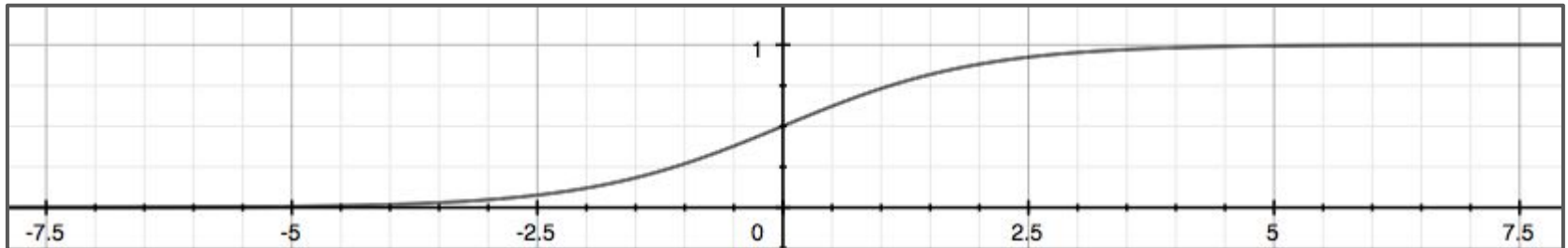# Logistic Regression

# Fade-In

- Machine Learning concepts and Training set
- Training set type : supervised Learning ( for prediction ) / unsupervised ( for clustering)
- Linear Regression (continuous value prediction)
- Logistic Regression (discrete value prediction ) - Classification
- Neural Network

# Logistic Regression (discrete value prediction ) - Classification

Classification:   y  =  0  or  1

$h_\theta(x)$ can be > 1 or < 0

# Reminder: Hypothesis for Linear Regression

$$h_\theta(x) = \theta_0 + \theta_1 x$$

$$h_\theta(x) = \theta^T x = \theta_0 x_0 + \theta_1 x_1 + \theta_2 x_2 + \cdots + \theta_n x_n$$

# Hypothesis representation

**Logistic Regression Model**
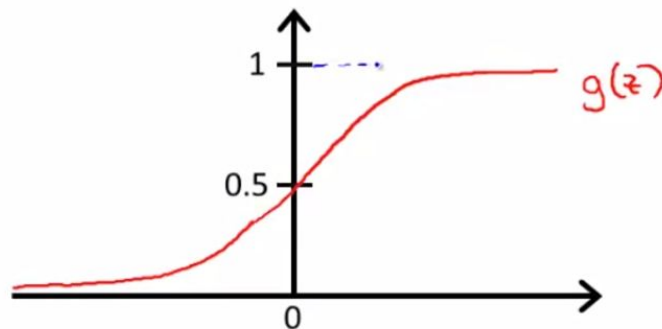
Want $0 \le h_\theta(x) \le 1$

$$h_\theta(x) = g(\theta^T x)$$

$$h_\theta(x) = \frac{1}{1 + e^{-\theta^T x}}$$
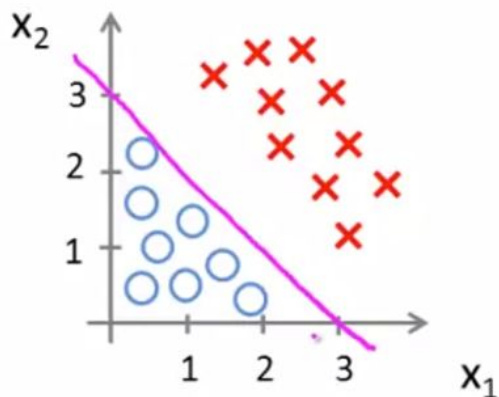
$$g(z) = \frac{1}{1 + e^{-z}}$$

$\theta^T x$

→ Sigmoid function

↳ Logistic function

$g(z)$

1

0.5

0

**Decision Boundary**



$$\theta = \begin{bmatrix} -3 \\ 1 \\ 1 \end{bmatrix} \leftarrow$$

$$\rightarrow h_\theta(x) = g(\theta_0 + \theta_1 x_1 + \theta_2 x_2)$$

$$\overset{\|}{-3} \qquad \overset{\|}{1} \qquad \overset{\|}{1}$$

$x_1, x_2$

Predict "$y = 1$" if $-3 + x_1 + x_2 \geq 0$

$\theta^T x$

$\rightarrow x_1 + x_2 \geq 3$

$x_1 + x_2 = 3$

# Cost function

$$cost(W, b) = \frac{1}{m} \sum_{i=1}^{m} (H(x^{(i)}) - y^{(i)})^2 \quad \text{when} \quad H(x) = Wx + b$$

$$cost(W, b) = \frac{1}{m} \sum_{i=1}^{m} (H(x^{(i)}) - y^{(i)})^2$$

$0 < \sim < 1$

$$H(x) = Wx + b$$

$$H(X) = \frac{1}{1 + e^{-W^T X}}$$

# New cost function for logistic

$$cost(W) = \frac{1}{m} \sum c(H(x), y)$$

$$c(H(x), y) = \begin{cases} -log(H(x)) & : y = 1 \\ -log(1 - H(x)) & : y = 0 \end{cases}$$

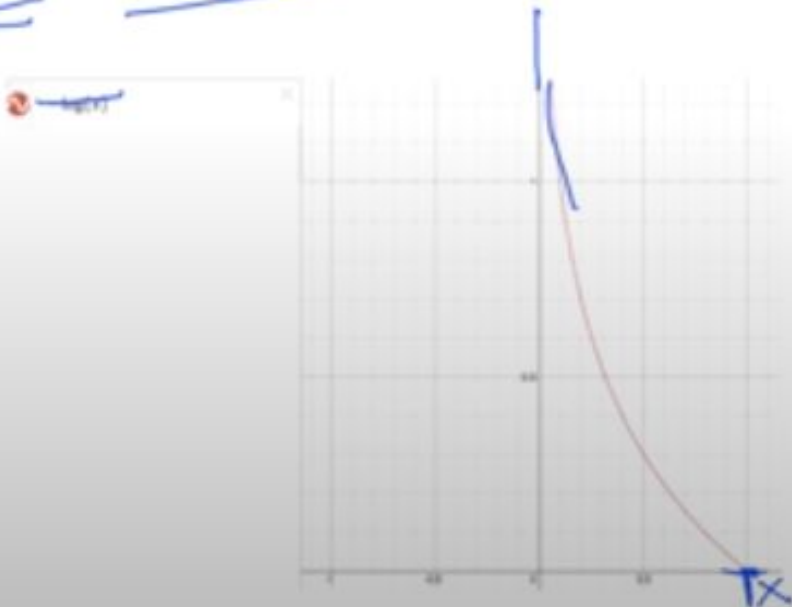$$c(H(x), y) = \begin{cases} -log(H(x)) & : y = 1 \\ -log(1 - H(x)) & : y = 0 \end{cases}$$

~~Cost~~ $\boxed{y = 1}$

$H(x) = 1 \rightarrow cost(1) = 0$

$H(x) = 0 \rightarrow cost = \infty$

Cost
$\shortparallel$

$g(2) = -log(2)$

$$c(H(x), y) = \begin{cases} -log(H(x)) & : y = 1 \\ -log(1 - H(x)) & : y = 0 \end{cases}$$

$y = 0$

$H(x) = 0$, $cost = 0$

$H(x) = 1$, $cost = \infty \uparrow$

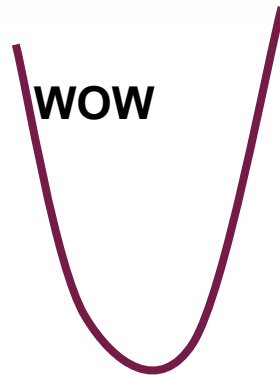$-\log(1-z)$

$\infty$

# Cost function of Logistic Regression

$$J(\theta) = \frac{1}{m} \sum_{i=1}^{m} \text{Cost}(h_\theta(x^{(i)}), y^{(i)})$$

$$\text{Cost}(h_\theta(x), y) = \begin{cases} -\log(h_\theta(x)) & \text{if } y = 1 \\ -\log(1 - h_\theta(x)) & \text{if } y = 0 \end{cases}$$

Note: $y = 0$ or $1$ always

# Simplified cost function

$$J(\theta) = \frac{1}{m} \sum_{i=1}^{m} \text{Cost}(h_\theta(x^{(i)}), y^{(i)})$$

$$= -\frac{1}{m} \left[ \sum_{i=1}^{m} y^{(i)} \log h_\theta(x^{(i)}) + (1 - y^{(i)}) \log (1 - h_\theta(x^{(i)})) \right]$$

$J(\theta)$

**WOW**

# Logistic Regression : Gradient descent

$$W := W - \alpha \frac{\partial}{\partial W} cost(W)$$

Same method -  Gradient descent of Linear Regression

# Multiclass Classification

**Multiclass classification**

Email foldering/tagging: Work, Friends, Family, Hobby

$$y=1 \quad y=2 \quad y=3 \quad y=4$$

Medical diagrams: Not ill, Cold, Flu

$$y=1 \quad 2 \quad 3$$

Weather: Sunny, Cloudy, Rain, Snow

$$y=1 \quad 2 \quad 3 \quad 4$$

$$0 \quad 1 \quad 2 \quad 3$$

**One-vs-all (one-vs-rest):**



Class 1: △
Class 2: □
Class 3: ✗

$$h_\theta^{(i)}(x) = P(y = i | x; \theta) \quad (i = 1, 2, 3)$$

$h_\theta^{(1)}(x) \quad P(y=1|x;\theta)$

$h_\theta^{(2)}(x)$

$h_\theta^{(3)}(x)$

# Multiclass Classification (cont'd)

**One-vs-all**

Train a logistic regression classifier $h_\theta^{(i)}(x)$ for each class $i$ to predict the probability that $y = i$.

On a new input $x$, to make a prediction, pick the class $i$ that maximizes

$$\max_i h_\theta^{(i)}(x)$$

# Solving the problem of overfitting

Example: Linear regression (housing prices)



$\rightarrow \theta_0 + \theta_1 x$

"Underfit"  "High bias"

$\rightarrow \theta_0 + \theta_1 x + \theta_2 x^2$

"Just right"

$\rightarrow \theta_0 + \theta_1 x + \theta_2 x^2 + \theta_3 x^3 + \theta_4 x^4$
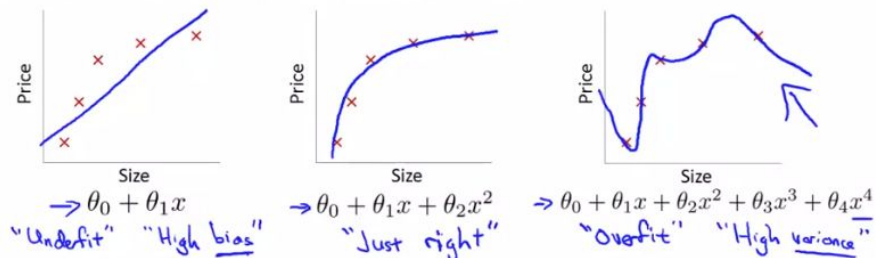
"Overfit"  "High variance"

**Overfitting:** If we have too many features, the learned hypothesis may fit the training set very well ($J(\theta) = \frac{1}{2m} \sum_{i=1}^{m} (h_\theta(x^{(i)}) - y^{(i)})^2 \approx 0$), but fail to generalize to new examples (predict prices on new examples).

Example: Logistic regression



$\rightarrow h_\theta(x) = g(\theta_0 + \theta_1 x_1 + \theta_2 x_2)$

( $g$ = sigmoid function)

"Underfit"

$g(\theta_0 + \theta_1 x_1 + \theta_2 x_2 + \theta_3 x_1^2 + \theta_4 x_2^2 + \theta_5 x_1 x_2)$

$g(\theta_0 + \theta_1 x_1 + \theta_2 x_1^2 + \theta_3 x_1^2 x_2 + \theta_4 x_1^2 x_2^2 + \theta_5 x_1^2 x_2^3 + \theta_6 x_1^3 x_2 + \dots)$

"Overfit"

# Solving the problem of overfitting (cont'd)

**Intuition**



$$\theta_0 + \theta_1 x + \theta_2 x^2$$

$$\theta_0 + \theta_1 x + \theta_2 x^2 + \theta_3 x^3 + \theta_4 x^4$$

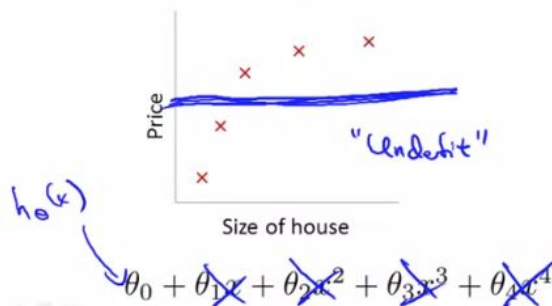Suppose we penalize and make $\theta_3, \theta_4$ really small.

$$\to \min_{\theta} \frac{1}{2m} \sum_{i=1}^{m} (h_\theta(x^{(i)}) - y^{(i)})^2 + 1000 \, \theta_3^2 + 1000 \, \theta_4^2$$

$$\theta_3 \approx 0 \qquad \theta_4 \approx 0$$

In regularized linear regression, we choose $\theta$ to minimize

$$J(\theta) = \frac{1}{2m} \left[ \sum_{i=1}^{m} (h_\theta(x^{(i)}) - y^{(i)})^2 + \lambda \sum_{j=1}^{n} \theta_j^2 \right]$$

What if $\lambda$ is set to an extremely large value (perhaps for too large for our problem, say $\lambda = 10^{10}$)?



"Underfit"

$$h_\theta(x)$$
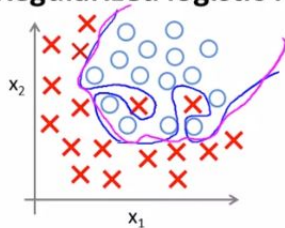$$\theta_0 + \theta_1 x + \theta_2 x^2 + \theta_3 x^3 + \theta_4 x^4$$

$$\theta_1, \theta_2, \theta_3, \theta_4$$
$$\theta_1 \approx 0, \ \theta_2 \approx 0$$
$$\theta_3 \approx 0, \ \theta_4 \approx 0$$

$$\boxed{h_\theta(x) = \theta_0}$$

# Solving the problem of overfitting (cont'd)

**Regularized logistic regression.**



$$h_\theta(x) = g(\theta_0 + \theta_1 x_1 + \theta_2 x_1^2$$
$$+ \theta_3 x_1^2 x_2 + \theta_4 x_1^2 x_2^2$$
$$+ \theta_5 x_1^2 x_2^3 + \dots)$$

Cost function:

$$J(\theta) = -\left[\frac{1}{m}\sum_{i=1}^{m} y^{(i)} \log h_\theta(x^{(i)}) + (1 - y^{(i)}) \log(1 - h_\theta(x^{(i)}))\right]$$
$$+ \frac{\lambda}{2m}\sum_{j=1}^{n} \theta_j^2 \qquad \boxed{\theta_1, \theta_2, \dots, \theta_n}$$

**Advanced optimization**

fminunc (@ costFunction)    $\theta = \begin{bmatrix} \theta_0 \\ \theta_1 \\ \vdots \\ \theta_n \end{bmatrix}$  — theta(1) ←
                                                   — theta(2)
                                                   theta(n+1)

```
function [jVal, gradient] = costFunction(theta)
```

jVal = [ code to compute $J(\theta)$ ] ;

$$J(\theta) = \left[-\frac{1}{m}\sum_{i=1}^{m} y^{(i)} \log(h_\theta(x^{(i)})) + (1 - y^{(i)}) \log 1 - h_\theta(x^{(i)})\right] + \frac{\lambda}{2m}\sum_{j=1}^{n} \theta_j^2$$

gradient(1) = [code to compute $\frac{\partial}{\partial \theta_0} J(\theta)$ ] ;

$$\frac{1}{m}\sum_{i=1}^{m}(h_\theta(x^{(i)}) - y^{(i)})x_0^{(i)} \leftarrow$$

gradient(2) = [code to compute $\frac{\partial}{\partial \theta_1} J(\theta)$ ] ;

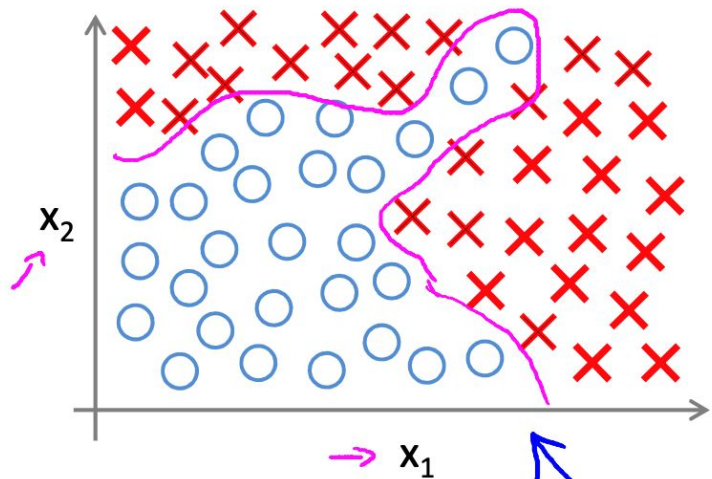$$\frac{1}{m}\sum_{i=1}^{m}(h_\theta(x^{(i)}) - y^{(i)})x_1^{(i)} + \frac{\lambda}{m}\theta_1 \leftarrow$$

gradient(3) = [code to compute $\frac{\partial}{\partial \theta_2} J(\theta)$ ] ;

$$\vdots \qquad \left(\frac{1}{m}\sum_{i=1}^{m}(h_\theta(x^{(i)}) - y^{(i)})x_2^{(i)}\right) + \frac{\lambda}{m}\theta_2$$

gradient(n+1) = [code to compute $\frac{\partial}{\partial \theta_n} J(\theta)$ ] ;

# Neural Networks

# Non-linear Classification



$$g(\theta_0 + \theta_1 x_1 + \theta_2 x_2$$
$$+\theta_3 x_1 x_2 + \theta_4 x_1^2 x_2$$
$$+\theta_5 x_1^3 x_2 + \theta_6 x_1 x_2^2 + \dots)$$

$x_2$

$x_1$

$x_1 = \text{size}$
$x_2 = \text{\# bedrooms}$
$x_3 = \text{\# floors}$
$x_4 = \text{age}$
$\dots$
$x_{100} -$

$n = 100$

$\rightarrow x_1^2, x_1 x_2, x_1 x_3, x_1 x_4 \dots x_1 x_{100}$
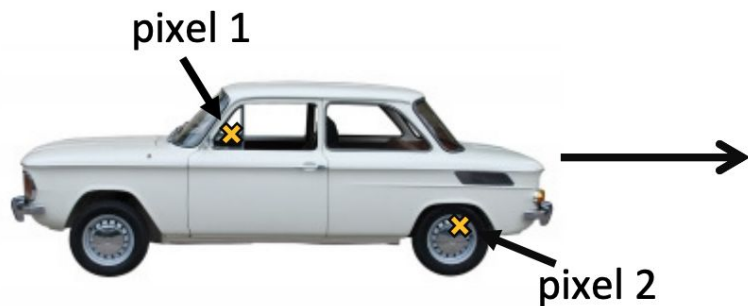$x_2^2, x_2 x_3 \dots$

$\sim 5000 \text{ feature}$ $\qquad O(n^2)$

$\rightarrow x_1^2, x_2^2, x_3^2, \dots, x_{100}^2 \qquad \sim \frac{n^2}{2}$

$\rightarrow x_1 x_2 x_3, x_1^2 x_2, x_{10} x_{11} x_{17}, \dots$
$O(n^3)$

$170,000$

pixel 1

pixel 2

Learning Algorithm

50 x 50 pixel images → 2500 pixels

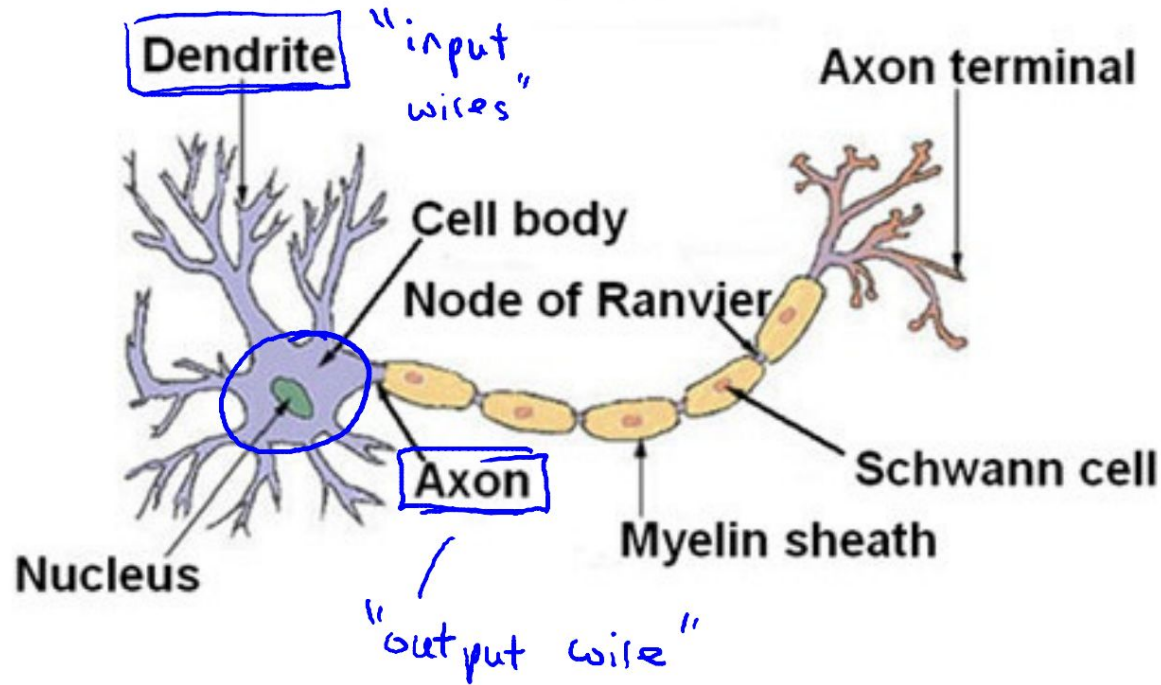$n = 2500$  (7500 if RGB)

pixel 2

pixel 1

+ Cars

– "Non"-Cars

0-255

$$x = \begin{bmatrix} \text{pixel 1 intensity} \\ \text{pixel 2 intensity} \\ \vdots \\ \text{pixel 2500 intensity} \end{bmatrix}$$

Quadratic features ($x_i \times x_j$): ≈3 million features

Andrew Ng

# Neural Networks

# Neuron in the brain

# Neuron model: Logistic unit



$$x = \begin{bmatrix} x_0 \\ x_1 \\ x_2 \\ x_3 \end{bmatrix} \quad \theta = \begin{bmatrix} \theta_0 \\ \theta_1 \\ \theta_2 \\ \theta_3 \end{bmatrix}$$

"bias unit"

$x_0 = 1$

"output"

$h_\theta(x)$

$h_\theta(x) = \dfrac{1}{1 + e^{-\Theta^T x}}$

"input wires"

"weights"
(parameters)

Sigmoid (logistic) activation function.

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

# Neural Network



$x_0$

$a_0^{(2)}$    "bias unit

$x_1$

$x_2$

$x_3$

$a_1^{(2)}$

$a_2^{(2)}$

$a_3^{(2)}$

$h_\Theta(x)$

Layer 1    $\times$    Layer 2    Layer 3

Input layer    Hidden layer    Output layer    $y$

Andrew Ng

# Other network architectures



$$h_\Theta(x)$$

Layer 1      Layer 2      Layer 3      Layer 4

Input      Hidden layer      Output

# Applications

# Applications

How a neural network can compute a compute non-linear function of the input

Non-linear 분류의 대표적인 예:
## XOR / XNOR

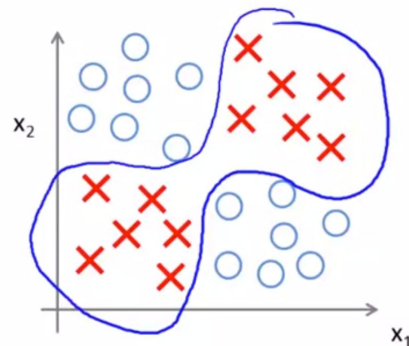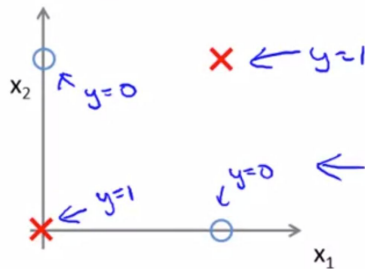| XOR | 0 | 1 |
|-----|---|---|
| 0 | 0 | 1 |
| 1 | 1 | 0 |

XOR 게이트는 exclusive OR 게이트의 줄임말로 한국어로는 상호 배제적인 OR 게이트이다. 상호 배제적 OR 게이트란 이름의 의미는 OR 게이트와 동일하게 작동하지만 입력값이 동일한 경우에는 1을 출력하지 않는다는 의미이다. 입력값이 서로 다르면 1을 출력하고, 같으면 0을 출력한다.

| XNOR | 0 | 1 |
|------|---|---|
| 0 | 1 | 0 |
| 1 | 0 | 1 |

XNOR 게이트(XNOR gate 또는 EXNOR, ENOR, NXOR, XAND gate)는 XOR 게이트 뒤에 NOT 게이트를 붙여 출력값을 반대로 만들어놓은 것이다. 입력값이 서로 같으면 1을 출력하고, 다르면 0을 출력한다. 그래서 비교 게이트나 일치 확인 게이트라고도 불린다.

논리회로: 출처 나무위키



**Non-linear classification example: XOR/XNOR**

$\rightarrow$ $x_1$, $x_2$ are binary (0 or 1).

$$y = x_1 \text{ XOR } x_2$$
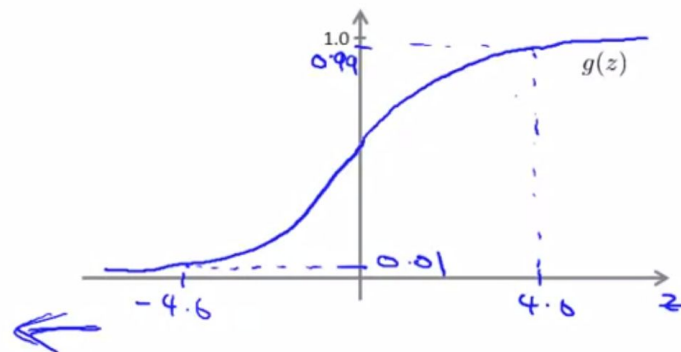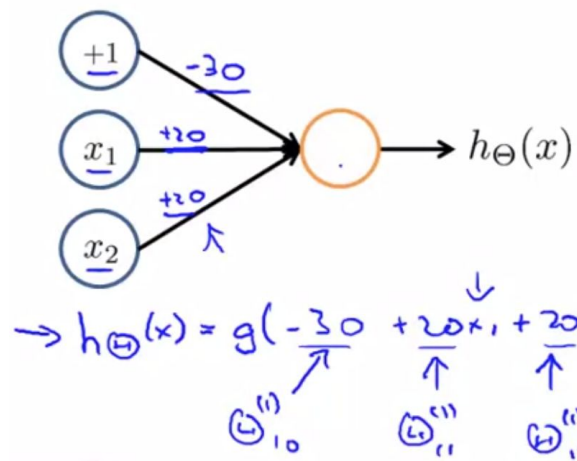$$x_1 \text{ XNOR } x_2$$
$$\text{NOT } (x_1 \text{ XOR } x_2)$$

# Simple example: And Function

And function 계산
(1 X -30
x1 X 20
x2 X 20)
에다가 Sigmoid
function 적용



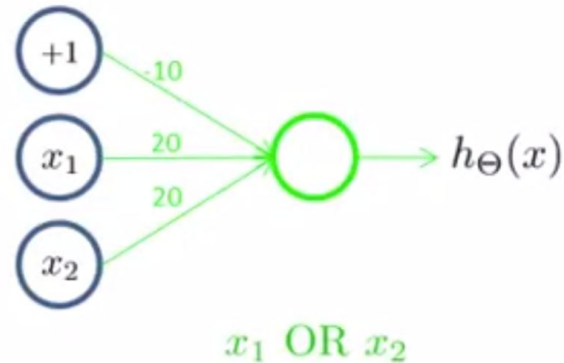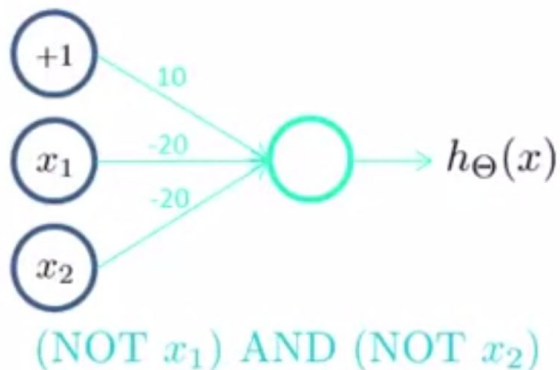## Simple example: AND

$\rightarrow x_1, x_2 \in \{0, 1\}$

$\rightarrow y = x_1 \text{ AND } x_2$

$+1$ $\xrightarrow{-30}$

$x_1$ $\xrightarrow{+20}$ $\bigcirc$ $\rightarrow h_\Theta(x)$

$x_2$ $\xrightarrow{+20}$

$\rightarrow h_\Theta(x) = g(-30 + 20x_1 + 20x_2) \leftarrow$

$\Theta^{(1)}_{10}$ $\quad$ $\Theta^{(1)}_{11}$ $\quad$ $\Theta^{(1)}_{12}$

$g(z)$

1.0
0.99

0.01

-4.6 $\qquad$ 4.6 $\quad$ z

| $x_1$ | $x_2$ | $h_\Theta(x)$ |
|-------|-------|---------------|
| 0 | 0 | $g(-30) \approx 0$ |
| 0 | 1 | $g(-10) \approx 0$ |
| 1 | 0 | $g(-10) \approx 0$ |
| 1 | 1 | $g(10) \approx 1$ |

$h_\Theta(x) \approx x_1 \text{ AND } x_2$

# Putting it together: $x_1$ XNOR $x_2$



$x_1$ AND $x_2$

(NOT $x_1$) AND (NOT $x_2$)

$x_1$ OR $x_2$

x1 XNOR x2 =
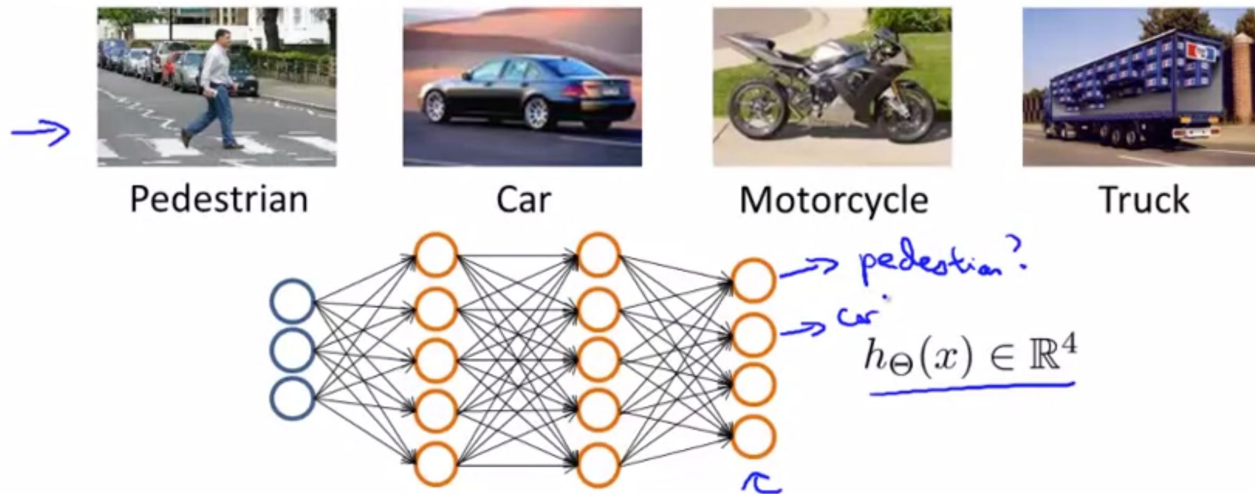
x1 AND x2
(NOT x1) AND (NOT x2)
x1 OR x2

XNOR 이나 XOR 과 같이
Linear function으로 Classification이 되지 않는 경우에도
Neural Network를 사용해 Classification이 가능함

| $x_1$ | $x_2$ | $a_1^{(2)}$ | $a_2^{(2)}$ | $h_\Theta(x)$ |
|---|---|---|---|---|
| 0 | 0 | 0 | 1 | 1 |
| 0 | 1 | 0 | 0 | 0 |
| 1 | 0 | 0 | 0 | 0 |
| 1 | 1 | 1 | 0 | 1 |

# Multiclass Classification

**Multiple output units: One-vs-all.**



Pedestrian      Car      Motorcycle      Truck

→ pedestrian?

→ car

$h_\Theta(x) \in \mathbb{R}^4$

Want $h_\Theta(x) \approx \begin{bmatrix} 1 \\ 0 \\ 0 \\ 0 \end{bmatrix}$, $h_\Theta(x) \approx \begin{bmatrix} 0 \\ 1 \\ 0 \\ 0 \end{bmatrix}$, $h_\Theta(x) \approx \begin{bmatrix} 0 \\ 0 \\ 1 \\ 0 \end{bmatrix}$, etc.

when pedestrian    when car    when motorcycle

One vs all 의 확장판쯤

Pedestrian? Yes or No
[1, 0, 0, 0]
Car? Yes or No
[0, 1, 0, 0]
Motocycle? Yes or No
[0, 0, 1, 0]
Truck? Yes or No
[0, 0, 0, 1]

다른 예시: MNIST 10 digits

```
print(Y_train[:5])
```

```
[[0. 1. 0. 0. 0. 0. 0. 0. 0. 0.]
 [1. 0. 0. 0. 0. 0. 0. 0. 0. 0.]
 [0. 1. 0. 0. 0. 0. 0. 0. 0. 0.]
 [0. 0. 0. 0. 1. 0. 0. 0. 0. 0.]
 [1. 0. 0. 0. 0. 0. 0. 0. 0. 0.]]
```