# Shallow Neural Networks
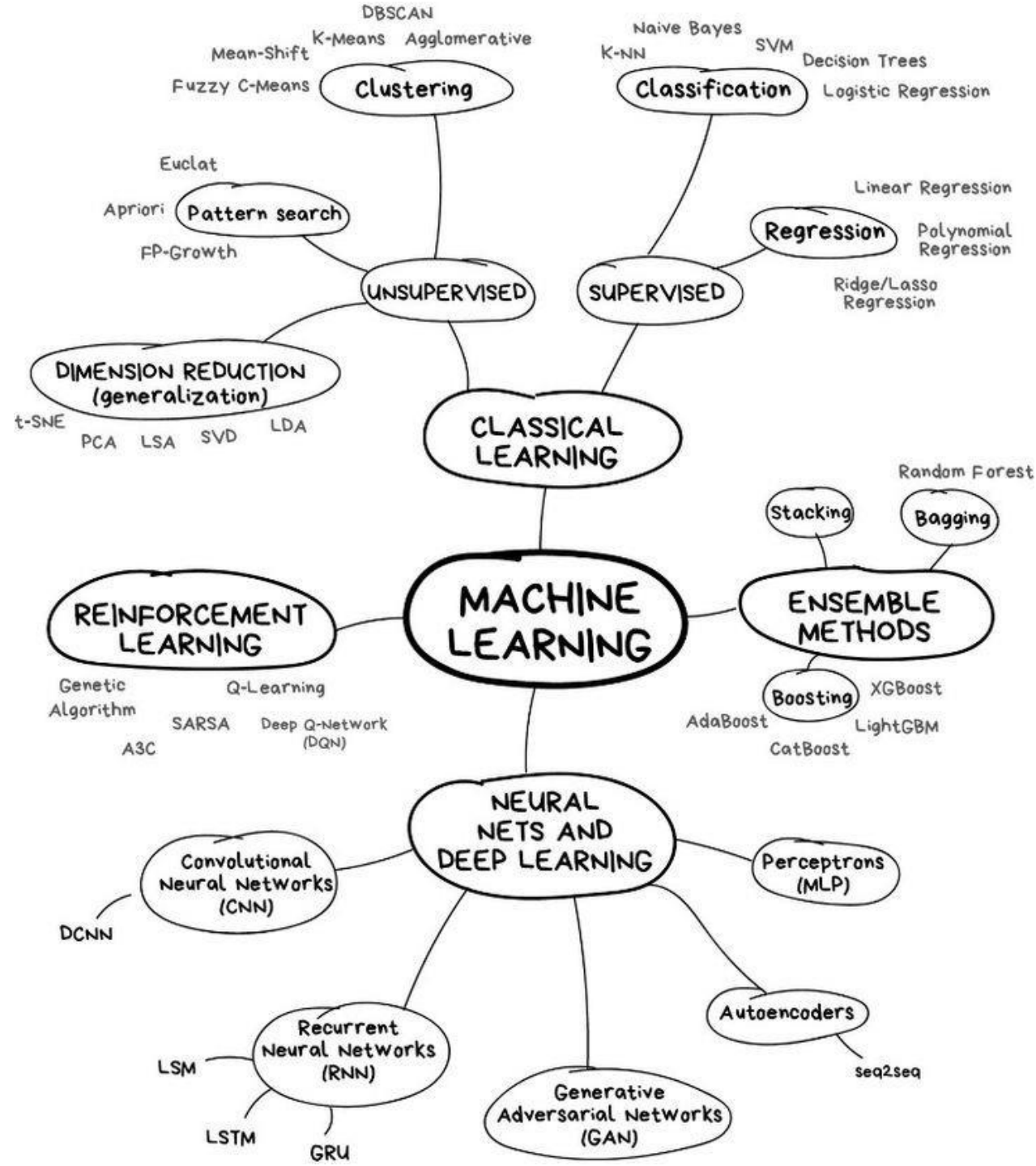
한양대학교 ERICA
소프트웨어융합대학
COLLEGE OF COMPUTING

인공지능학과
Department of
Artificial Intelligence
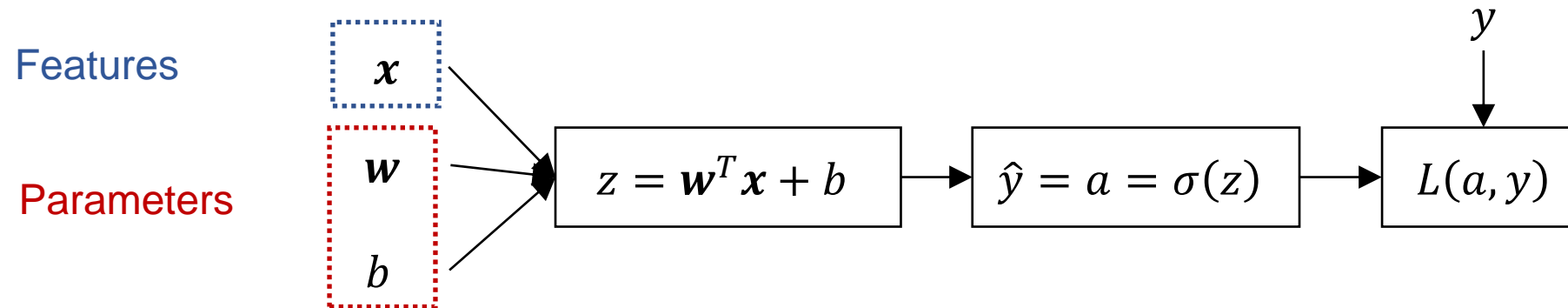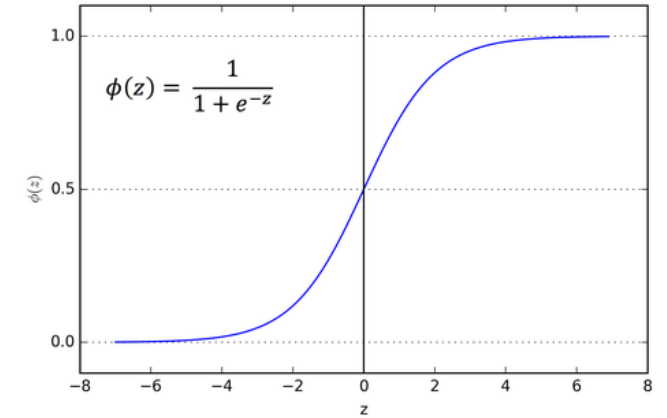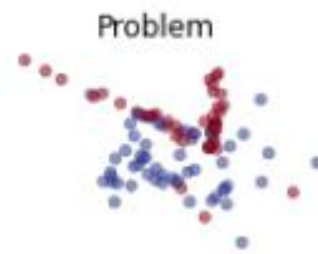
정 우 환 (whjung@hanyang.ac.kr)

**Fall 2021**

# Logistic Regression



- Output: $\hat{y} = \sigma(\boldsymbol{w}^\top \boldsymbol{x} + b)$ where $\sigma(z) = \dfrac{1}{1+e^{-z}}$

- Loss: $L(\hat{y}, y) = -y \log \hat{y} - (1-y) \log(1-\hat{y})$

Features

$\boldsymbol{x}$

Parameters

$\boldsymbol{w}$

$b$

$y$

$z = \boldsymbol{w}^T \boldsymbol{x} + b$ → $\hat{y} = a = \sigma(z)$ → $L(a, y)$

# Decision Boundaries

Problem

kNN, k=5

kNN, k=15

Logistic Regression simple

Logistic Regression basic polynomials

Decision tree

Random forest

SVM

SVM

# (Simple) XOR problem: linearly separable?



**Solution: make it more complicated**

# What is a Neural Network?
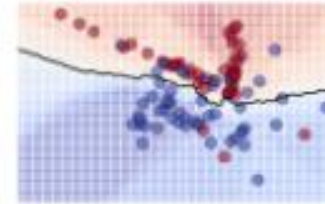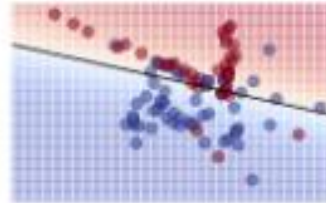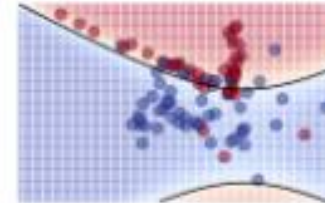


$x_1$
$x_2$
$\hat{y}$
$x_3$

$$z = \boldsymbol{w}^T \boldsymbol{x} + b \rightarrow \hat{y} = a = \sigma(z) \rightarrow L(a, y)$$

$\boldsymbol{x}$
$\boldsymbol{w}$
$b$
$y$

$1^{st}$ layer  $2^{nd}$ layer

$x_1$
$x_2$
$\hat{y}$
$x_3$

$$\boldsymbol{z}^{[1]} = \boldsymbol{W}^{[1]}\boldsymbol{x} + \boldsymbol{b}^{[1]} \rightarrow \boldsymbol{a}^{[1]} = \sigma(\boldsymbol{z}^{[1]}) \rightarrow \boldsymbol{z}^{[2]} = \boldsymbol{W}^{[2]}\boldsymbol{a}^{[1]} + \boldsymbol{b}^{[2]} \rightarrow \boldsymbol{a}^{[2]} = \sigma(\boldsymbol{z}^{[2]}) \rightarrow L(\boldsymbol{a}^{[2]}, y)$$

$\boldsymbol{x}$
$\boldsymbol{W}^{[1]}$
$\boldsymbol{b}^{[1]}$
$\boldsymbol{W}^{[2]}$
$\boldsymbol{b}^{[2]}$
$y$

# Neural Network Representation

- 2-layer neural network



Input layer

Hidden layer

Output layer

$a^{[0]} = x$

$x_1$

$x_2$

$x_3$

$a_1^{[1]}$

$a_2^{[1]}$

$a_3^{[1]}$

$a_4^{[1]}$

$a^{[2]}$

$\hat{y} = a^{[2]}$

$W^{[1]}, b^{[1]}$

$a^{[1]} \in \mathbb{R}^4$

$W^{[2]}, b^{[2]}$

# Neural Network Representation

$x_1$

$x_2$

$x_3$

$$w^T x + b \mid \sigma(z)$$

$z$     $a$

$$a = \hat{y}$$

$$z = w^T x + b$$

$$a = \sigma(z)$$

$$z_1^{[1]} = \boldsymbol{w}_1^{[1]^T} \boldsymbol{x} + b_1^{[1]}$$

$$a_1^{[1]} = \sigma\left(z_1^{[1]}\right)$$

$x_1$

$x_2$

$x_3$

$$\hat{y}$$

$$z_2^{[1]} = \boldsymbol{w}_2^{[1]^T} \boldsymbol{x} + b_2^{[1]}$$

$$a_2^{[1]} = \sigma\left(z_2^{[1]}\right)$$

# Neural Network Representation



$$z_1^{[1]} = \boldsymbol{w}_1^{[1]^T} \boldsymbol{x} + b_1^{[1]} \qquad a_1^{[1]} = \sigma\left(z_1^{[1]}\right)$$

$$z_2^{[1]} = \boldsymbol{w}_2^{[1]^T} \boldsymbol{x} + b_2^{[1]} \qquad a_2^{[1]} = \sigma\left(z_2^{[1]}\right)$$

$$z_3^{[1]} = \boldsymbol{w}_3^{[1]^T} \boldsymbol{x} + b_3^{[1]} \qquad a_3^{[1]} = \sigma\left(z_3^{[1]}\right)$$

$$z_4^{[1]} = \boldsymbol{w}_4^{[1]^T} \boldsymbol{x} + b_4^{[1]} \qquad a_4^{[1]} = \sigma\left(z_4^{[1]}\right)$$

$$z^{[2]} = \boldsymbol{w}^{[2]^T} \boldsymbol{a}^{[1]} + b^{[2]} \qquad a^{[2]} = \sigma\left(z^{[2]}\right)$$

# Neural Network Representation: Vectorization

Hidden layer

$$z_i^{[1]} = \boldsymbol{W}_i^{[1]} \boldsymbol{x} + b_i^{[1]} \qquad a_i^{[1]} = \sigma\left(z_i^{[1]}\right)$$

For $1 \le i \le h$ where $h$: # of hidden nodes

Output layer

$$z^{[2]} = \boldsymbol{w}^{[2]^T} \boldsymbol{a}^{[1]} + b^{[2]} \qquad a^{[2]} = \sigma\left(z^{[2]}\right)$$

Let $\quad \boldsymbol{x} = \begin{bmatrix} x_1 \\ x_2 \\ \dots \\ x_m \end{bmatrix} \qquad \boldsymbol{w}_i^{[1]} \in \mathbb{R}^m$

$$\boldsymbol{z}^{[1]} = \begin{bmatrix} z_1^{[1]} \\ z_2^{[1]} \\ \dots \\ z_h^{[1]} \end{bmatrix} \qquad \boldsymbol{b}^{[1]} = \begin{bmatrix} b_1^{[1]} \\ b_2^{[1]} \\ \dots \\ b_h^{[1]} \end{bmatrix} \qquad \boldsymbol{W}^{[1]} = \begin{bmatrix} \boldsymbol{w}_1^{[1]^T} \\ \boldsymbol{w}_2^{[1]^T} \\ \dots \\ \boldsymbol{w}_h^{[1]^T} \end{bmatrix} = \begin{bmatrix} w_{1,1}^{[1]} & \cdots & w_{1,m}^{[1]} \\ \vdots & \ddots & \vdots \\ w_{h,1}^{[1]} & \cdots & w_{h,m}^{[1]} \end{bmatrix} \qquad \boldsymbol{a}^{[1]} = \begin{bmatrix} a_1^{[1]} \\ a_2^{[1]} \\ \dots \\ a_h^{[1]} \end{bmatrix} \qquad \boldsymbol{w}^{[2]} = \begin{bmatrix} w_1^{[2]} \\ w_2^{[2]} \\ \dots \\ w_h^{[2]} \end{bmatrix}$$

# Neural Network with a Hidden Layer: Almost Done!



Input:

$$x \in \mathbb{R}^m$$

Parameters:

$$W^{[1]} \in \mathbb{R}^{h \times m} \qquad w^{[2]} \in \mathbb{R}^h$$

$$b^{[1]} \in \mathbb{R}^h \qquad b^{[2]} \in \mathbb{R}$$

Forward pass:

$$z^{[1]} = W^{[1]}x + b^{[1]}$$

$$a^{[1]} = \sigma(z^{[1]})$$

$$z^{[2]} = {w^{[2]}}^T a^{[1]} + b^{[2]}$$

$$\hat{y} = a^{[2]} = \sigma(z^{[2]})$$

# Activation Functions

# Neural Network with a Hidden Layer: Almost Done!



Input:

$$\boldsymbol{x} \in \mathbb{R}^m$$

Parameters:

$$\boldsymbol{W}^{[1]} \in \mathbb{R}^{h \times m} \qquad \boldsymbol{w}^{[2]} \in \mathbb{R}^h$$

$$\boldsymbol{b}^{[1]} \in \mathbb{R}^h \qquad b^{[2]} \in \mathbb{R}$$

Forward pass:

$$\boldsymbol{z}^{[1]} = \boldsymbol{W}^{[1]}\boldsymbol{x} + \boldsymbol{b}^{[1]}$$

$$\boldsymbol{a}^{[1]} = g\big(\boldsymbol{z}^{[1]}\big) \text{ where } g(.) \text{ is an activation function}$$

$$z^{[2]} = \boldsymbol{w}^{[2]}{}^T \boldsymbol{a}^{[1]} + b^{[2]}$$

$$\hat{y} = a^{[2]} = \sigma\big(z^{[2]}\big)$$

# Why we need to use non-linear activation functions?

$$x \longrightarrow \boxed{z^{[1]} = W^{[1]}x + b^{[1]}} \longrightarrow \boxed{z^{[2]} = W^{[2]}z^{[1]} + b^{[2]}} \longrightarrow z^{[2]}$$

$$z^{[2]} = W^{[2]}(W^{[1]}x + b^{[1]}) + b^{[2]}$$

$$= W^{[2]}W^{[1]}x + (W^{[2]}b^{[1]} + b^{[2]})$$

$$= W'x + b'$$

Where $W' = W^{[2]}W^{[1]}$ and $b' = W^{[2]}b^{[1]} + b^{[2]}$

Composition of linear functions => linear function

# Activation functions

There are so many activation functions ..

We'll cover some important and currently widely used activation functions
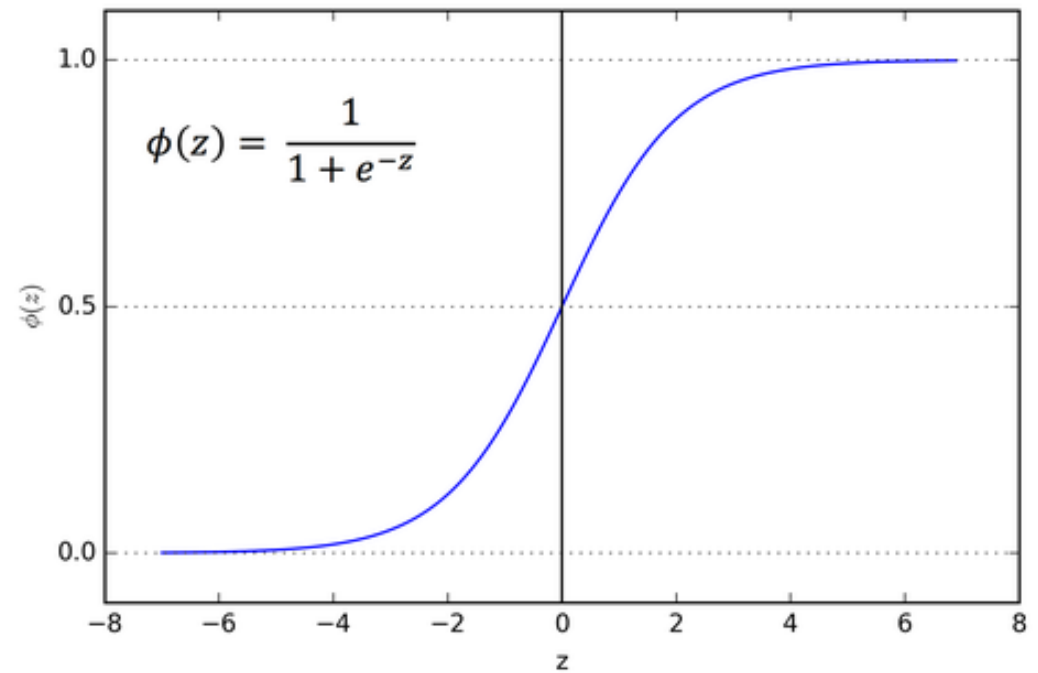
Sigmoid

tanh

ReLU

LeakyReLU

| Name | Plot | Function, $f(x)$ | Derivative of $f$, $f'(x)$ |
|---|---|---|---|
| Identity | | $x$ | $1$ |
| Binary step | | $\begin{cases} 0 & \text{if } x < 0 \\ 1 & \text{if } x \geq 0 \end{cases}$ | $\begin{cases} 0 & \text{if } x \neq 0 \\ \text{undefined} & \text{if } x = 0 \end{cases}$ |
| Logistic, sigmoid, or soft step | | $\sigma(x) = \dfrac{1}{1+e^{-x}}$[1] | $f(x)(1 - f(x))$ |
| tanh | | $\tanh(x) = \dfrac{e^x - e^{-x}}{e^x + e^{-x}}$ | $1 - f(x)^2$ |
| Rectified linear unit (ReLU)[11] | | $\begin{cases} 0 & \text{if } x \leq 0 \\ x & \text{if } x > 0 \end{cases}$ $= \max\{0, x\} = x\mathbf{1}_{x>0}$ | $\begin{cases} 0 & \text{if } x < 0 \\ 1 & \text{if } x > 0 \\ \text{undefined} & \text{if } x = 0 \end{cases}$ |
| Gaussian Error Linear Unit (GELU)[6] | | $\dfrac{1}{2}x\left(1 + \operatorname{erf}\left(\dfrac{x}{\sqrt{2}}\right)\right)$ $= x\Phi(x)$ | $\Phi(x) + x\phi(x)$ |
| Softplus[12] | | $\ln(1 + e^x)$ | $\dfrac{1}{1+e^{-x}}$ |
| Exponential linear unit (ELU)[13] | | $\begin{cases} \alpha(e^x - 1) & \text{if } x \leq 0 \\ x & \text{if } x > 0 \end{cases}$ with parameter $\alpha$ | $\begin{cases} \alpha e^x & \text{if } x < 0 \\ 1 & \text{if } x > 0 \\ 1 & \text{if } x = 0 \text{ and } \alpha = 1 \end{cases}$ |
| Square linear unit (SQLU)[14] | | $\begin{cases} x & \text{if } x > 0.0 \\ \alpha(x + \frac{x^2}{4}) & \text{if } -2.0 \leq x \leq 0 \\ -\alpha & \text{if } x < -2.0 \end{cases}$ | $\begin{cases} 1 & \text{if } x > 0.0 \\ 1 + \frac{x}{2} & \text{if } -2.0 \leq x \leq 0 \\ 0 & \text{if } x < -2.0 \end{cases}$ |
| Scaled exponential linear unit (SELU)[15] | | $\lambda\begin{cases} \alpha(e^x - 1) & \text{if } x < 0 \\ x & \text{if } x \geq 0 \end{cases}$ with parameters $\lambda = 1.0507$ and $\alpha = 1.67326$ | $\lambda\begin{cases} \alpha e^x & \text{if } x < 0 \\ 1 & \text{if } x \geq 0 \end{cases}$ |
| Leaky rectified linear unit (Leaky ReLU)[16] | | $\begin{cases} 0.01x & \text{if } x < 0 \\ x & \text{if } x \geq 0 \end{cases}$ | $\begin{cases} 0.01 & \text{if } x < 0 \\ 1 & \text{if } x \geq 0 \end{cases}$ |
| Parameteric rectified linear unit (PReLU)[17] | | $\begin{cases} \alpha x & \text{if } x < 0 \\ x & \text{if } x \geq 0 \end{cases}$ with parameter $\alpha$ | $\begin{cases} \alpha & \text{if } x < 0 \\ 1 & \text{if } x \geq 0 \end{cases}$ |
| ElliotSig,[18][19] softsign[20][21] | | $\dfrac{x}{1 + |x|}$ | $\dfrac{1}{(1 + |x|)^2}$ |

# Sigmoid

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

- Range: $(0,1)$

- Derivative

  $$\frac{d\sigma(x)}{dx} = \sigma(x)\big(1 - \sigma(x)\big)$$

  - $0 < \dfrac{d\sigma(x)}{dx} \leq 0.25$
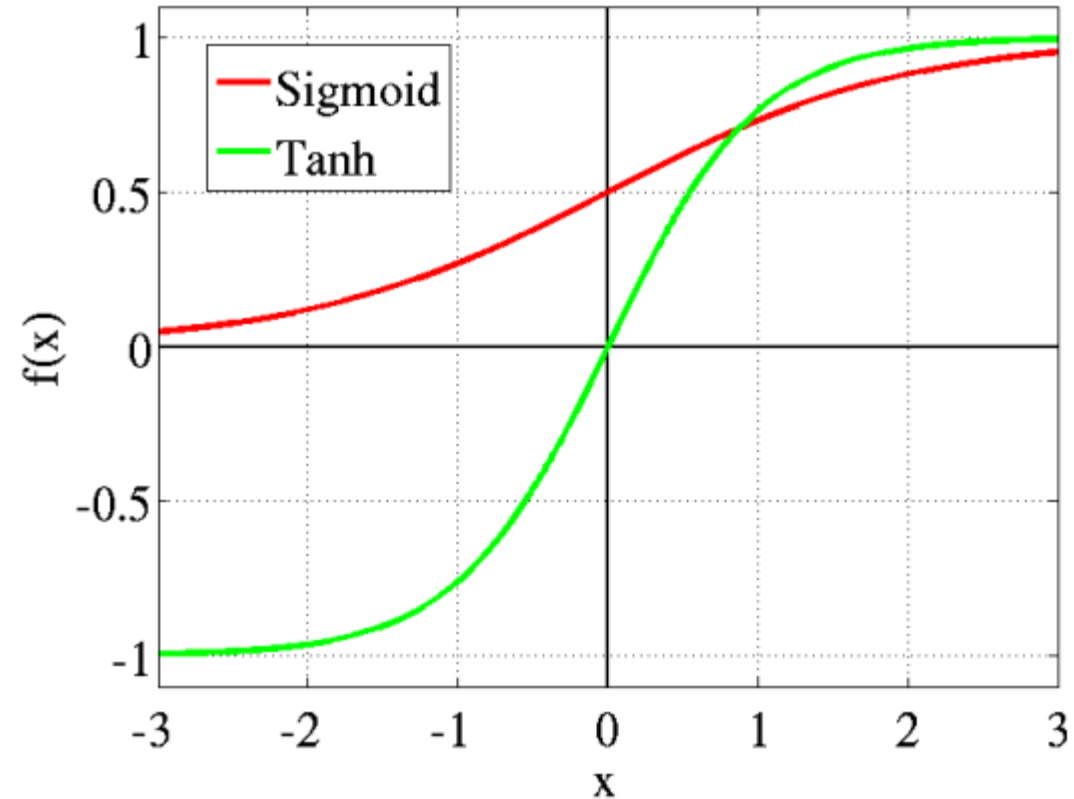
# Hyperbolic tangent: tanh

- $\tanh x = \dfrac{e^x - e^{-x}}{e^x + e^{-x}}$

$$= 2\sigma(2x) - 1$$

- Range: $(-1,1)$
- Derivative

$$\frac{d \tanh x}{dx} = 1 - \tanh^2 x$$

- $0 < \dfrac{d \tanh x}{dx} \le 1$

# Vanishing gradient

"The term vanishing gradient refers to the fact that in a feedforward network (FFN) the backpropagated error signal typically decreases (or increases) exponentially as a function of the distance from the final layer" by Jason Brownlee

https://machinelearningmastery.com/how-to-fix-vanishing-gradients-using-the-rectified-linear-activation-function/

- Chain rule (for a Deep NN)

$$\frac{\partial L}{\partial z^{[n]}} \frac{\partial z^{[n]}}{\partial z^{[n-1]}} \cdots \frac{\partial z^{[3]}}{\partial z^{[2]}} \frac{\partial z^{[2]}}{\partial z^{[1]}} \frac{\partial z^{[1]}}{\partial w}$$

# Rectified linear unit: ReLU

- $f(x) = \max\{0, x\}$
- Range: $[0, \infty)$
- Derivative

$$\frac{df(x)}{dx} = \begin{cases} 0 & if\ x < 0 \\ 1 & if\ x > 0 \\ undefined & if\ x = 0 \end{cases}$$

sigmoid

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

ReLU

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

# Leaky ReLU

- $f(x) = \begin{cases} ax & if\ x < 0 \\ x & if\ x \geq 0 \end{cases}$
  - $a \ll 1$ (e.g, $a = 0.01$)
- Range: $(-\infty, \infty)$
- Derivative

$$\frac{df(x)}{dx} = \begin{cases} a & if\ x < 0 \\ 1 & if\ x > 0 \\ undefined & if\ x = 0 \end{cases}$$

ReLU



Leaky ReLU

# Activation functions

- There is no rule but … in many cases

- Output layer
  - Sigmoid

- Hidden layer
  - tanh, ReLU, LeakyReLU

My personal opinion:
If the performance of two things are similar,
**Use the simpler one! (Use ReLU!)**

# Neural Network with a Hidden Layer: Done!



$$z^{[1]} = W^{[1]}x + b^{[1]}$$

$a^{[1]} = g(z^{[1]})$ where $g(.)$ is an activation function

$$z^{[2]} = w^{[2]^T}a^{[1]} + b^{[2]}$$

$$\hat{y} = a^{[2]} = \sigma(z^{[2]})$$

You may use tanh, ReLU, Leaky ReLU for $g$

# Gradient Descent

Shallow Neural Network with tanh

# Shallow Neural Network



Parameters: $\boldsymbol{\theta} = \{\boldsymbol{W}^{[1]}, \boldsymbol{b}^{[1]}, \boldsymbol{w}^{[2]}, b^{[2]}\}$

Loss $L(\hat{y}, y)$
Cost $J(\theta)$

Gradient descent

$$\boldsymbol{\theta} := \boldsymbol{\theta} - \eta \nabla J = \boldsymbol{\theta} - \eta \frac{\nabla L}{m}$$

$$\boldsymbol{z}^{[1]} = \boldsymbol{W}^{[1]}\boldsymbol{x} + \boldsymbol{b}^{[1]}$$

$$\boldsymbol{a}^{[1]} = \tanh \boldsymbol{z}^{[1]}$$

$$z^{[2]} = \boldsymbol{w}^{[2]^T}\boldsymbol{a}^{[1]} + b^{[2]}$$

$$\hat{y} = a^{[2]} = \sigma(z^{[2]})$$

# Partial Derivatives: 2nd layer

$x$

$y$

$W^{[1]}$

$b^{[1]}$

$$z^{[1]} = W^{[1]}x + b^{[1]}$$

$$a^{[1]} = \tanh z^{[1]}$$

$w^{[2]}$
$b^{[2]}$

$$z^{[2]} = w^{[2]^T}a^{[1]} + b^{[2]}$$

$$a^{[2]} = \sigma(z^{[2]})$$

$$L(a^{[2]}, y)$$

$$L(a^{[2]}, y) = -y \log a^{[2]} - (1-y)\log(1 - a^{[2]})$$

- $\dfrac{\partial L(a^{[2]},y)}{\partial b^{[2]}} = \dfrac{\partial L(a^{[2]},y)}{\partial a^{[2]}} \dfrac{\partial a^{[2]}}{\partial z^{[2]}} \dfrac{\partial z^{[2]}}{\partial b^{[2]}}$

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

$$= \left(\frac{-y}{a^{[2]}} + \frac{1-y}{1-a^{[2]}}\right)\sigma(z^{[2]})\left(1 - \sigma(z^{[2]})\right)$$

$$= \left(\frac{-y}{a^{[2]}} + \frac{1-y}{1-a^{[2]}}\right)a^{[2]}(1 - a^{[2]})$$

$$= -y(1 - a^{[2]}) + a^{[2]}(1-y) = a^{[2]} - y$$

Do you remember?

# Partial Derivatives: 2<sup>nd</sup> layer

$\boldsymbol{x}$

$\boldsymbol{W}^{[1]}$

$\boldsymbol{b}^{[1]}$

$y$

$$\boxed{\boldsymbol{z}^{[1]} = \boldsymbol{W}^{[1]}\boldsymbol{x} + \boldsymbol{b}^{[1]}} \rightarrow \boxed{\boldsymbol{a}^{[1]} = \tanh \boldsymbol{z}^{[1]}} \rightarrow \boxed{z^{[2]} = \boldsymbol{w}^{[2]^T}\boldsymbol{a}^{[1]} + b^{[2]}} \rightarrow \boxed{a^{[2]} = \sigma(z^{[2]})} \rightarrow \boxed{L(a^{[2]}, y)}$$

$\boldsymbol{W}^{[2]}$

$b^{[2]}$

$$L(a^{[2]}, y) = -y \log a^{[2]} - (1 - y) \log(1 - a^{[2]})$$

- $\dfrac{\partial L(a^{[2]}, y)}{\partial w_i^{[2]}} = \dfrac{\partial L(a^{[2]}, y)}{\partial a^{[2]}} \dfrac{\partial a^{[2]}}{\partial z^{[2]}} \dfrac{\partial z^{[2]}}{\partial w_i^{[2]}}$

$$= \left(a^{[2]} - y\right) \frac{\partial z^{[2]}}{\partial w_i^{[2]}}$$

$$= \left(a^{[2]} - y\right) a_i^{[1]}$$

Parameters $\boldsymbol{W}^{[1]}, \boldsymbol{b}^{[1]}, \boldsymbol{w}^{[2]}, b^{[2]}$

# Partial Derivatives: 1$^{\text{st}}$ layer

$x$

$\boldsymbol{W}^{[1]}$

$\boldsymbol{b}^{[1]}$

$y$

$$\boxed{\boldsymbol{z}^{[1]} = \boldsymbol{W}^{[1]}\boldsymbol{x} + \boldsymbol{b}^{[1]}} \rightarrow \boxed{\boldsymbol{a}^{[1]} = \tanh \boldsymbol{z}^{[1]}} \rightarrow \boxed{z^{[2]} = \boldsymbol{w}^{[2]^T}\boldsymbol{a}^{[1]} + b^{[2]}} \rightarrow \boxed{a^{[2]} = \sigma(z^{[2]})} \rightarrow \boxed{L(a^{[2]}, y)}$$

$\boldsymbol{W}^{[2]}$

$b^{[2]}$

$$L(a^{[2]}, y) = -y \log a^{[2]} - (1 - y)\log(1 - a^{[2]})$$

- $\dfrac{\partial L(a^{[2]}, y)}{\partial b_i^{[1]}} = \dfrac{\partial L(a^{[2]}, y)}{\partial a^{[2]}} \dfrac{\partial a^{[2]}}{\partial z^{[2]}} \dfrac{\partial z^{[2]}}{\partial a_i^{[1]}} \dfrac{\partial a_i^{[1]}}{\partial z_i^{[1]}} \dfrac{\partial z_i^{[1]}}{\partial b_i^{[1]}}$

$$\dfrac{d \tanh x}{dx} = 1 - \tanh^2 x$$

$$= (a^{[2]} - y) \dfrac{\partial z^{[2]}}{\partial a_i^{[1]}} \dfrac{\partial a_i^{[1]}}{\partial z_i^{[1]}} \dfrac{\partial z_i^{[1]}}{\partial b_i^{[1]}}$$

$$= (a^{[2]} - y)w_i^{[2]}\left(1 - \tanh^2 z_i^{[1]}\right) \cdot 1$$

$$= (a^{[2]} - y)w_i^{[2]}\left(1 - a_i^{[1]^2}\right)$$

# Partial Derivatives: 1ˢᵗ layer

$x$

$\boldsymbol{W}^{[1]}$

$\boldsymbol{b}^{[1]}$

$$\boxed{\boldsymbol{z}^{[1]} = \boldsymbol{W}^{[1]}\boldsymbol{x} + \boldsymbol{b}^{[1]}} \rightarrow \boxed{\boldsymbol{a}^{[1]} = \tanh \boldsymbol{z}^{[1]}} \rightarrow \boxed{z^{[2]} = \boldsymbol{w}^{[2]^T}\boldsymbol{a}^{[1]} + b^{[2]}} \rightarrow \boxed{a^{[2]} = \sigma(z^{[2]})} \rightarrow \boxed{L(a^{[2]}, y)}$$

$\boldsymbol{W}^{[2]}$

$b^{[2]}$

$y$

$$L(a^{[2]}, y) = -y \log a^{[2]} - (1-y)\log(1 - a^{[2]})$$
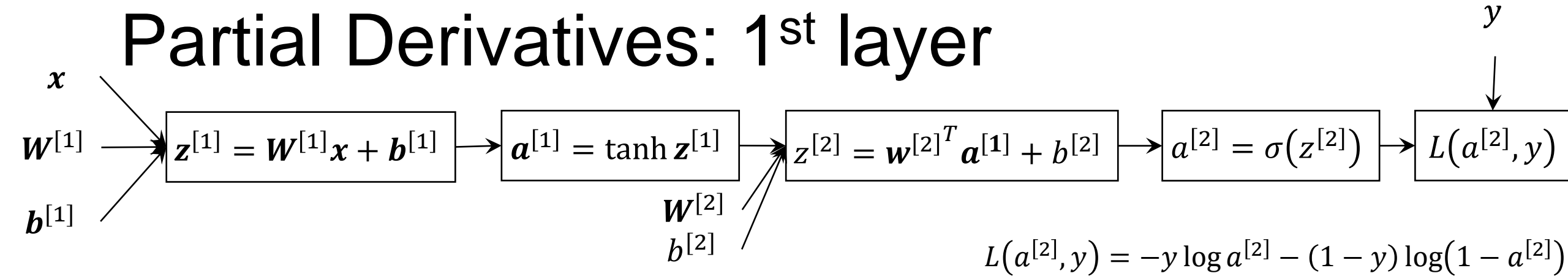
$$\bullet \quad \frac{\partial L(a^{[2]}, y)}{\partial W_{ij}^{[1]}} = \frac{\partial L(a^{[2]}, y)}{\partial a^{[2]}} \frac{\partial a^{[2]}}{\partial z^{[2]}} \frac{\partial z^{[2]}}{\partial a_i^{[1]}} \frac{\partial a_i^{[1]}}{\partial z_i^{[1]}} \frac{\partial z_i^{[1]}}{\partial W_{ij}^{[1]}}$$

$$= \left(a^{[2]} - y\right) w_i^{[2]} \left(1 - a_i^{[1]^2}\right) \frac{\partial z_i^{[1]}}{\partial W_{ij}^{[1]}}$$

$$= \left(a^{[2]} - y\right) w_i^{[2]} \left(1 - a_i^{[1]^2}\right) x_j$$

# Partial derivatives

- $\dfrac{\partial L\left(a^{[2]},y\right)}{\partial b^{[2]}} = a^{[2]} - y$

- $\dfrac{\partial L\left(a^{[2]},y\right)}{\partial w_i^{[2]}} = \left(a^{[2]} - y\right) a_i^{[1]}$

- $\dfrac{\partial L\left(a^{[2]},y\right)}{\partial b_i^{[1]}} = \left(a^{[2]} - y\right) w_i^{[2]} \left(1 - a_i^{[1]^2}\right)$

- $\dfrac{\partial L\left(a^{[2]},y\right)}{\partial W_{ij}^{[1]}} = \left(a^{[2]} - y\right) w_i^{[2]} \left(1 - a_i^{[1]^2}\right) x_j$

# Gradient Descent

- Let $\boldsymbol{x} \in \mathbb{R}^n$

- We have $h$ hidden units

- Randomly initialize the parameters
  - $\boldsymbol{W}^{[1]} \in \mathbb{R}^{h \times n}$
  - $\boldsymbol{b}^{[1]} \in \mathbb{R}^h$
  - $\boldsymbol{w}^{[2]} \in \mathbb{R}^h$
  - $b^{[2]} \in \mathbb{R}$

- For each epoch
  - $d\boldsymbol{W}^{[1]} = \mathbf{0} \in \mathbb{R}^{h \times n}, d\boldsymbol{b}^{[1]} = \mathbf{0} \in \mathbb{R}^h$
  - $d\boldsymbol{w}^{[2]} = \mathbf{0} \in \mathbb{R}^h, db^{[2]} = 0 \in \mathbb{R}$
  - For $(x, y) \in D$
    - $dw_i^{[2]} += \frac{\partial L(a^{[2]}, y)}{\partial w_i^{[2]}}$ for $1 \le i \le h$
    - $db^{[2]} += \frac{\partial L(a^{[2]}, y)}{\partial b^{[2]}}$
    - $dW_{ij}^{[1]} += \frac{\partial L(a^{[2]}, y)}{\partial W_{ij}^{[1]}}$ for $1 \le i \le n, 1 \le j \le h$
    - $db_i^{[1]} += \frac{\partial L(a^{[2]}, y)}{\partial b_i^{[1]}}$ for $1 \le i \le h,$
  - $\boldsymbol{W}^{[1]} -= \eta \cdot d\boldsymbol{W}^{[1]}/|D|$
  - $\boldsymbol{b}^{[1]} -= \eta \cdot d\boldsymbol{b}^{[1]}/|D|$
  - $\boldsymbol{w}^{[2]} -= \eta \cdot d\boldsymbol{w}^{[1]}/|D|$
  - $b^{[2]} -= \eta \cdot db^{[2]}/|D|$

# Initialize with the same value?



$$\frac{\partial L(a^{[2]}, y)}{\partial b_i^{[1]}} = \left(1 - a_i^{[1]^2}\right) w_i^{[2]} \left(a^{[2]} - y\right)$$

$$\frac{\partial L(a^{[2]}, y)}{\partial b_1^{[1]}} = \frac{\partial L(a^{[2]}, y)}{\partial b_2^{[1]}} = \frac{\partial L(a^{[2]}, y)}{\partial b_3^{[1]}} = \cdots$$

# Programming Exercise

Shallow neural network for XOR

# Slicing a numpy array

Slicing

```
In [77]: a = np.array([2,4,6,8,10])
         b = np.array([1,3])
```

```
In [78]: a[b]
```

```
Out[78]: array([4, 8])
```

# NumPy where

- numpy.where(*condition*[, x, y])
  - Parameters
    - Condition: array_like, bool
    - x,y: array_like
  - Returns:
    - Out: ndarray
      - An array with elements from *x* where *condition* is True, and elements from *y* elsewhere.

Examples

```
>>> a = np.arange(10)
>>> a
array([0, 1, 2, 3, 4, 5, 6, 7, 8, 9])
>>> np.where(a < 5, a, 10*a)
array([ 0,  1,  2,  3,  4, 50, 60, 70, 80, 90])
```

# Slicing a numpy array with condition

```
In [84]: a = np.array([2,4,6,8,10,12])
         b = np.array([1,3])
```

```
In [85]: a[a>5]
```

```
Out[85]: array([ 6,  8, 10, 12])
```

```
In [86]: a[a%3==0]
```

```
Out[86]: array([ 6, 12])
```

# Data preparation

**XOR data**

In [4]:
```python
x_seeds = np.array([(0,0),(1,0),(0,1),(1,1)],dtype=np.float)
y_seeds = np.array([0,1,1,0])
```

In [5]:
```python
N = 1000
idxs = np.random.randint(0,4,N)
```

In [6]:
```python
X = x_seeds[idxs]
Y = y_seeds[idxs]
```

In [7]:
```python
X += np.random.normal(scale = 0.25, size = X.shape)
```
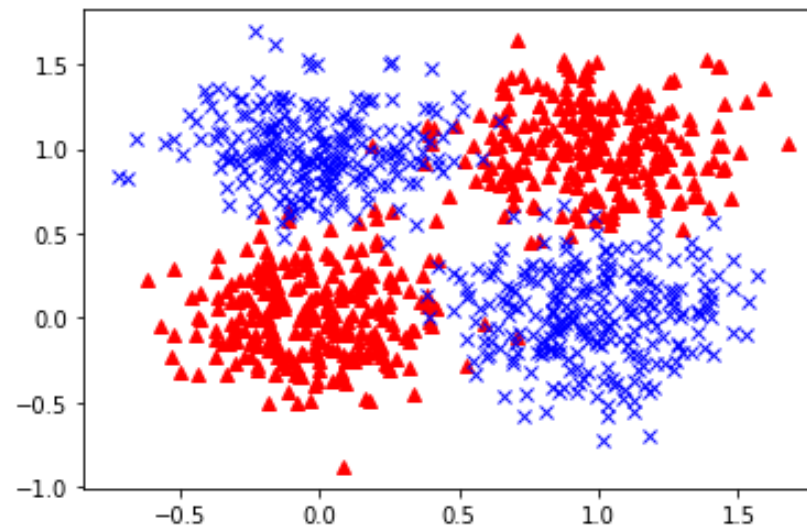
# Data plotting

- import matplotlib.pyplot as plt

**Plot data** ¶

```
In [8]: idxs_1 = np.where(Y==1)
        idxs_0 = np.where(Y==0)

In [9]: X_0 = X[idxs_0]
        Y_0 = Y[idxs_0]

n [10]: X_1 = X[idxs_1]
        Y_1 = Y[idxs_1]

n [11]: #plt.clf()
        plt.plot(X_0[:,0],X_0[:,1],"r^")
        plt.plot(X_1[:,0],X_1[:,1],"bx")
        plt.show()
```

# Model

```python
class shallow_neural_network():
    def __init__(self, num_input_features, num_hiddens):
        self.num_input_features = num_input_features
        self.num_hiddens = num_hiddens

        self.W1 = np.random.normal(size = (num_hiddens, num_input_features))
        self.b1 = np.random.normal(size = num_hiddens)
        self.W2 = np.random.normal(size = num_hiddens)
        self.b2 = np.random.normal(size = 1)


    def sigmoid(self,z):
        return 1/(1 + np.exp(-z))

    def predict(self,x):
        z1 = np.matmul(self.W1,x) + self.b1
        a1 = np.tanh(z1)
        z2 = np.matmul(self.W2,a1)+ self.b2
        a2 = self.sigmoid(z2)
        return a2, (z1,a1,z2, a2)
```

```python
model = shallow_neural_network(2,3)
```

# Train

```python
def train(X, Y, model, lr = 0.1):
    dW1 = np.zeros_like(model.W1)
    db1 = np.zeros_like(model.b1)
    dW2 = np.zeros_like(model.W2)
    db2 = np.zeros_like(model.b2)
    m = len(X)
    cost = 0.0
    for x,y in zip(X,Y):
        a2, (z1,a1,z2, _) = model.predict(x)
        if y == 1:
            cost -= np.log(a2)
        else:
            cost -= np.log(1-a2)

        diff = a2-y
        # layer 2
        # db2
        db2 += diff

        # dw2 - todo: remove for-loops
        for i in range(model.num_hiddens):
            dW2[i] += a1[i]*diff
        #layer 1
        # db1 - todo: remove for-loops
        for i in range(model.num_hiddens):
            db1[i] += (1-a1[i]**2)*model.W2[i]*diff
        # db2 - todo: remove for-loops
        for i in range(model.num_hiddens):
            for j in range(model.num_input_features):
                dW1[i,j] += x[j]*(1-a1[i]**2)*model.W2[i]*diff

    cost /= m
    model.W1 -= lr * dW1/m
    model.b1 -= lr * db1/m
    model.W2 -= lr * dW2/m
    model.b2 -= lr * db2/m

    return cost
```

```python
for epoch in range(100):
    cost = train(X,Y, model, 1.0)
    if epoch %10 == 0:
        print(epoch, cost)
```

```
0 [1.16857084]
10 [0.68480409]
20 [0.65975895]
30 [0.60186859]
40 [0.52126503]
50 [0.4437902]
60 [0.3854394]
70 [0.34575853]
80 [0.31936125]
90 [0.30160166]
```

# Test

In [73]: `model.predict((1,1))[0].item()`

Out [73]: 0.0720145593561928

In [74]: `model.predict((1,0))[0].item()`

Out [74]: 0.8872553390238738

In [75]: `model.predict((0,1))[0].item()`

Out [75]: 0.85885430328791

In [76]: `model.predict((0,0))[0].item()`

Out [76]: 0.0755082353589167

# 숙제

- Train code에서 for-loops 없이 동작하도록 수정
- Hint
  - Layer 2 (Logistic regression참고)
  - Layer 1

numpy.**outer**(*a, b, out=None*)  [source]

Compute the outer product of two vectors.

Given two vectors, a = [a0, a1, ..., aM] and b = [b0, b1, ..., bN], the outer product [1] is:

```
[[a0*b0  a0*b1 ... a0*bN ]
 [a1*b0     .
 [ ...            .
 [aM*b0           aM*bN ]]
```

**Train**

```python
def train(X, Y, model, lr = 0.1):
    dW1 = np.zeros_like(model.W1)
    db1 = np.zeros_like(model.b1)
    dW2 = np.zeros_like(model.W2)
    db2 = np.zeros_like(model.b2)
    m = len(X)
    cost = 0.0
    for x,y in zip(X,Y):
        a2, (z1,a1,z2, _) = model.predict(x)
        if y == 1:
            cost -= np.log(a2)
        else:
            cost -= np.log(1-a2)

        diff = a2-y
        # layer 2
        # db2
        db2 += diff

        # dw2 - todo: remove for-loops
        for i in range(model.num_hiddens):
            dW2[i] += a1[i]*diff
        #layer 1
        # db1 - todo: remove for-loops
        for i in range(model.num_hiddens):
            db1[i] += (1-a1[i]**2)*model.W2[i]*diff
        # db2 - todo: remove for-loops
        for i in range(model.num_hiddens):
            for j in range(model.num_input_features):
                dW1[i,j] += x[j]*(1-a1[i]**2)*model.W2[i]*diff

    cost /= m
    model.W1 -= lr * dW1/m
    model.b1 -= lr * db1/m
    model.W2 -= lr * dW2/m
    model.b2 -= lr * db2/m

    return cost
```