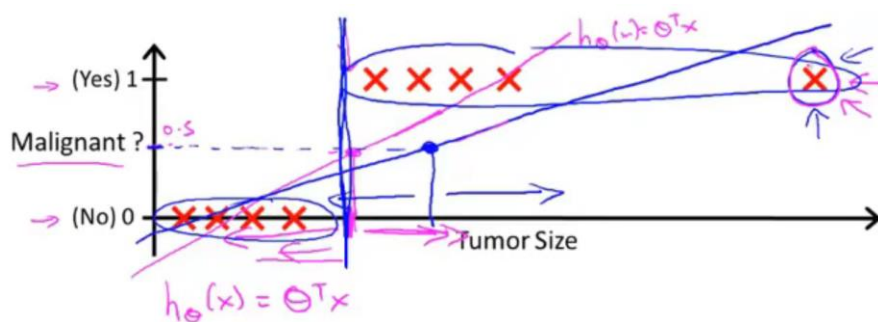# Logistic Regression

## 1 Classification

**Classification**

→ Email: Spam / Not Spam?
→ Online Transactions: Fraudulent (Yes / No)?
→ Tumor: Malignant / Benign ?

→ $y \in \{0, 1\}$
　　　→ 0: "Negative Class" (e.g., benign tumor)
　　　→ 1: "Positive Class" (e.g., malignant tumor)

→ $y \in \{0, 1, 2, 3\}$



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

→ Threshold classifier output $h_\theta(x)$ at 0.5:

　　→ If $h_\theta(x) \geq 0.5$, predict "y = 1"
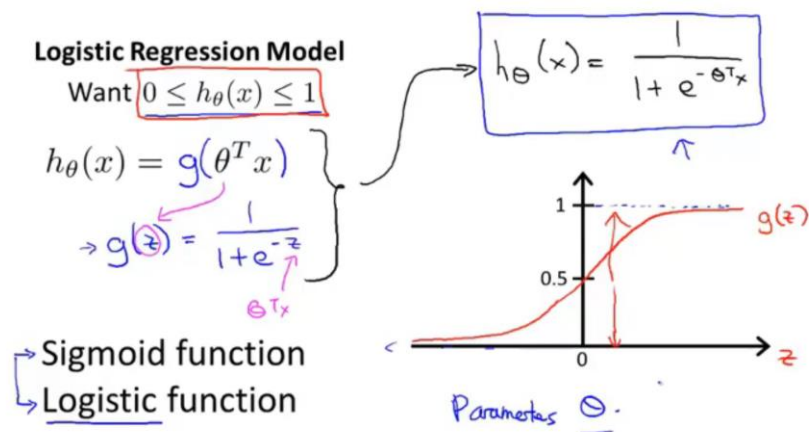
　　If $h_\theta(x) < 0.5$, predict "y = 0"

Classification:　$y \;=\; 0 \;\; or \;\; 1$

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

Logistic Regression:　$0 \leq h_\theta(x) \leq 1$

　　　Classification

# 2 Hypothesis representation

**Logistic Regression Model**

Want $0 \le h_\theta(x) \le 1$

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

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

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

$$\theta^T x$$

→ Sigmoid function
↳ Logistic function



Parameters $\theta$.

**Interpretation of Hypothesis Output**     $h_\theta(x)$

$h_\theta(x)$ = estimated probability that $y = 1$ on input x ←

Example: If $x = \begin{bmatrix} x_0 \\ x_1 \end{bmatrix} = \begin{bmatrix} 1 \\ \text{tumorSize} \end{bmatrix}$ ←

$$h_\theta(x) = 0.7 \qquad y=1$$

Tell patient that 70% chance of tumor being malignant

$$h_\theta(x) = P(y=1|x;\theta)$$

"probability that y = 1, given x, parameterized by $\theta$"

$$y = 0 \quad \text{or} \quad 1$$

$$\to P(y=0|x;\theta) + P(y=1|x;\theta) = 1$$
$$P(y=0|x;\theta) = 1 - P(y=1|x;\theta)$$

Andrew Ng

# 3 Decision boundary

**Logistic regression**

$$\to h_\theta(x) = g(\theta^T x) = P(y=1|x;\theta)$$

$$\to g(z) = \frac{1}{1+e^{-z}}$$
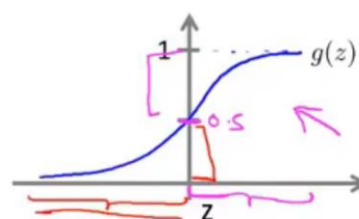
Suppose predict "$y = 1$" if $h_\theta(x) \ge 0.5$

$$\to \theta^T x \ge 0$$

predict "$y = 0$" if $h_\theta(x) < 0.5$

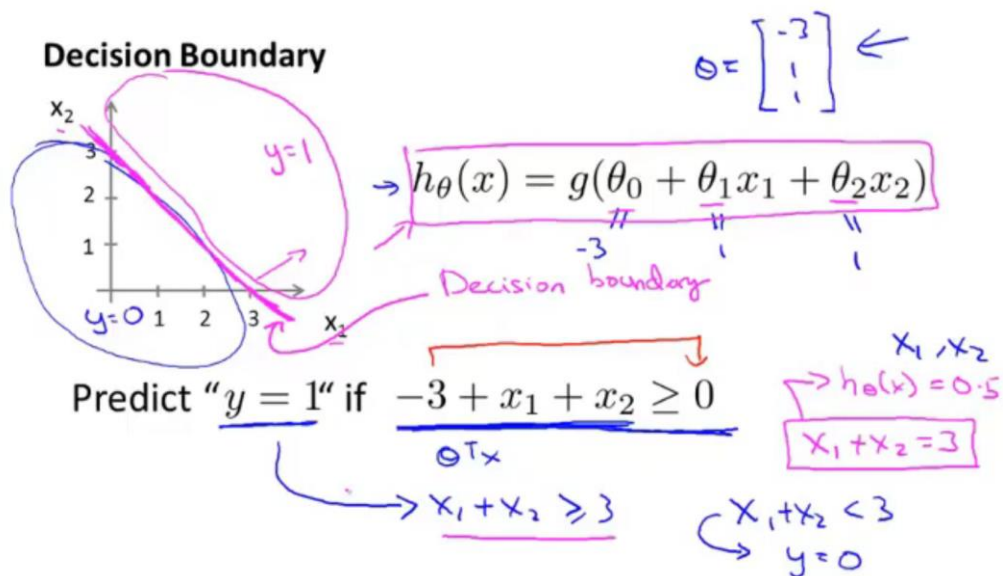$$h_\theta(x) = g(\theta^T x)$$
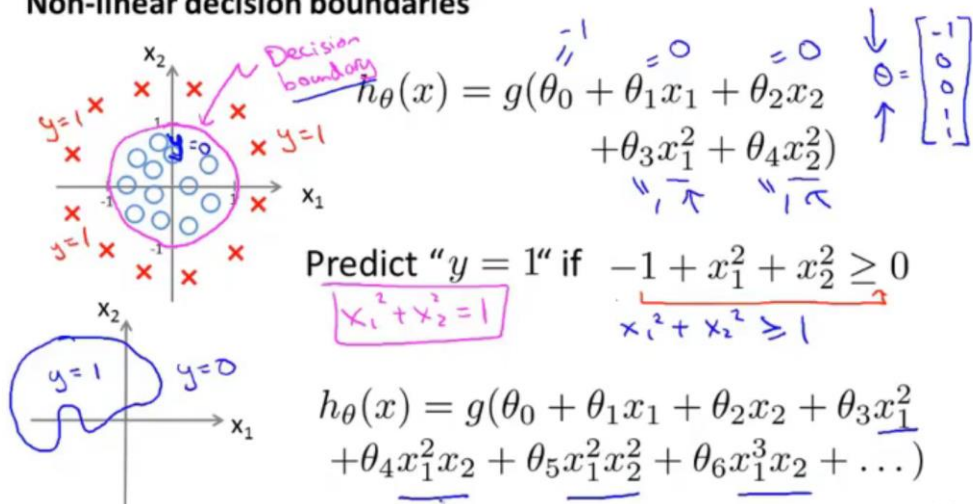$$\to \theta^T x < 0$$



$g(z) \ge 0.5$
when $z \ge 0$

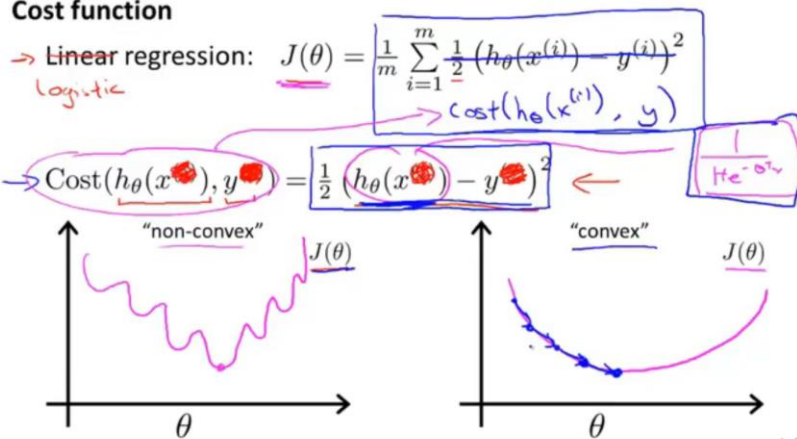$h_\theta(x) = g(\theta^T x) \ge 0.5$
whenever $\theta^T x \ge 0$

$g(z) < 0.5$

**Decision Boundary**

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

$x_2$

$y=1$

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

-3

Decision boundary

$y=0$ 1  2  3

$x_1$

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

$\theta^T x$

$\rightarrow x_1 + x_2 \geq 3$

$x_1, x_2$
$\rightarrow h_\theta(x) = 0.5$

$x_1 + x_2 = 3$

$x_1 + x_2 < 3$
$\rightarrow y = 0$

---

**Non-linear decision boundaries**

$x_2$

Decision boundary

$y=1$

$y=1$

$y=1$

$x_1$

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

$\theta = \begin{bmatrix} -1 \\ 0 \\ 0 \\ 1 \\ 1 \end{bmatrix}$

Predict "$y = 1$" if $-1 + x_1^2 + x_2^2 \geq 0$

$x_1^2 + x_2^2 = 1$

$x_1^2 + x_2^2 \geq 1$

$x_2$

$y=1$   $y=0$

$x_1$

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

Andrew Ng

# 4 Cost function

**Cost function**

$\rightarrow$ ~~Linear~~ regression:  $J(\theta) = \dfrac{1}{m} \sum\limits_{i=1}^{m} \dfrac{1}{2} \left( h_\theta(x^{(i)}) - y^{(i)} \right)^2$

Logistic

$\rightarrow \text{Cost}(h_\theta(x^{(i)}), y)$

$\rightarrow \text{Cost}(h_\theta(x), y) = \dfrac{1}{2} \left( h_\theta(x) - y \right)^2 \leftarrow$

$\dfrac{1}{1 + e^{-\theta^T x}}$

"non-convex"

$J(\theta)$

$\theta$

"convex"

$J(\theta)$

$\theta$

Andrew Ng

**Logistic regression cost function**

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

If y = 0

$-\log(1-z)$

# 5 Simplified cost function and gradient descent

**Logistic regression cost function**

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

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

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

$$\text{Cost}(h_\theta(x), y) = -y \log(h_\theta(x)) - (1-y) \log(1 - h_\theta(x))$$

$\overset{y=1}{=0} \qquad \overset{=0}{=1}$

If $y=1$: $\text{Cost}(h_\theta(x), y) = -\log h_\theta(x)$

If $y=0$: $\text{Cost}(h_\theta(x), y) = -\log(1 - h_\theta(x))$

**Logistic regression cost function**

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

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

To fit parameters $\theta$:

$$\min_\theta J(\theta) \qquad \text{Get } \theta$$

To make a prediction given new $x$:

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

$p(y=1 \mid x; \theta)$

**Gradient Descent**

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

Want $\min_\theta J(\theta)$:

$$\Theta = \begin{bmatrix} \theta_0 \\ \theta_1 \\ \theta_2 \\ \vdots \\ \theta_n \end{bmatrix} \leftarrow \quad \text{for } i=0 \text{ to } n$$

Repeat $\{$

$$\theta_j := \theta_j - \alpha \sum_{i=1}^{m} \left(h_\theta(x^{(i)}) - y^{(i)}\right)x_j^{(i)}$$

(simultaneously update all $\theta_j$)

$\}$

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

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

Algorithm looks identical to linear regression!

Andrew Ng

# 6 Advanced optimization

**Optimization algorithm**

Cost function $J(\theta)$. Want $\min_\theta J(\theta)$.

Given $\theta$, we have code that can compute

- $J(\theta)$ $\leftarrow$
- $\frac{\partial}{\partial \theta_j} J(\theta)$ $\leftarrow$ (for $j = 0, 1, \ldots, n$)

Gradient descent:

Repeat $\{$

$$\theta_j := \theta_j - \alpha \frac{\partial}{\partial \theta_j} J(\theta)$$

$\}$

**Optimization algorithm**

Given $\theta$, we have code that can compute

- $J(\theta)$ $\leftarrow$
- $\frac{\partial}{\partial \theta_j} J(\theta)$ $\leftarrow$ (for $j = 0, 1, \ldots, n$)

Optimization algorithms:
- Gradient descent
- Conjugate gradient
- BFGS
- L-BFGS

Advantages:
- No need to manually pick $\alpha$
- Often faster than gradient descent.

Disadvantages:
- More complex $\leftarrow$

**Example:** $\min\limits_{\theta} J(\theta)$

$\theta = \begin{bmatrix} \theta_1 \\ \theta_2 \end{bmatrix}$  $\theta_1 = 5, \theta_2 = 5.$

$J(\theta) = (\theta_1 - 5)^2 + (\theta_2 - 5)^2$

$\frac{\partial}{\partial \theta_1} J(\theta) = 2(\theta_1 - 5)$

$\frac{\partial}{\partial \theta_2} J(\theta) = 2(\theta_2 - 5)$

```
function [jVal, gradient]
           = costFunction(theta)
jVal = (theta(1)-5)^2 + ...
         (theta(2)-5)^2;
gradient = zeros(2,1);
gradient(1) = 2*(theta(1)-5);
gradient(2) = 2*(theta(2)-5);
```

```
options = optimset('GradObj', 'on', 'MaxIter', '100');
initialTheta = zeros(2,1);
[optTheta, functionVal, exitFlag] ...
     = fminunc(@costFunction, initialTheta, options);
```

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

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

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

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

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

⋮

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

# 7 Multi-class classification: One-vs-all

**Multiclass classification**

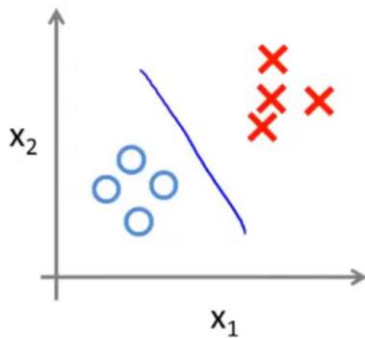Email foldering/tagging: Work, Friends, Family, Hobby

$y=1$  $y=2$  $y=3$  $y=4$

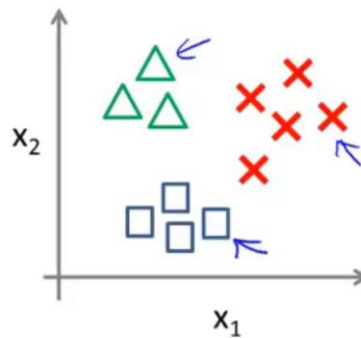Medical diagrams: Not ill, Cold, Flu

$y=1$  2  3

Weather: Sunny, Cloudy, Rain, Snow

$y=1$  2  3  4
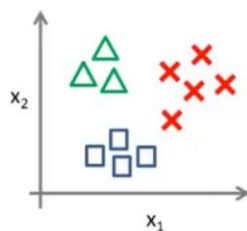
0  1  2  3

**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

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