

# Deep Learning Programming

Lecture 3.1: Neural network classifier

Sangryul Jeon

School of Computer Science and Engineering

[srjeonn@pusan.ac.kr](mailto:srjeonn@pusan.ac.kr)

# 1.1 Linear Model

Linear Model

이미지  $\mathbf{x}$



$$f(\mathbf{x}, \mathbf{W}) = \mathbf{W}\mathbf{x} + \mathbf{b}$$

가중치 혹은 파라미터  $\mathbf{W}$



편향

Linear Model

# 1.1 Linear Model

Linear Model

이미지  $\mathbf{x}$

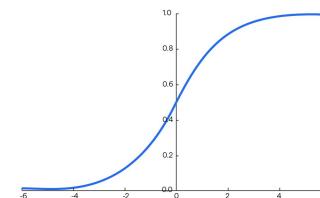


$$f(\mathbf{x}, \mathbf{W}) = \mathbf{W}\mathbf{x} + \mathbf{b}$$

Linear Model

Softmax Function

$$p(y = c_i | x) = \frac{e^{s_i}}{\sum_j e^{s_j}}$$



Softmax Loss



Softmax 함수 식

# 1.1 Linear Model

Linear Model

이미지  $\mathbf{x}$

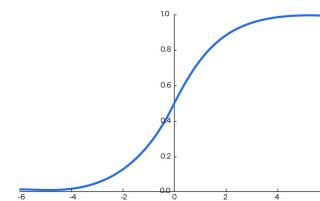


$$f(\mathbf{x}, \mathbf{W}) = \mathbf{W}\mathbf{x} + \mathbf{b}$$

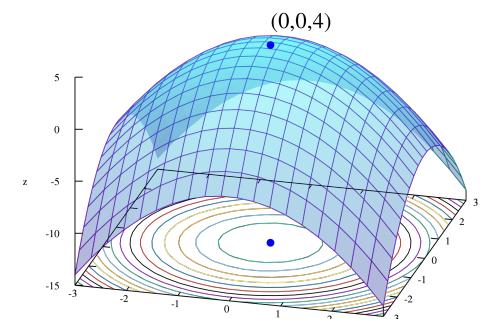
Linear Model

$$p(y = c_i | x) = \frac{e^{s_i}}{\sum_j e^{s_j}}$$

Softmax Function



Softmax Loss



[https://en.wikipedia.org/wiki/File:Max\\_paraboloid.svg](https://en.wikipedia.org/wiki/File:Max_paraboloid.svg)

매개 변수들은 (Stochastic) Gradient Descent에 의해 조정됩니다

## 1.2 Issues with Linear Classifiers

Linear Model

Linear Classifier은 강력하지 않습니다. 왜냐하면

- 시각적으로) 각 클래스 당 하나의 template(카테고리)만 학습할 수 있습니다.

$$f(\mathbf{x}) = \mathbf{Wx}$$



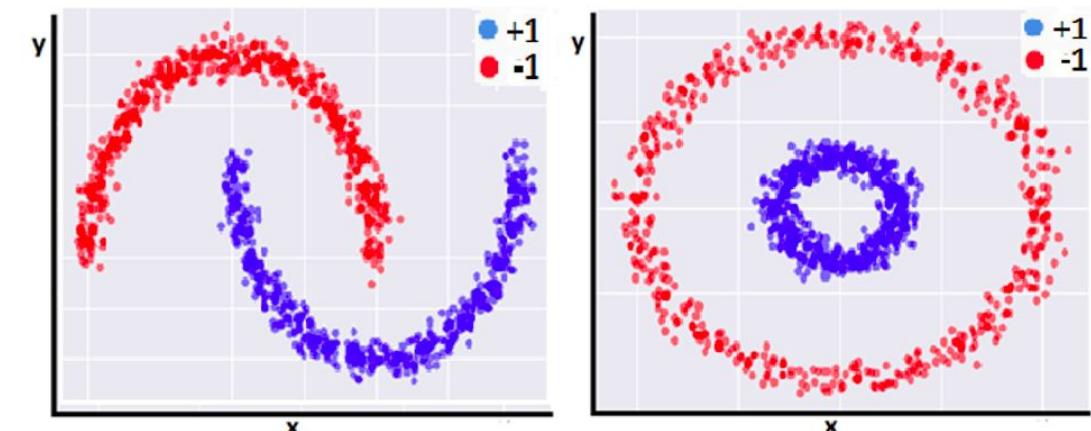
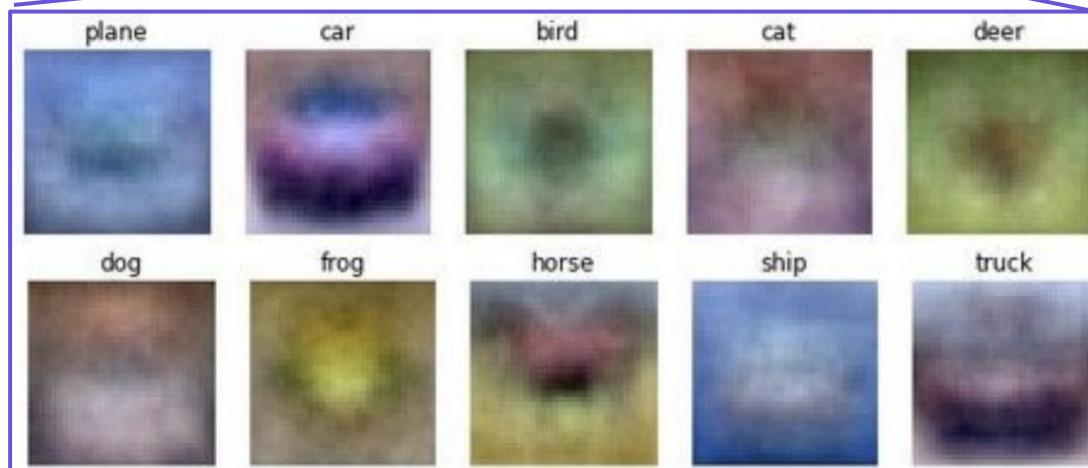
## 1.2 Issues with Linear Classifiers

Linear Model

Linear Classifier은 강력하지 않습니다. 왜냐하면

- 시각적으로) 각 클래스 당 하나의 template(카테고리)만 학습할 수 있습니다.
- 기하학적으로) 직선 형태의 decision boundary들만 그릴 수 있습니다.
  - 따라서 복잡한 관계로 이루어진 두 개의 클래스들은 완벽하게 분리할 수 없습니다.

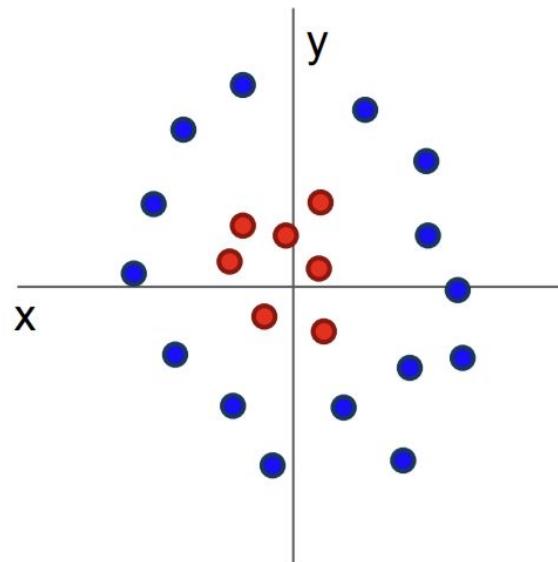
$$f(\mathbf{X}) = \mathbf{W}\mathbf{X}$$



## 1.3 Features

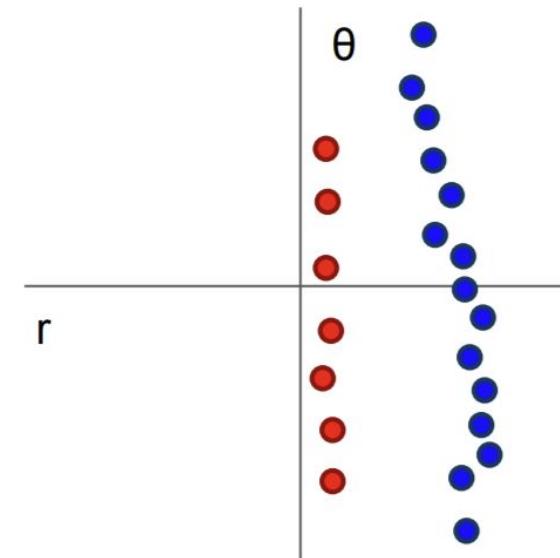
### Linear Model

- 입출력 관계를 linear classifier로 직접 연결(mapping)하는 대신 입력을 표현하기 위해 몇 가지의 특징들을 추출할 수 있습니다.
- 만약 feature space에서 입력이 선형 분류가 가능하다면 linear classifier이 잘 작동할 수 있습니다!



Original Space

$$f(x, y) = (r(x, y), \theta(x, y))$$

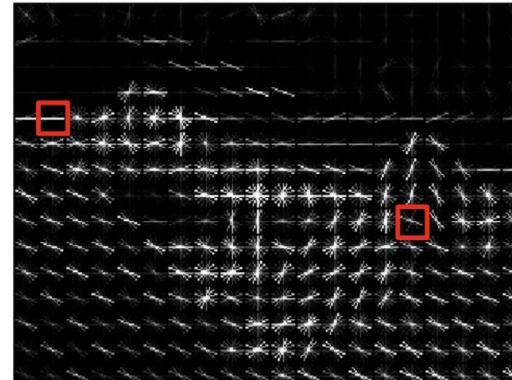
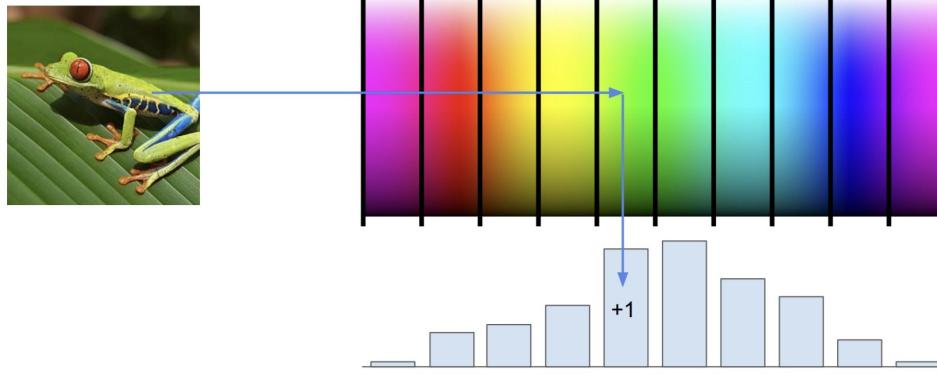


Feature Space

## 1.4 Image Features

# Linear Model

- 실제 이미지의 특징을 표현하는 방법은 복잡합니다
  - 직접 입출력(pixel-class) 관계를 mapping하는 대신 몇가지 특징을 추출하여 입력(image)를 표현합니다.
  - 예시
    - Color histogram
    - Histogram of oriented gradients (HoG)
    - Bag of words (BoW) with a pre-defined dictionary (codebook)

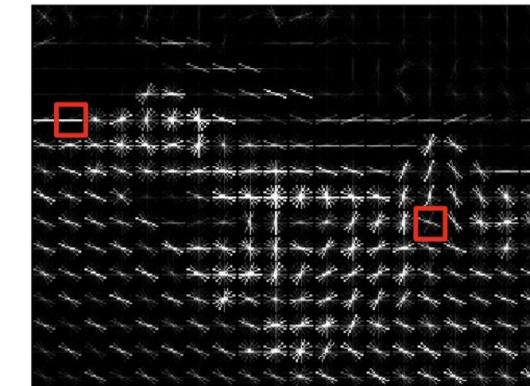
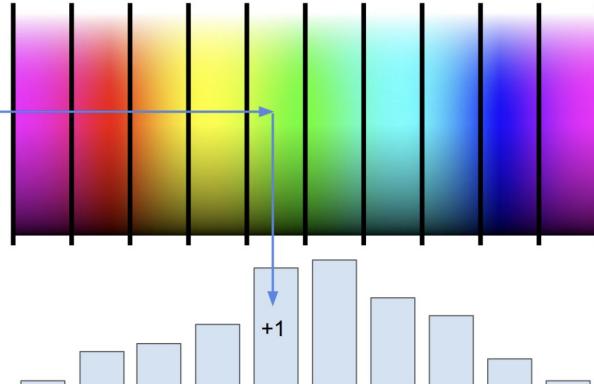


## 1.4 Image Features

Linear Model

- 예시

- Color histogram
- Histogram of oriented gradients (HoG)
- Bag of words (BoW) with a pre-defined dictionary (codebook)



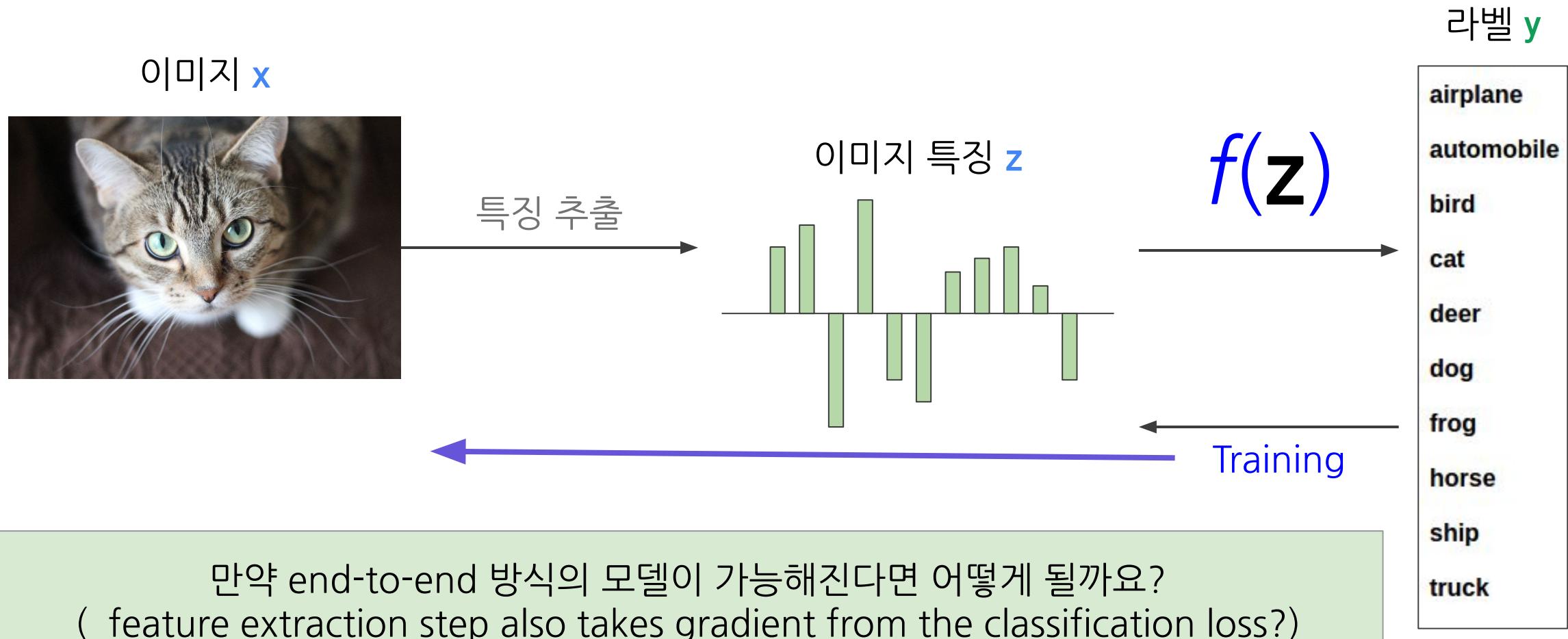
패치(pixel)당 그래디언트 방향



Clustered codebook

# 1.5 Image Classifier with pre-extracted Features

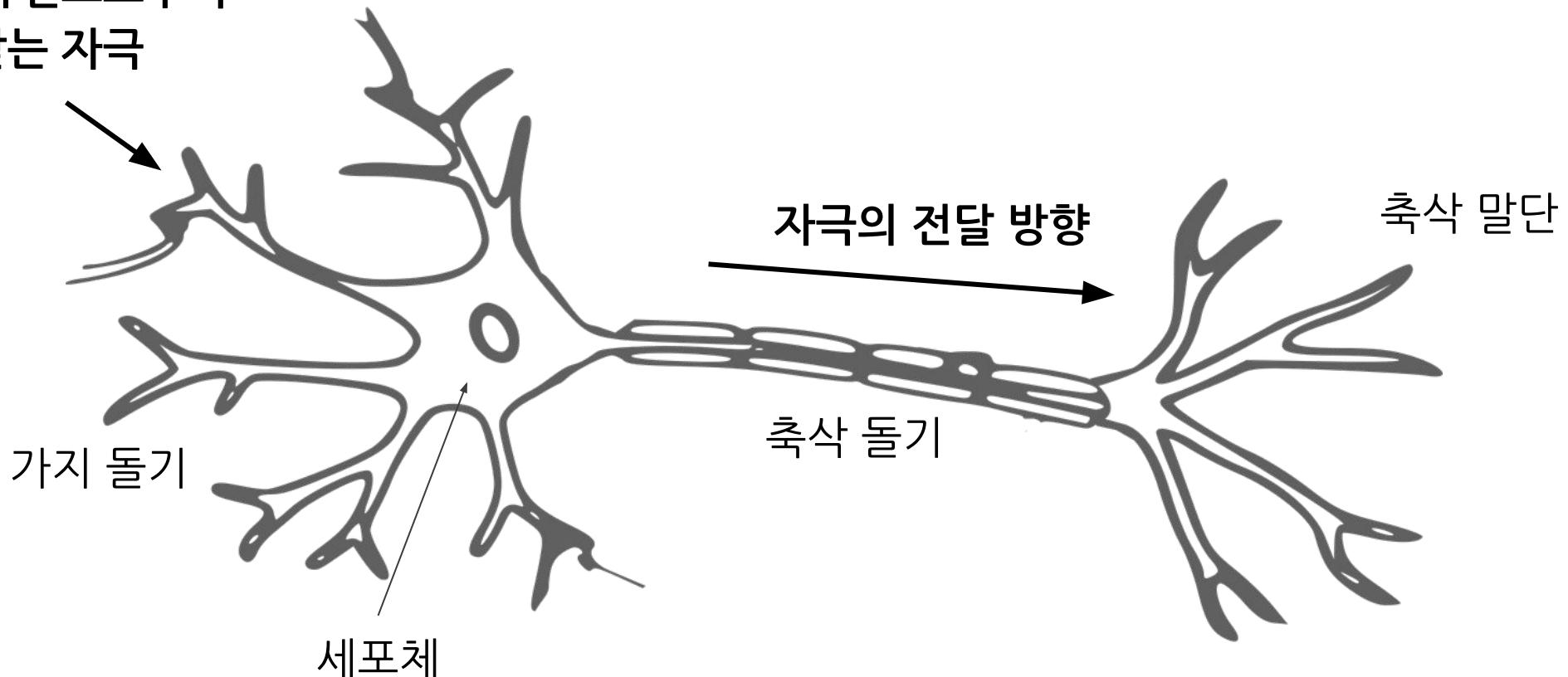
Linear Model



## 2.1 Perceptron

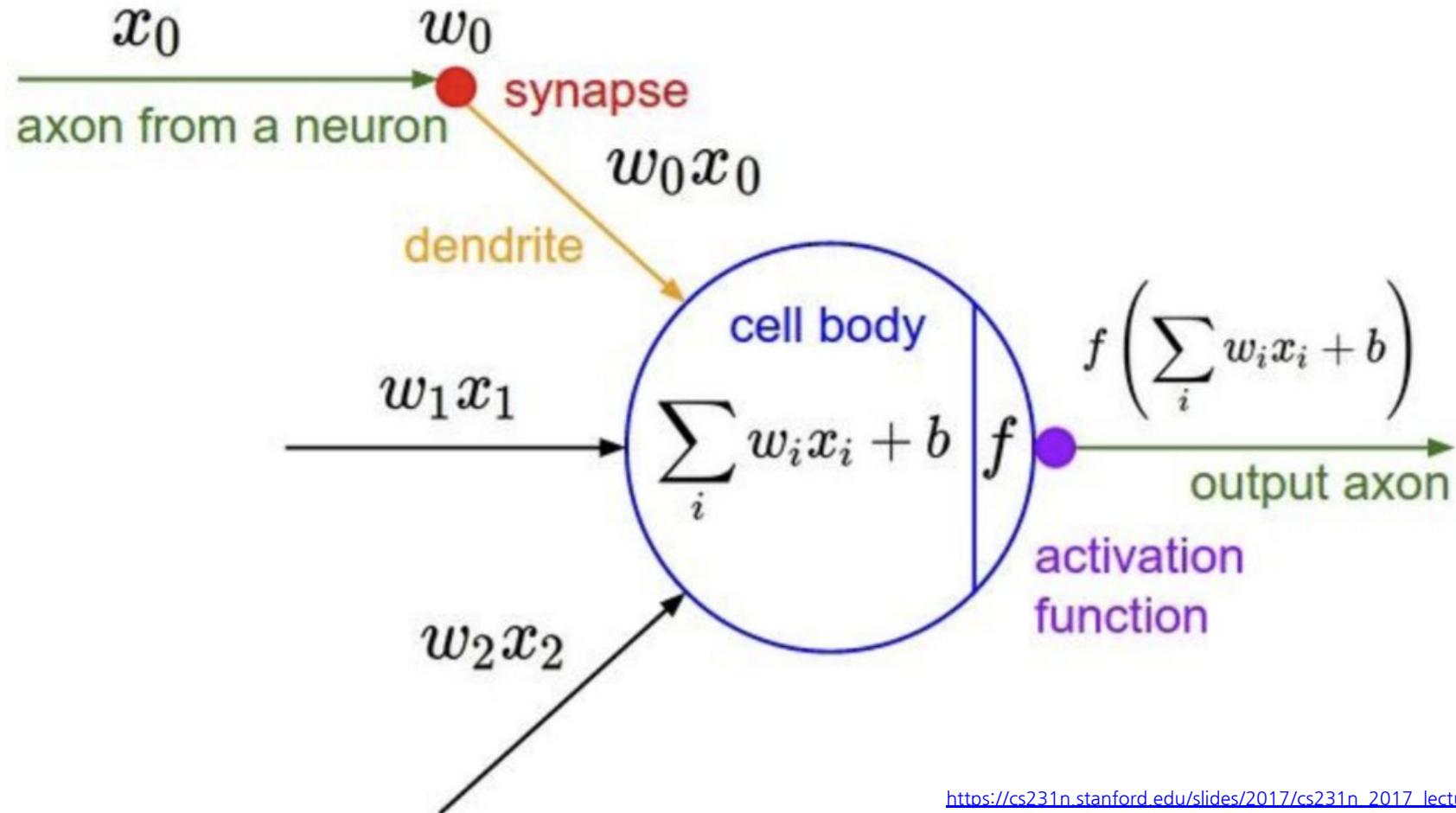
Neural Network

이전 뉴런으로부터  
받는 자극



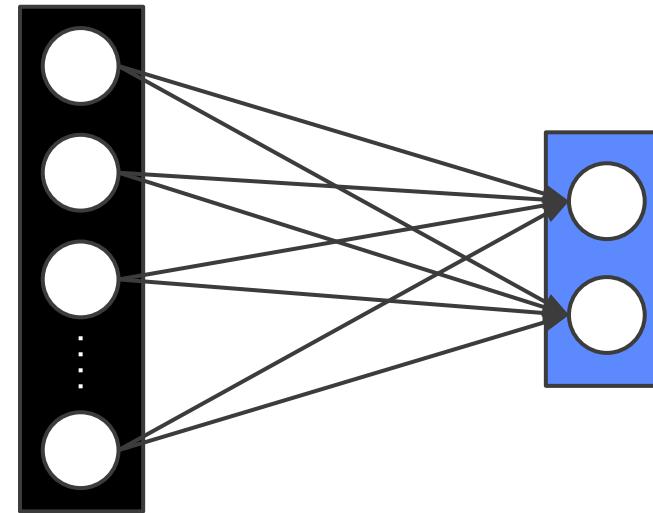
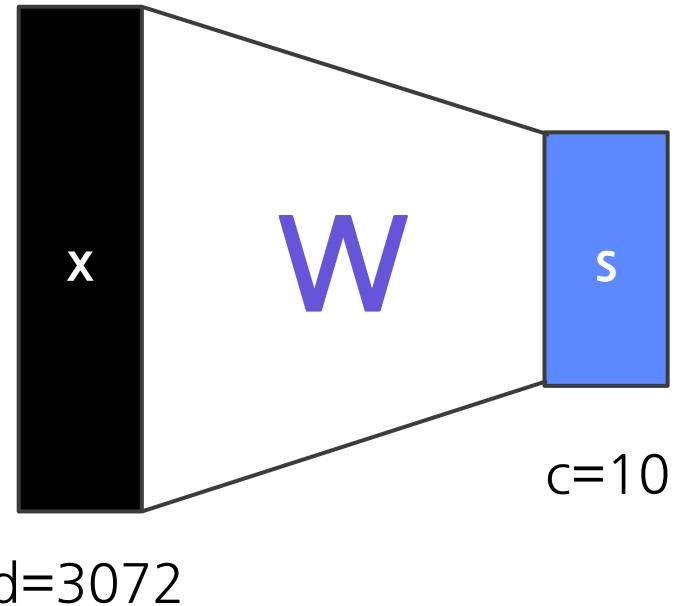
## 2.1 Perceptron

Neural Network



## 2.2 Neural Network with a Single Layer

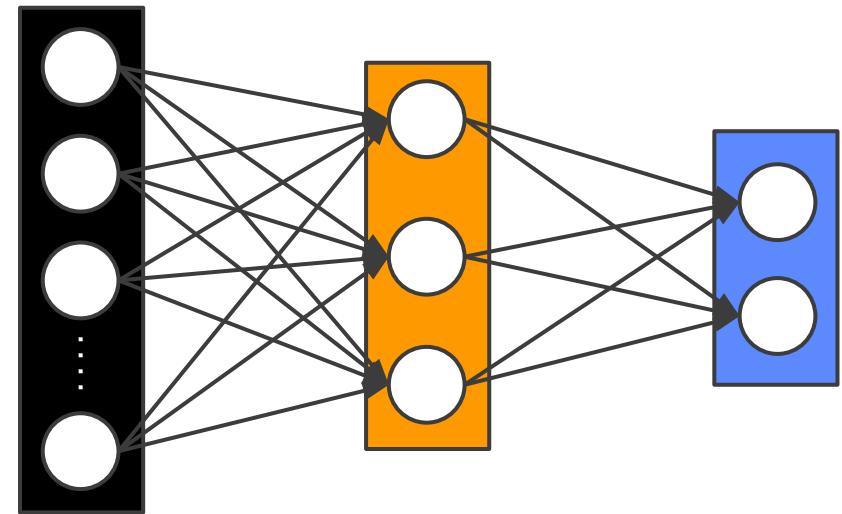
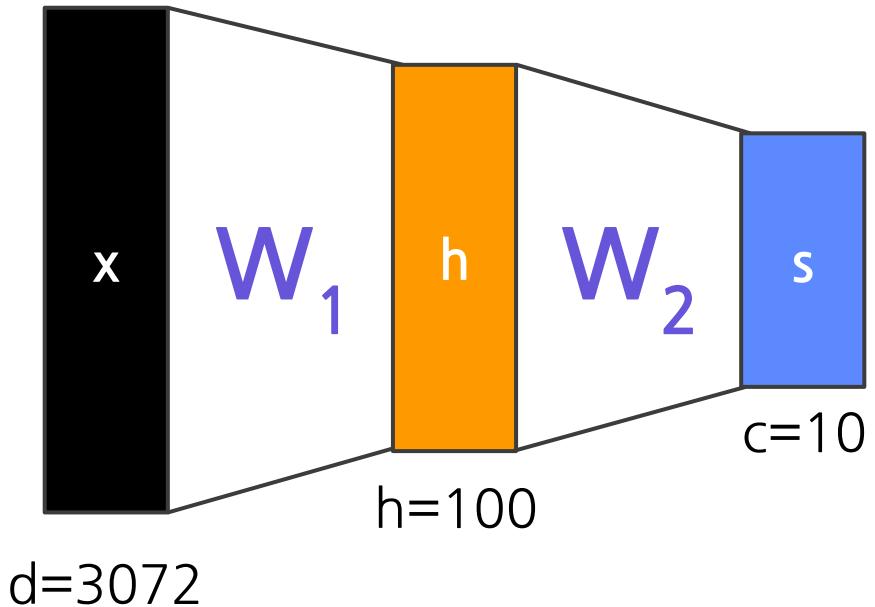
Neural Network



$$\mathbf{x} \in \mathbb{R}^d \quad \mathbf{W} \in \mathbb{R}^{c \times d} \quad \mathbf{s} \in \mathbb{R}^c$$

## 2.3 Multilayer Perceptron (MLP)

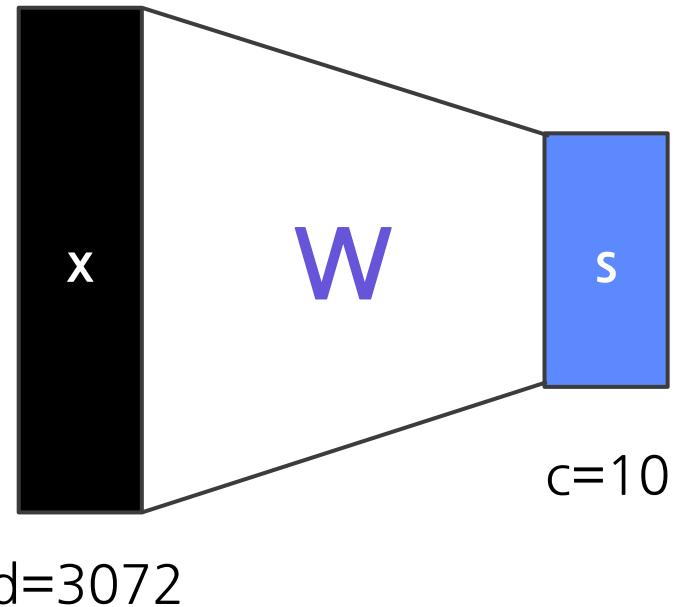
Neural Network



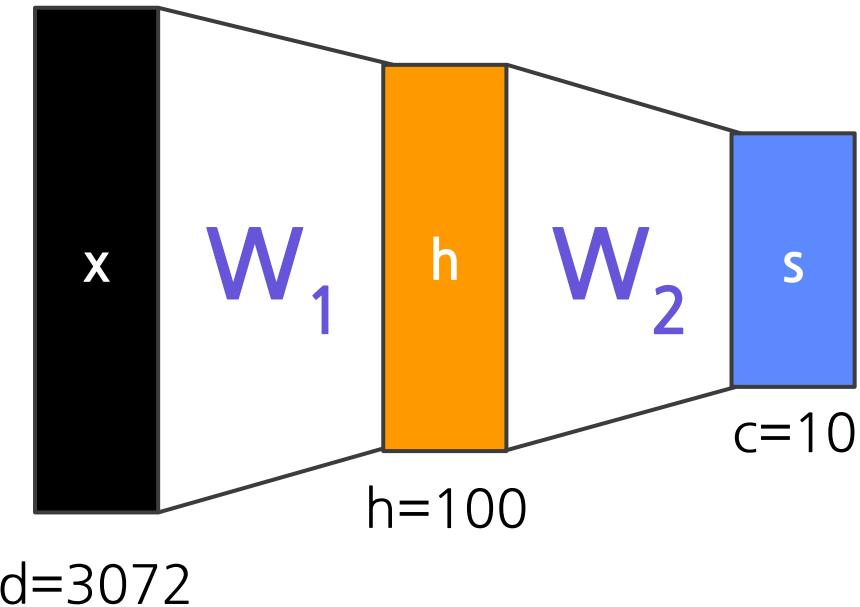
$$\mathbf{x} \in \mathbb{R}^d \quad \mathbf{W}_1 \in \mathbb{R}^{h \times d} \quad \mathbf{W}_2 \in \mathbb{R}^{c \times h} \quad \mathbf{s} \in \mathbb{R}^c$$

## 2.3 Multilayer Perceptron (MLP)

Neural Network



$$f(\mathbf{x}) = \mathbf{Wx}$$



$$f(\mathbf{x}) = \mathbf{W}_2(\mathbf{W}_1\mathbf{x})$$

## 2.3 Multilayer Perceptron (MLP)

Neural Network

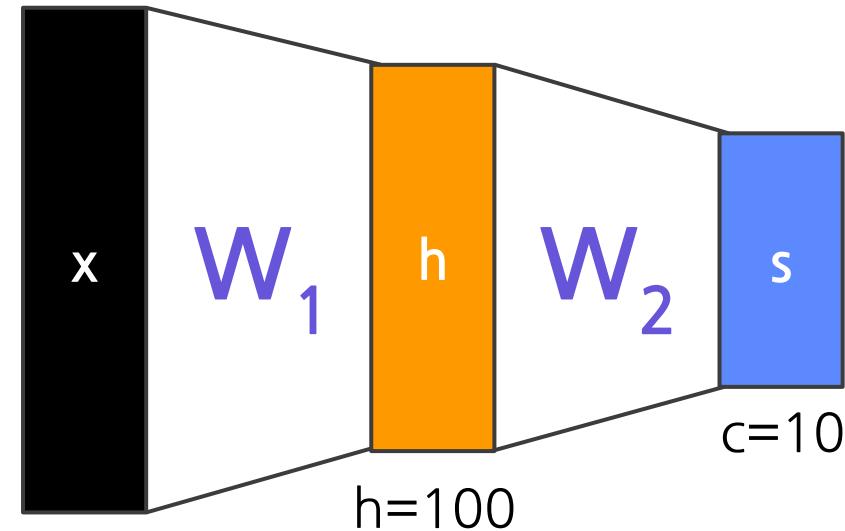
Multi-linear layers은 여전히 linear합니다.

- $\mathbf{W} = \mathbf{W}_2 \mathbf{W}_1$

어떻게 하면 non-linear하게 만들 수 있을까요?

→ Activation functions!

$$f(\mathbf{x}) = a_2(\mathbf{W}_2 a_1(\mathbf{W}_1 \mathbf{x}))$$



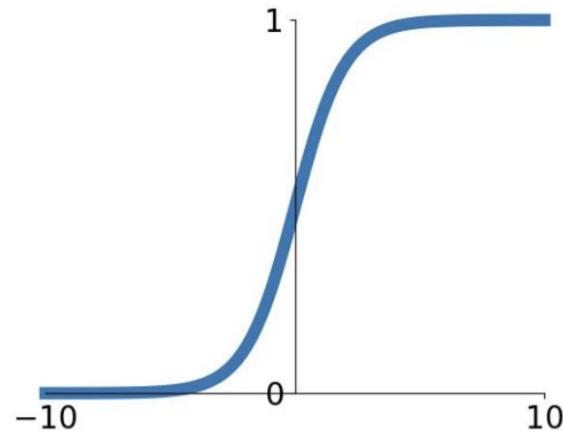
$$d=3072$$

$$f(\mathbf{x}) = \mathbf{W}_2(\mathbf{W}_1 \mathbf{x})$$

## 2.4 Activation Functions

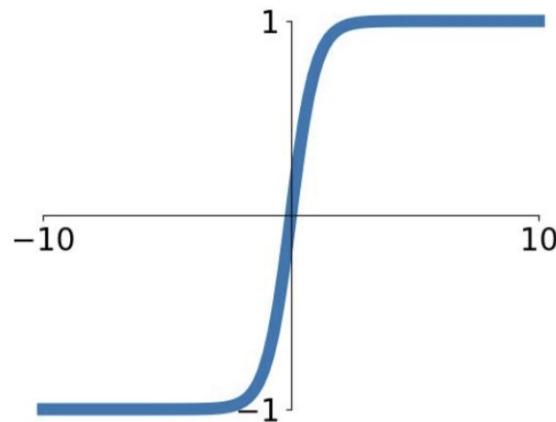
Neural Network

- Sigmoid



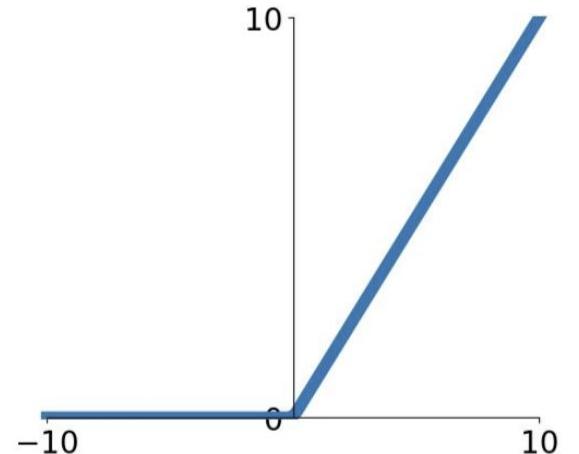
$$\sigma(x) = \frac{1}{1+e^{-x}}$$

- tanh



$$\tanh(x)$$

- ReLU



$$\max(0, x)$$

## 2.4 Sigmoid Functions

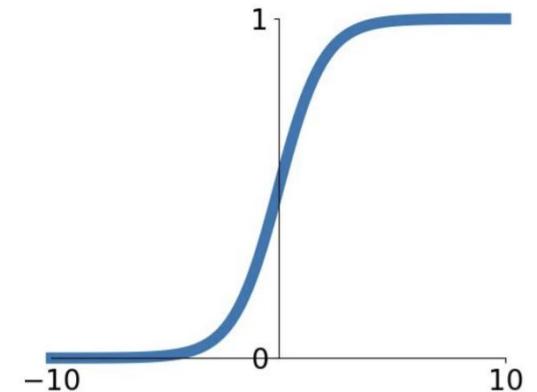
### Activation Functions

#### 특징

- 출력 값의 범위: [0, 1]
- 가장 많이 사용되었던 함수

#### 단점

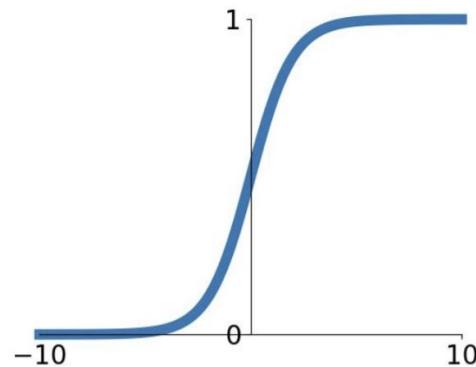
- Vanishing Gradient: input이 크거나 작을때 기울기가 0에 가까워집니다.
- 출력 값들이 zero-centered 하지 않습니다.
- $\exp()$ 의 연산이 비쌉니다.



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

## 2.4 Sigmoid Functions: Killed Gradients

Activation Functions

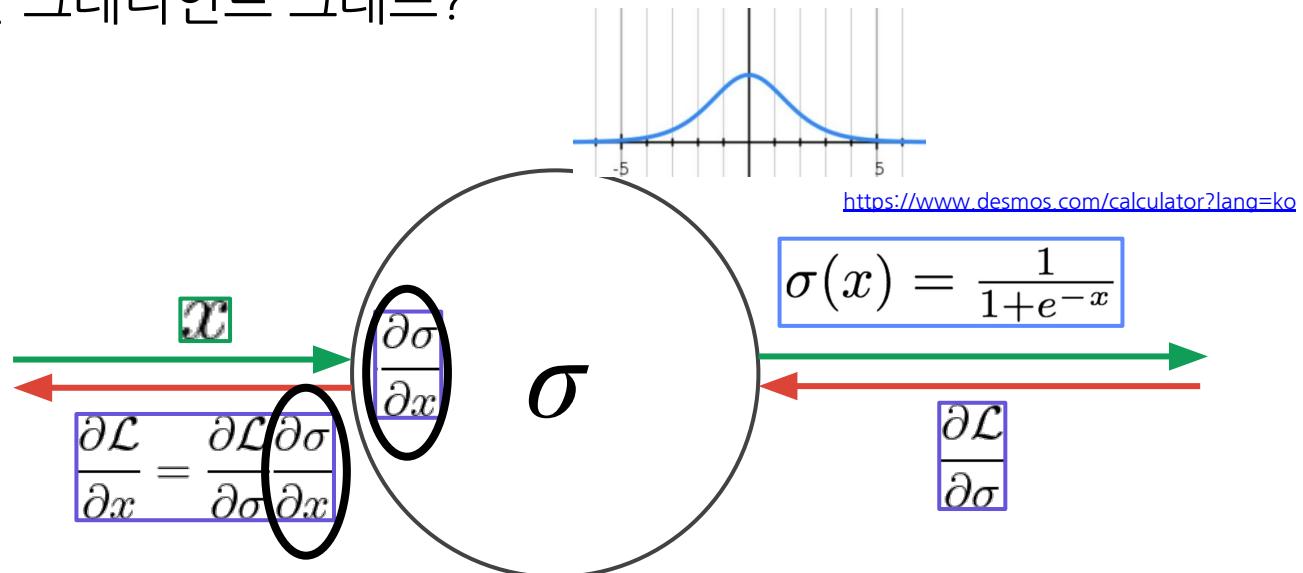


$$\sigma(x) = \frac{1}{1+e^{-x}}$$

$$\frac{\partial \sigma(x)}{\partial x} = \sigma(x)(1 - \sigma(x))$$

$$\begin{aligned}\sigma'(x) &= \frac{\partial}{\partial x} \frac{1}{1+e^{-x}} \\ &= \frac{e^{-x}}{(1+e^{-x})^2} \\ &= \frac{1}{1+e^{-x}} \cdot \frac{e^{-x}}{1+e^{-x}} \\ &= \sigma(x)(1 - \sigma(x))\end{aligned}$$

- $x$ 가 10일 때  
 $\sigma(x) \approx 1$  이기 때문에  $1 - \sigma(x) \approx 0$ , 그레디언트는 0에 수렴합니다.
- $x$ 가 -10일 때  
 $\sigma(x) \approx 0$  이기 때문에 그레디언트는 0에 수렴합니다.
- 입력  $x$ 에 대한 그레디언트 그래프?



따라서 만약 입력  $x$ 가 크다면 이 노드를 통과했을 때 기울기에  $\approx 0$ 이 곱해진다.

## 2.4 Tanh Functions

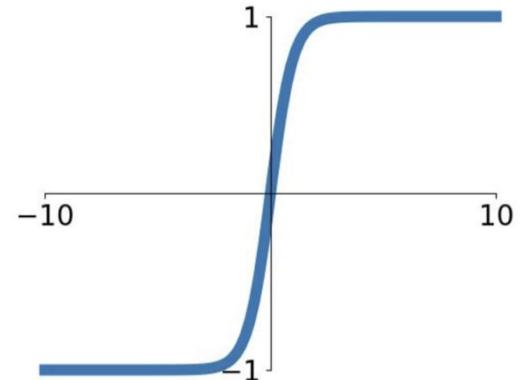
### Activation Functions

#### 특징

- 출력 값의 범위: [-1, 1].
- Zero-centered.

#### 단점

- Vanishing Gradient 문제가 여전히 발생합니다.
  - 커진 sigmoid 뉴런과 비슷하다.  
:  $\tanh(x) = 2\sigma(2x) - 1$



$\tanh(x)$

## 2.4 ReLU (Rectified Linear Unit)

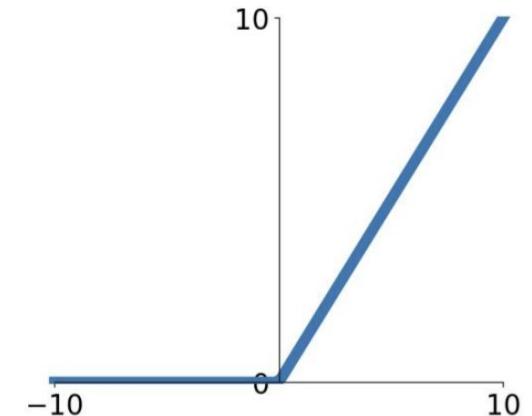
### Activation Functions

#### 특징

- + 영역에서 saturate 되지 않습니다.
- 연산이 효율적입니다.
- sigmoid나 tanh 보다 빨리 수렴합니다.

#### 단점

- 출력 값이 zero-centered 되지 않습니다.
- **Dead ReLU problem:** 출력 값이 음수라면 saturated 되는 문제가 발생합니다.
- $x=0$ 일 때 미분 불가능합니다. 0.01과 같은 약간의 기울기는 ReLU 뉴런들을 초기화하는데 도움을 줍니다.



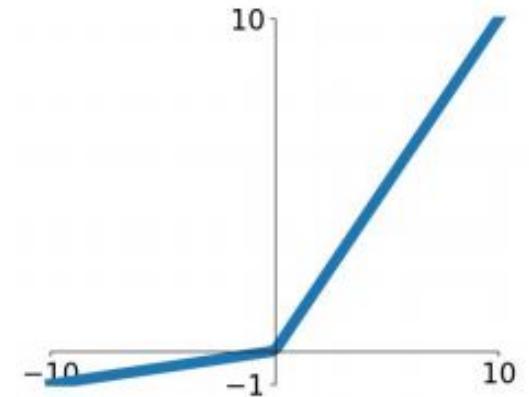
$$\max(0, x)$$

## 2.4 Leaky ReLU (Rectified Linear Unit)

Activation Functions

### 특징

- 모든 영역에서 saturate되지 않습니다.
- 연산에 효율적입니다.
- sigmoid나 tanh 보다 수렴 속도가 빠릅니다.
- No Dead ReLU 문제: gradient vanishing 되지 않습니다.



### 단점

- 추가적인 하이퍼파라미터( x가 0 미만일 때의 기울기)

$$\max(0.01x, x)$$

## 2.4 ELU (Exponential Linear Unit)

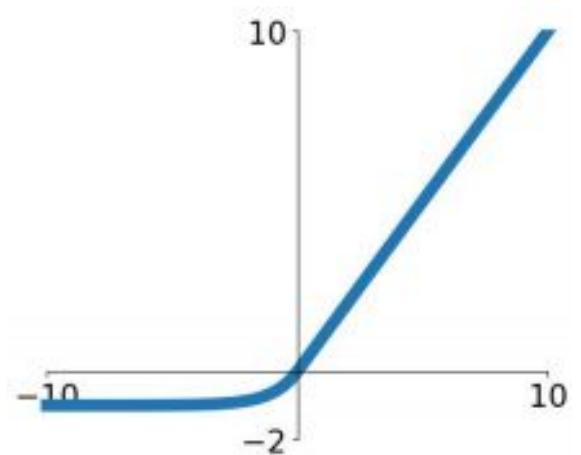
### Activation Functions

#### 특징

- ReLU의 모든 장점
- (Leaky) ReLU에 비해 saturated된 음수 지역은 견고성을 더합니다.

#### 단점

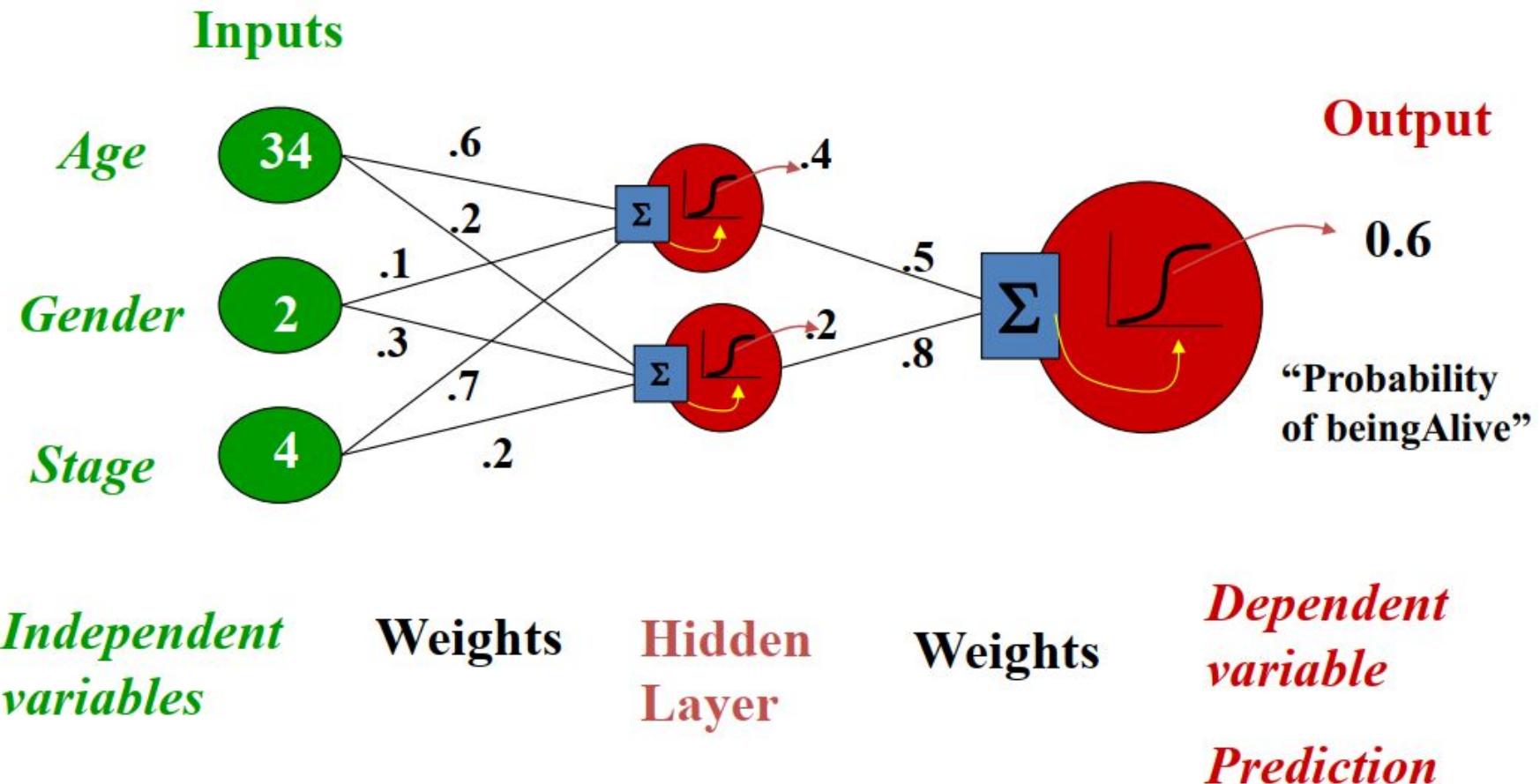
- $\exp()$ 의 연산이 비쌉니다.



$$\begin{cases} x & \text{if } x \geq 0 \\ \alpha(e^x - 1) & \text{if } x < 0 \end{cases}$$

## 2.5 An Example of Neural Network

Neural Network



## 2.6 Computing Gradients

Neural Network

(Stochastic) GD에서는 classification loss에 대한 그레디언트를 필요로 한다.

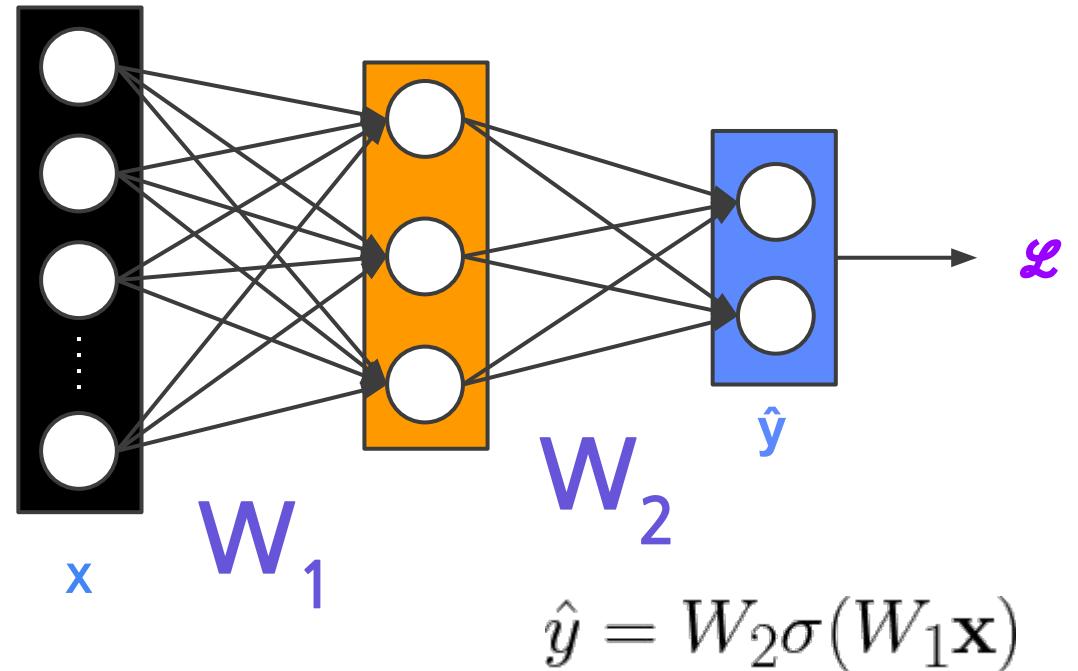
- w.r.t. 각 파라미터 (weight)

$$\frac{\partial \mathcal{L}}{\partial W_1}, \frac{\partial \mathcal{L}}{\partial W_2}$$

- 각각 파라미터를 업데이트 하기 위해 변수가 loss에 얼마나 큰 영향을 미치는지를 나타냅니다.

$$\Theta^{new} = \Theta^{old} - \alpha \nabla_{\Theta} J(\Theta)$$

$$\theta_i^{new} = \theta_i^{old} - \alpha \frac{\partial}{\partial \theta_i^{old}} J(\Theta)$$



## 2.6 Computing Gradients

Neural Network

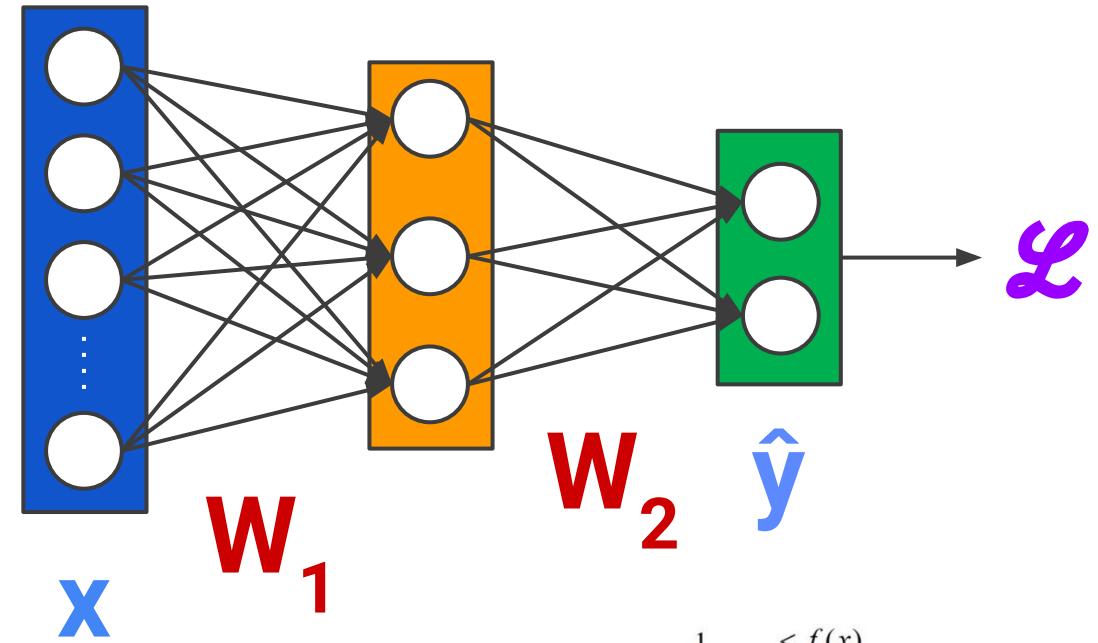
$$\mathcal{L} = (\hat{y} - y)^2$$

$$\hat{y} = W_2 \sigma(W_1 \mathbf{x})$$

$$\frac{\partial \mathcal{L}}{\partial \hat{y}} = 2(\hat{y} - y)$$

$$\frac{\partial \mathcal{L}}{\partial W_2} = \frac{\partial \mathcal{L}}{\partial \hat{y}} \frac{\partial \hat{y}}{\partial W_2} = 2(\hat{y} - y) \cdot \sigma(W_1 \mathbf{x})$$

$$\begin{aligned} \frac{\partial \mathcal{L}}{\partial W_1} &= \frac{\partial \mathcal{L}}{\partial \hat{y}} \frac{\partial \hat{y}}{\partial h} \frac{\partial h}{\partial W_1} \\ &= 2(\hat{y} - y) \cdot W_2 \cdot h(1 - h) \cdot \mathbf{x} \end{aligned}$$



$$\begin{aligned} F(x) &= \frac{1}{1+e^{-x}} &< f(x) \\ &< g(x) \\ F'(x) &= \frac{f'(x)g(x) - f(x)g'(x)}{g^2(x)} \\ &= \frac{0 - (-e^{-x})}{(1+e^{-x})^2} = \frac{1}{1+e^{-x}} \frac{e^{-x}}{1+e^{-x}} \\ &= F(x) \frac{1+e^{-x}-1}{1+e^{-x}} = F(x) \left( \frac{1+e^{-x}}{1+e^{-x}} - \frac{1}{1+e^{-x}} \right) \\ &= F(x)(1-F(x)) \end{aligned}$$

## 2.6 Implementation: 2-layer MLP

Neural Network

```
import numpy as np
from numpy.random import randn

n, d, h, c = 64, 1000, 100, 10
x, y = randn(n, d), randn(n, c)
w1, w2 = randn(d, h), randn(h, c)
learning_rate = 1e-4
```

```
for t in range(1000):
    y_0 = x.dot(w1)
    h_0 = 1 / (1 + np.exp(-y_0))
    y_pred = h_0.dot(w2)
    loss = np.square(y_pred - y).sum()
    print(t, loss)
```

```
grad_y_pred = 2.0 * (y_pred - y)
grad_w2 = h.T.dot(grad_y_pred)
grad_h = grad_y_pred.dot(w2.T)
grad_w1 = x.T.dot(grad_h * h * (1 - h))
```

```
w1 -= learning_rate * grad_w1
w2 -= learning_rate * grad_w2
```

모델 정의

(n: #examples, d: input dim,  
d: hidden dim, c: #classes)

**Forward Pass:** 현재 모델을 이용하여 예측합니다.

해석적 그래디언트들을 계산합니다.

**Gradient Descent**

## 2.6 Computing Gradients

Neural Network

- Output node:

$$\frac{\partial E_n}{\partial w_{kj}} = \frac{\partial}{\partial w_{kj}} \frac{1}{2} \sum_{k'} (y_{k'} - t_{k'})^2 = \frac{\partial}{\partial w_{kj}} \frac{1}{2} \sum_{k'} \left( \sum_j w_{kj} z_j - t_{k'} \right)^2$$

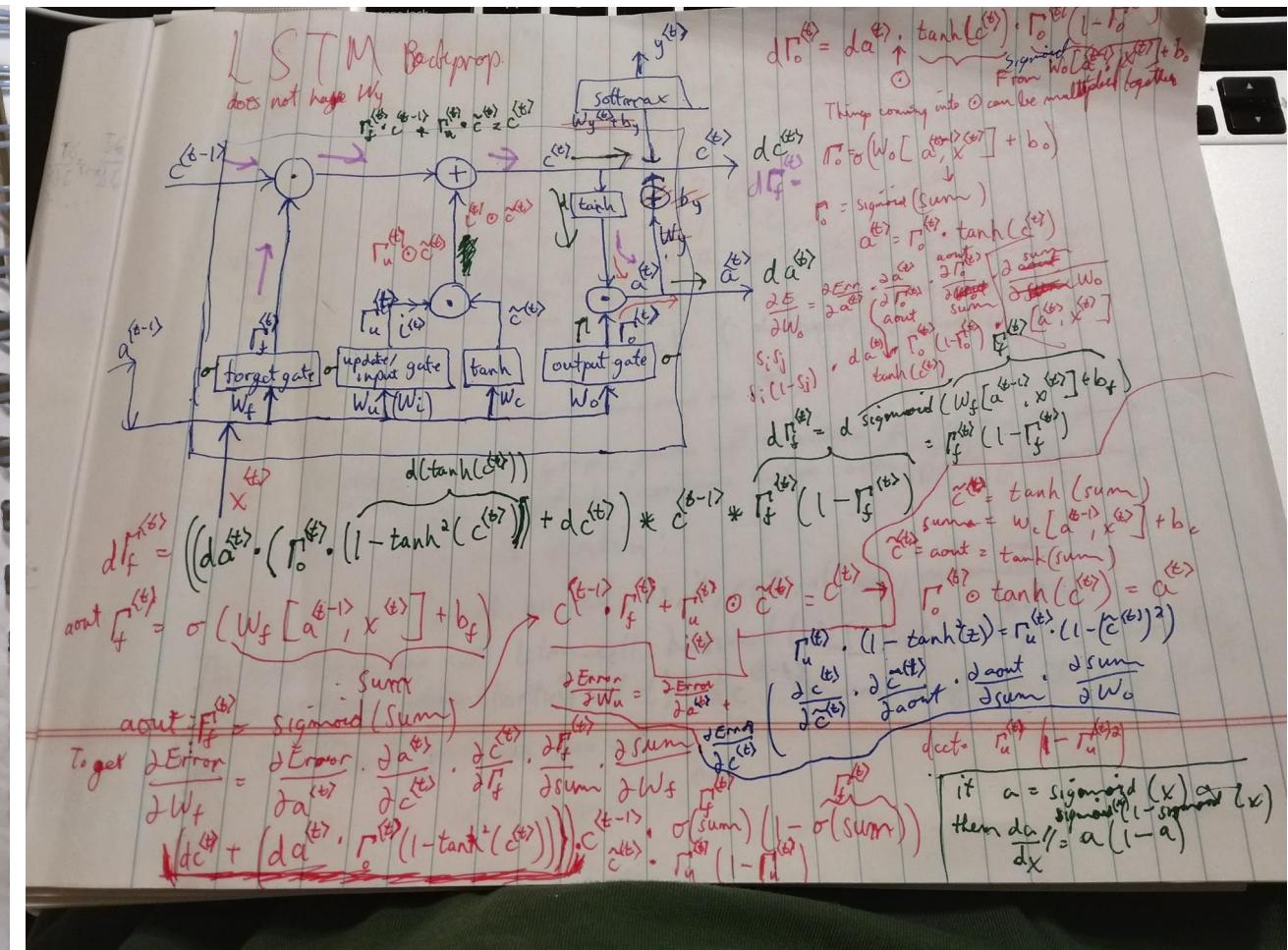
$$= (\sum_j w_{kj} z_j - t_{k'}) z_j = \underbrace{(\sum_j w_{kj} z_j - t_{k'})}_{s_k} \underbrace{z_j}_{\text{output is input}}$$

- Hidden node:

$$\frac{\partial E_n}{\partial w_{ji}} = \frac{\partial}{\partial w_{ji}} \frac{1}{2} \sum_{k'} \left( \sum_j w_{kj} h \left( \sum_i w_{ji} x_i \right) - t_{k'} \right)^2$$

$$= \sum_{k'} \left( \underbrace{\sum_j w_{kj} h \left( \sum_i w_{ji} x_i \right)}_{y_{k'}} - t_{k'} \right) \left( w_{kj} h' \left( \sum_i w_{ji} x_i \right) x_i \right)$$

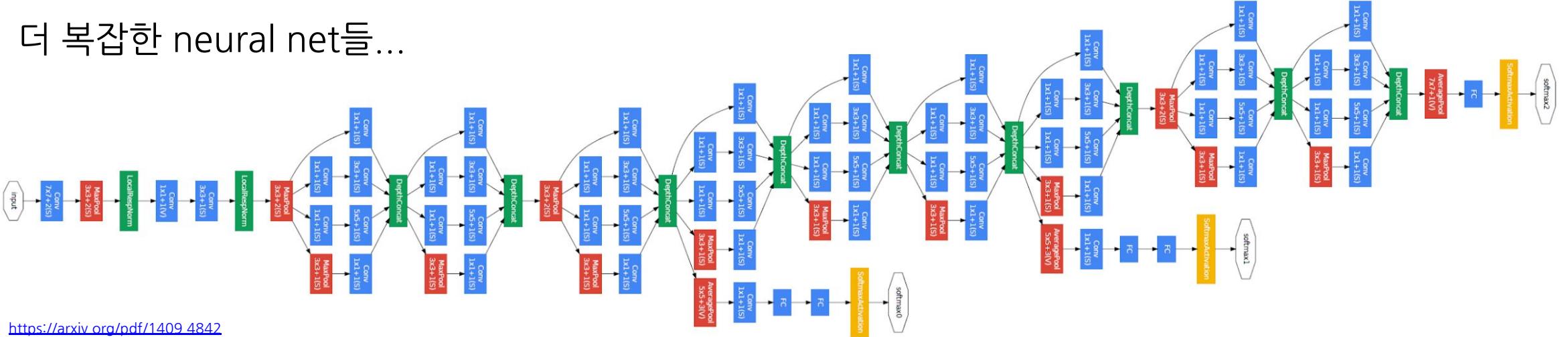
$$= \sum_{k'} \underbrace{(y_{k'} - t_{k'})}_{s_k} \underbrace{w_{kj} h'(a_j) x_i}_{\text{output is input}}$$



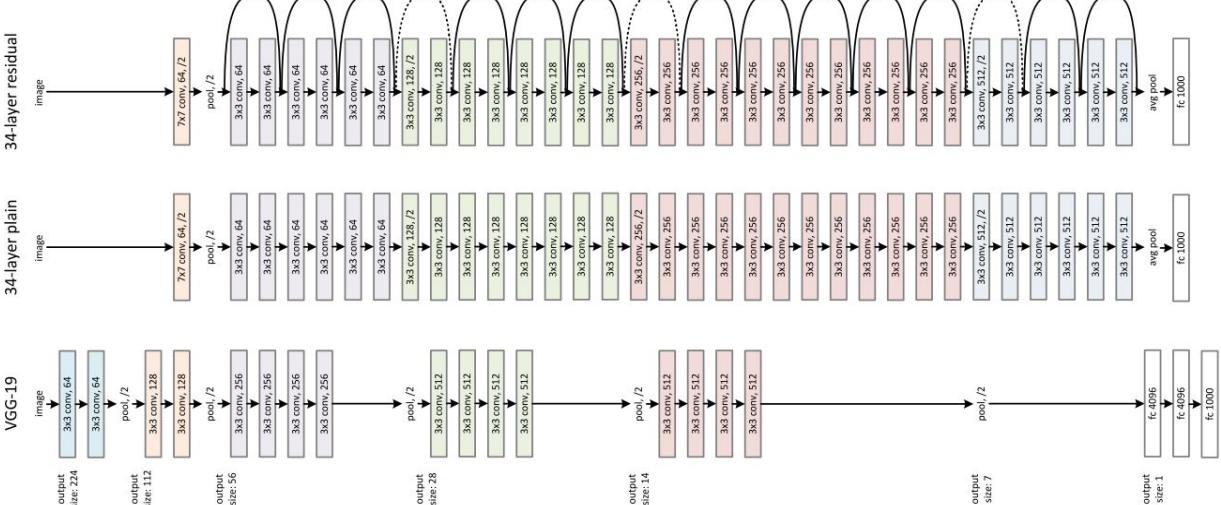
# 2.6 Computing Gradients

Neural Network

더 복잡한 neural net들...



<https://arxiv.org/pdf/1409.4842>



<https://arxiv.org/pdf/1512.03385>

# Thank you

**Prof. Jeon, Sangryul**

Computer Vision Lab.

Pusan National University, Korea

Tel: +82-51-510-2423

Web: <http://sr-jeon.github.io/>

E-mail: [srjeonn@pusan.ac.kr](mailto:srjeonn@pusan.ac.kr)