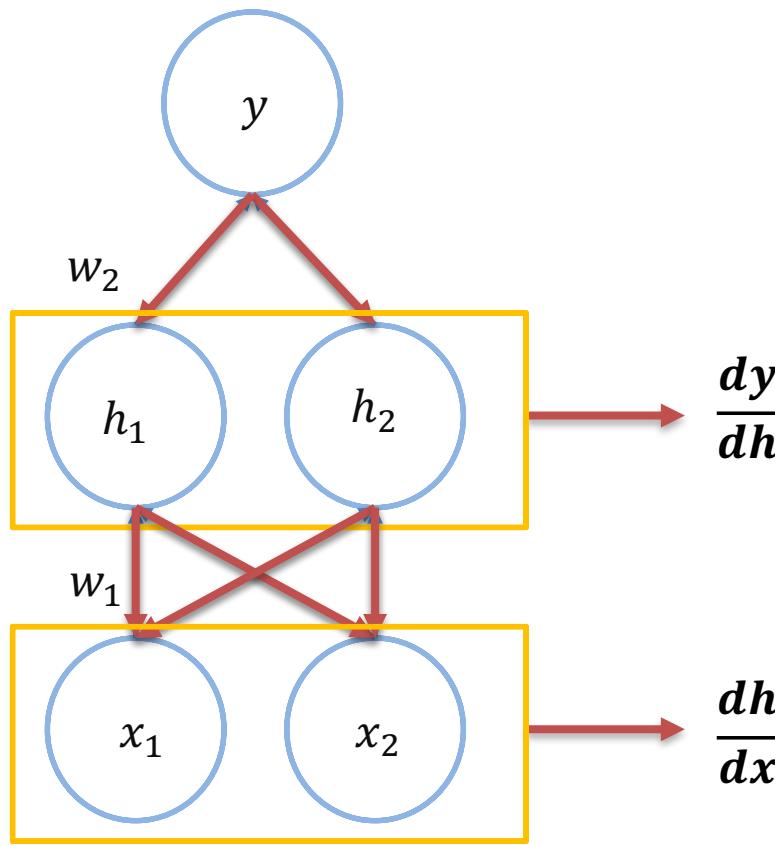
# Back-Propagation

- Method to calculate the gradient of the loss function with respect to the weights in an artificial neural network



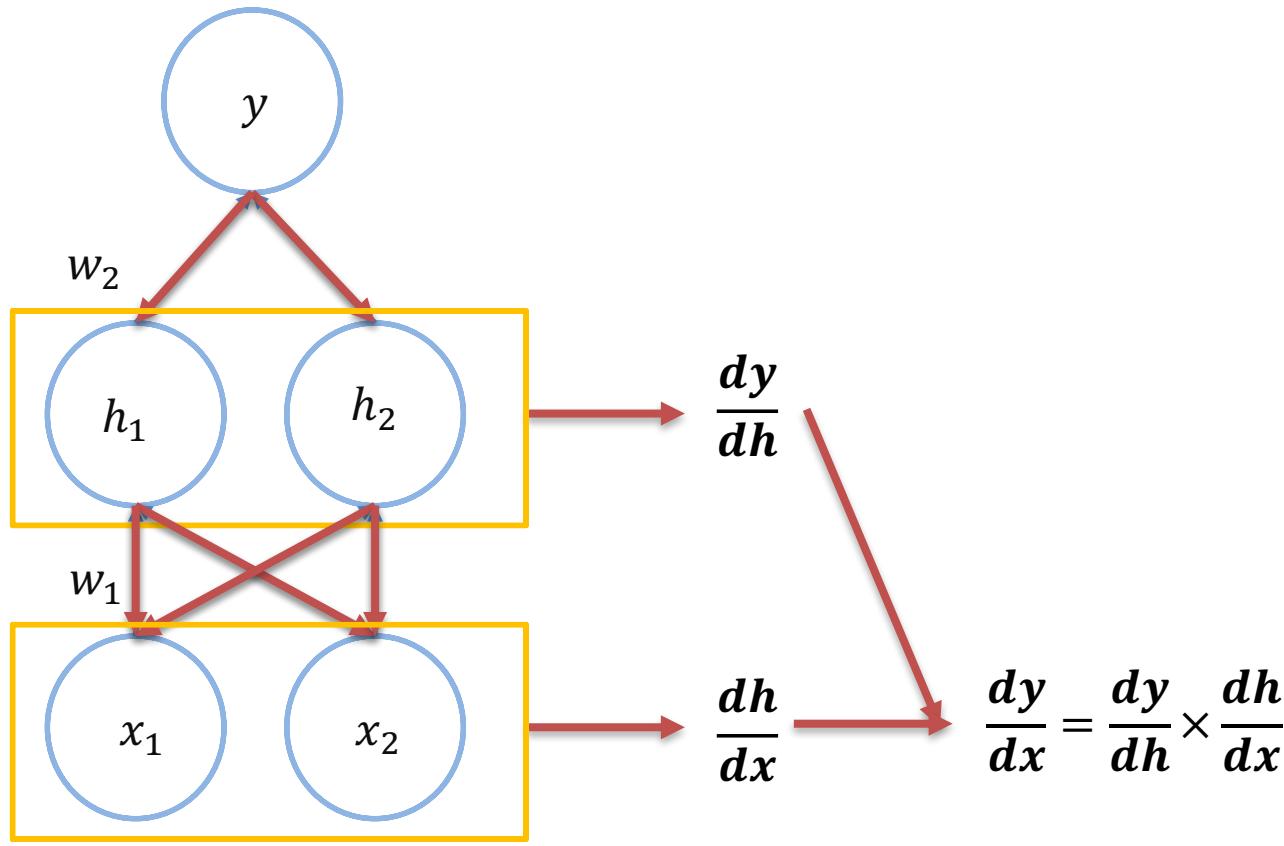
# Back-Propagation

- Method to calculate the gradient of the loss function with respect to the weights in an artificial neural network

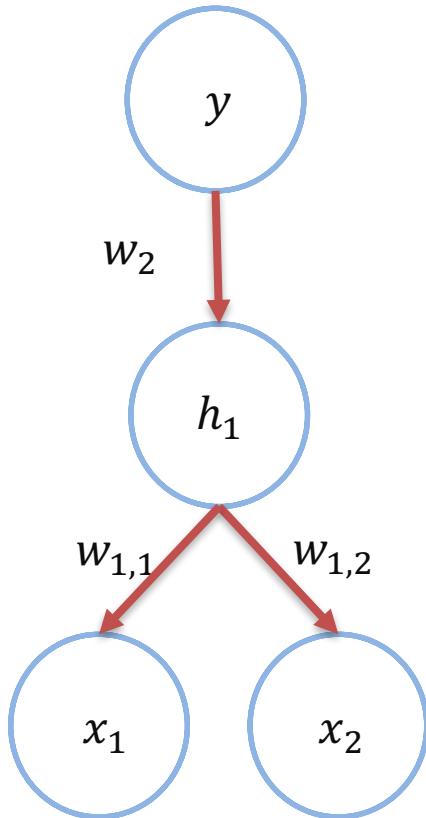


# Back-Propagation

- Method to calculate the gradient of the loss function with respect to the weights in an artificial neural network



# Back-Propagation



$$\frac{\delta y}{\delta w_{1,1}} = \frac{\delta y}{\delta x_1} \times \frac{\delta x_1}{\delta h_1} \times \frac{\delta h_1}{\delta w_{1,1}}$$

Update  $w_{1,1}$

$$w_{1,1} = w_{1,1} - \eta \frac{\delta y}{\delta w_{1,1}}$$

( $\eta$ : learning weight)

Using same way,  
all weight update

# Next Deep learning seminar

## Chapter 7. Regularization for Deep Learning

**7.1 Parameter Norm Penalties**

**7.2 Norm Penalties as Constrained Optimizaiton**

**7.3 Regularization and Under-Constrained Problems**

**7.4 Dataset Augmentation**

**7.5 Noise Robustness**

**7.6 Semi-Supervised Learning**

**7.7 Multitask Learning**



## Q&A

Thank you

