# Logistic Regression Classification
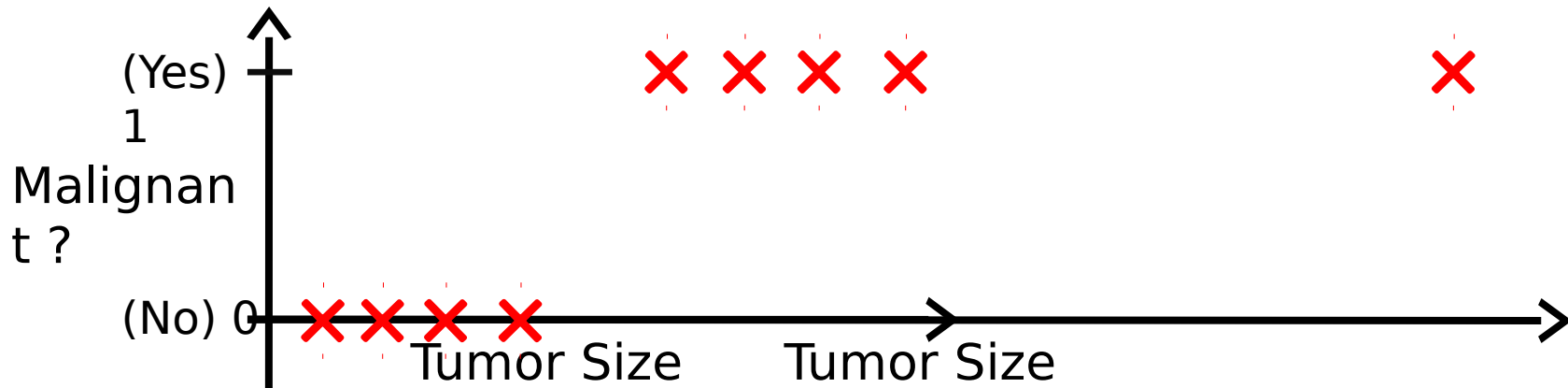
Machine Learning

# 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)

(Yes) 1

Malignant ?

(No) 0

Tumor Size

Tumor Size

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

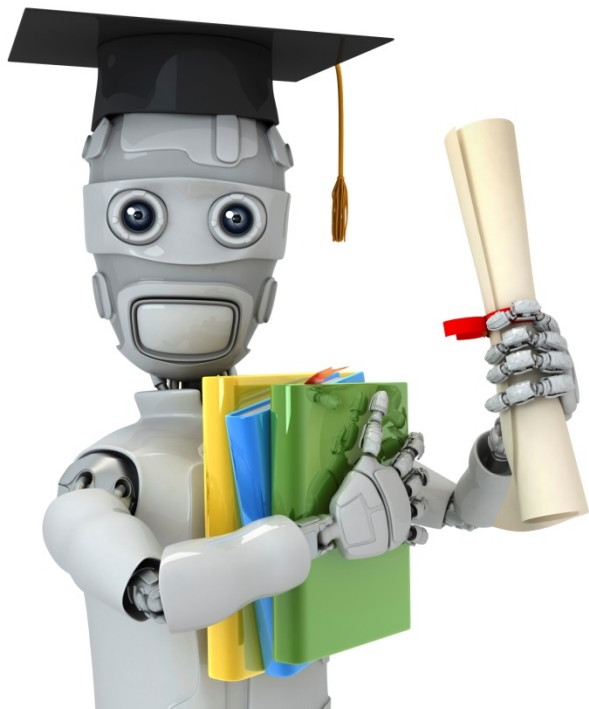If $h_\theta(x) \geq 0.5$ , predict "y = 1"

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

Classification:   y   =   0 or  1

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

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

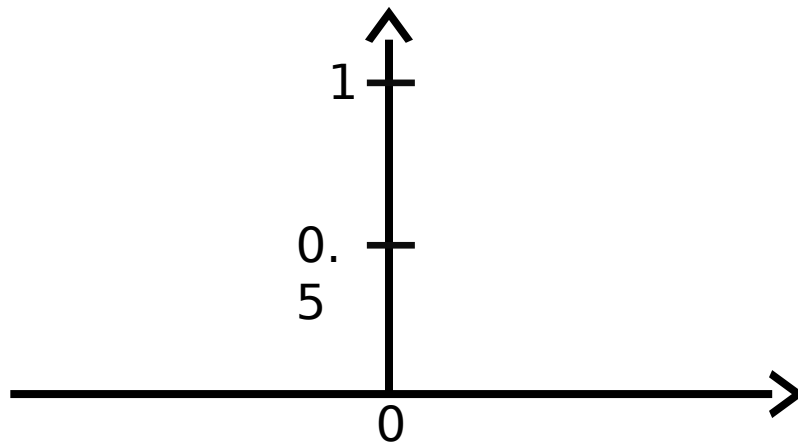Logistic Regression Hypothesis Representation

Machine Learning

# Logistic Regression Model

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

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

Sigmoid function
Logistic function

## Interpretation of Hypothesis Output

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

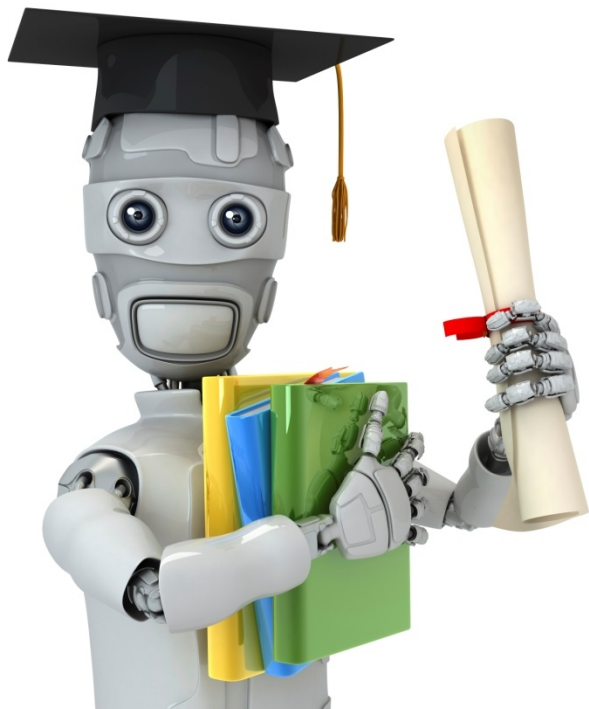Example: $x = \begin{bmatrix} x_0 \\ x_1 \end{bmatrix} = \begin{bmatrix} 1 \\ \text{tumorSize} \end{bmatrix}$

If

$$h_\theta(x) = 0.7$$

Tell patient that 70% chance of tumor being malignant

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

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

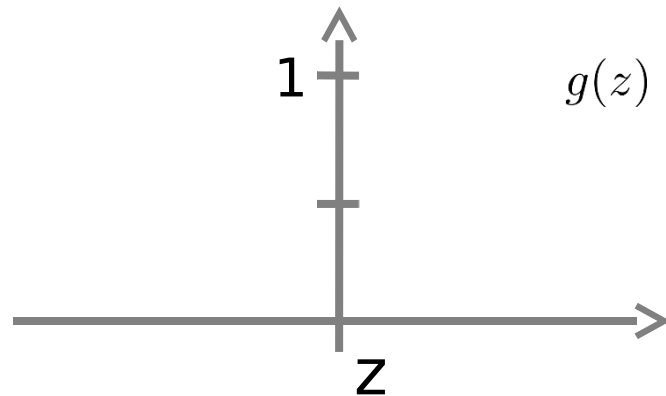$$P(y = 0|x;\theta) = 1 - P(y = 1|x;\theta)$$

# Logistic Regression

## Decision boundary

Machine Learning

# Logistic regression

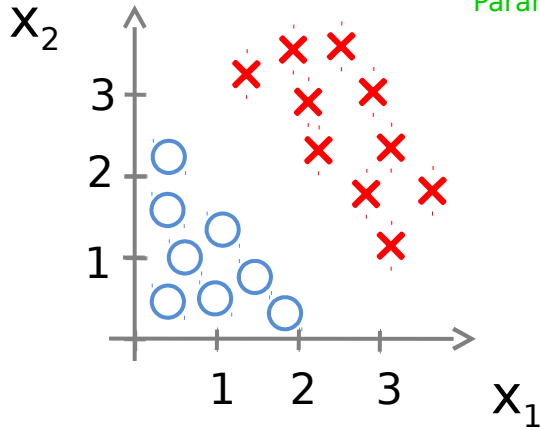$$h_\theta(x) = g(\theta^T x)$$

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



Suppose predict " $y = 1$ " if $h_\theta(x) \geq 0.5$

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

# Decision Boundary
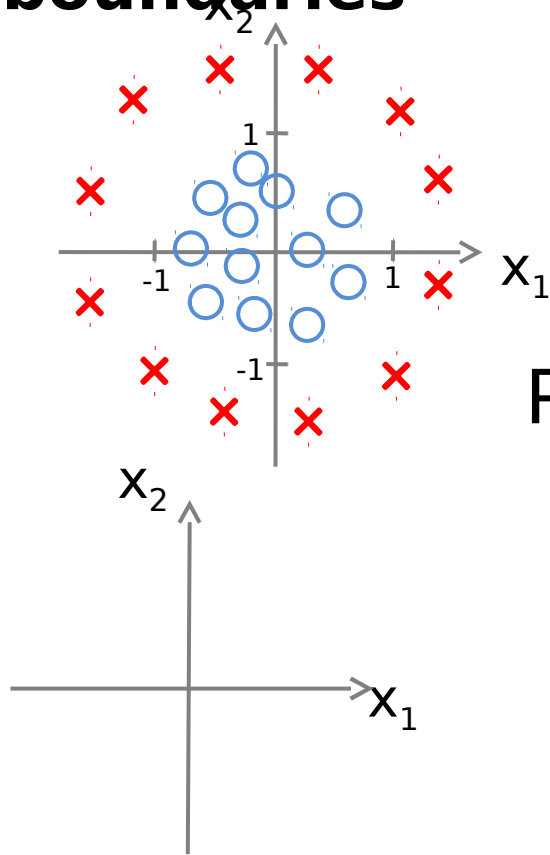
Parameter θ (θ0, θ1, θ2) defines the decision boundary not the training set. Training set may be used to find the Parameter θ



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

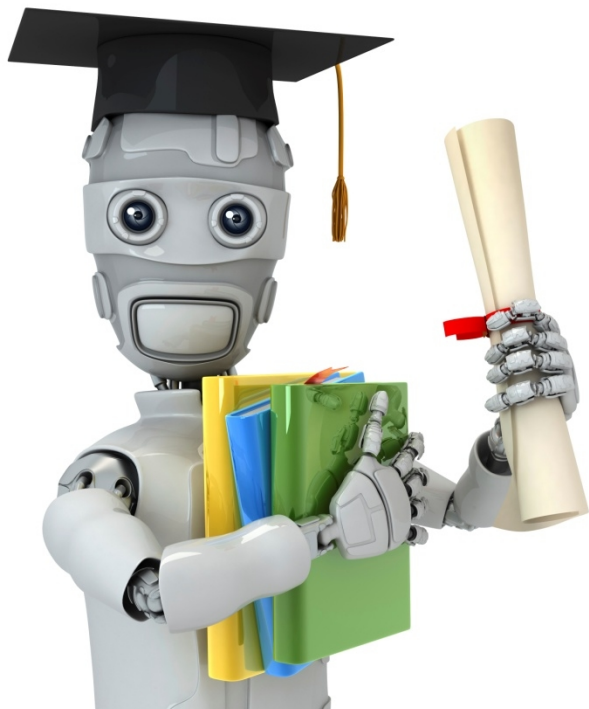Predict $y = 1$ "if" $-3 + x_1 + x_2 \geq 0$

Andrew Ng

# Non-linear decision boundaries



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

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

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

# Logistic Regression

## Cost function

Machine Learning

Training set: $\{(x^{(1)}, y^{(1)}), (x^{(2)}, y^{(2)}), \cdots, (x^{(m)}, y^{(m)})\}$

m examples
$$x \in \begin{bmatrix} x_0 \\ x_1 \\ \dots \\ x_n \end{bmatrix} \qquad x_0 = 1, y \in \{0, 1\}$$

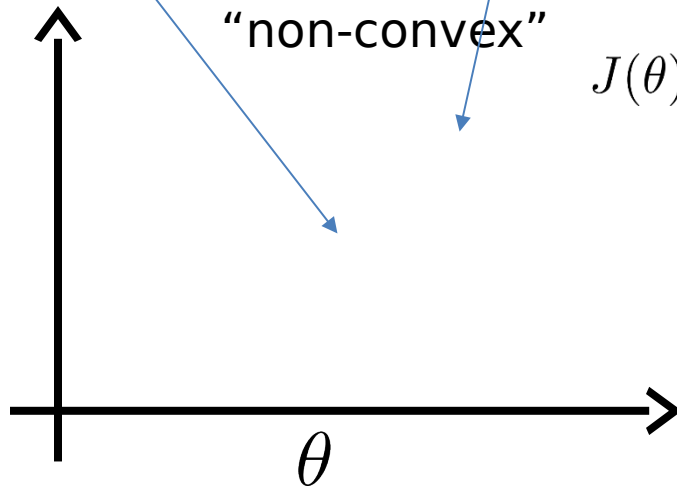$$h_\theta(x) = \frac{1}{1 + e^{-\theta^T x}}$$
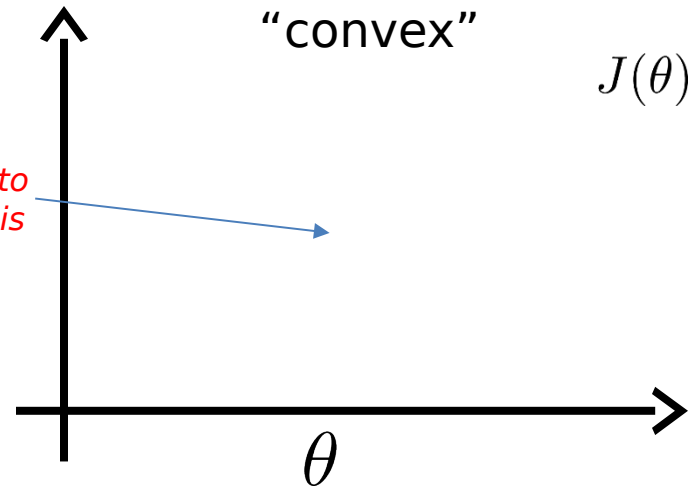
How to choose parameters $\theta$ ?

# Cost function

Linear regression: $J(\theta) = \frac{1}{m} \sum_{i=1}^{m} \frac{1}{2} \left( h_\theta(x^{(i)}) - y^{(i)} \right)^2$

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

"non-convex"

$J(\theta)$

"convex"

$J(\theta)$

$\theta$

$\theta$

Andrew Ng

# Logistic regression cost function

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

If y = 1



Cost $= 0$ if $y = 1, h_\theta(x) = 1$
But as $\quad h_\theta(x) \to 0$
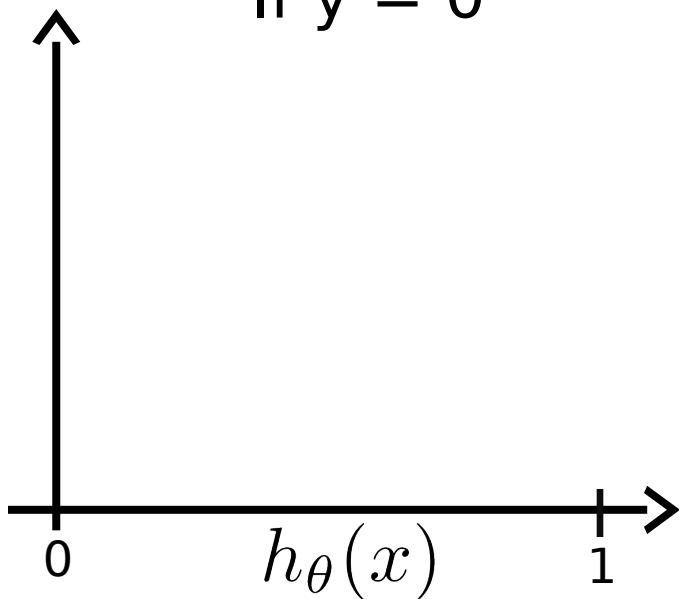$\qquad\qquad Cost \to \infty$

Captures intuition that if $h_\theta(x) = 0$, (predict $P(y = 1|x; \theta) = 0$), but $y = 1$, we'll penalize learning algorithm by a very large cost.

# Logistic regression cost function

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

If y = 0



$h_\theta(x)$

0        1

In logistic regression, the cost function for our hypothesis outputting (predicting) $h_\theta(x)$ on a training example that has label $y \in \{0, 1\}$ is:

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

Which of the following are true? Check all that apply.

☑ If $h_\theta(x) = y$, then $\mathrm{cost}(h_\theta(x), y) = 0$ (for $y = 0$ and $y = 1$).

> **Well done!**

☑ If $y = 0$, then $\mathrm{cost}(h_\theta(x), y) \to \infty$ as $h_\theta(x) \to 1$.

> **Well done!**
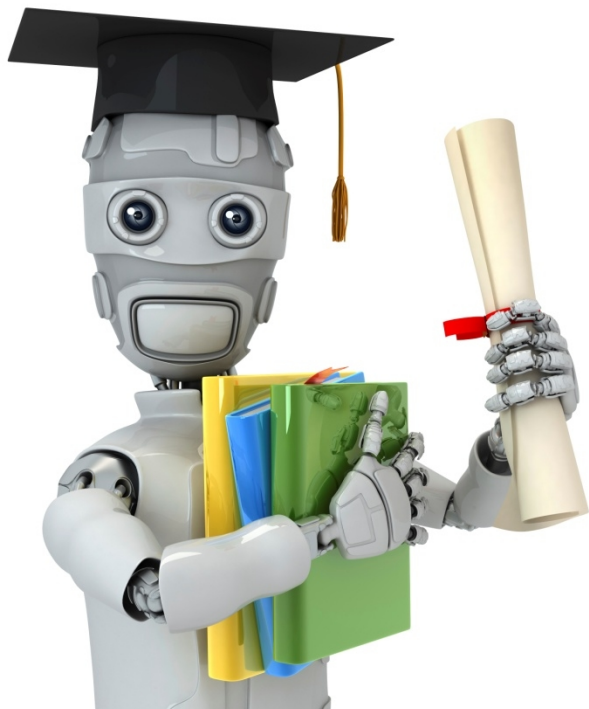
☐ If $y = 0$, then $\mathrm{cost}(h_\theta(x), y) \to \infty$ as $h_\theta(x) \to 0$.

> **Well done!**

☑ Regardless of whether $y = 0$ or $y = 1$, if $h_\theta(x) = 0.5$, then $\mathrm{cost}(h_\theta(x), y) > 0$.

> **Well done!**

# Logistic Regression

## Simplified cost function and gradient descent

Machine Learning

# 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

# Logistic regression cost function

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

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

To fit parameters $\theta$ :

$$\min_\theta J(\theta)$$

To make a prediction given new $x$ :

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

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

Repeat {

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

(simultaneously update all $\theta_j$ )

}

# Gradient Descent

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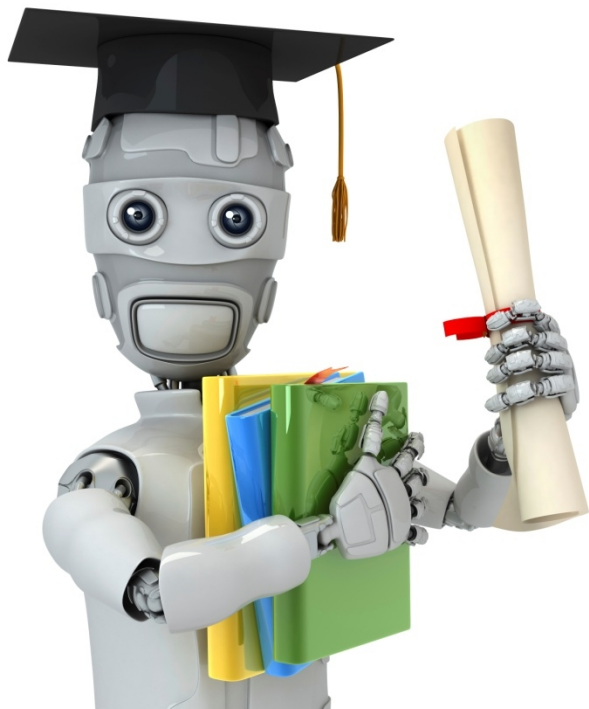
Want $\min_\theta J(\theta)$ :

Repeat $\{$

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

$\}$         (simultaneously update all $\theta_j$ )

Algorithm looks identical to linear regression!

Logistic Regression

Advanced optimization

Machine Learning

# Optimization algorithm

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

Given $\theta$, we have code that can compute
$J(\theta)$

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

# Example:

$$\theta = \begin{bmatrix} \theta_1 \\ \theta_2 \end{bmatrix}$$

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

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

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

```
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);
```

$$\textbf{theta} = \begin{bmatrix} \theta_0 \\ \theta_1 \\ \vdots \\ \theta_n \end{bmatrix}$$
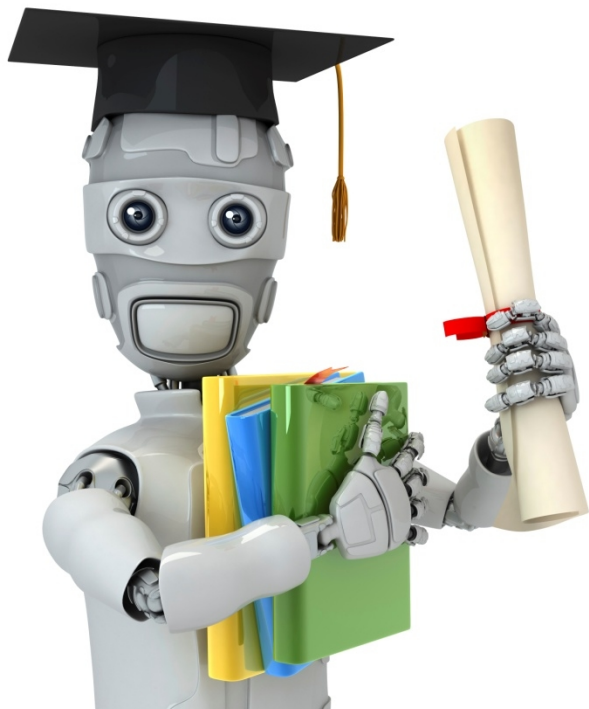
**function [jVal, gradient] = costFunction(theta)**

     **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 $\frac{\partial}{\partial \theta_n} J(\theta)$ **];**

Andrew Ng

# Logistic Regression

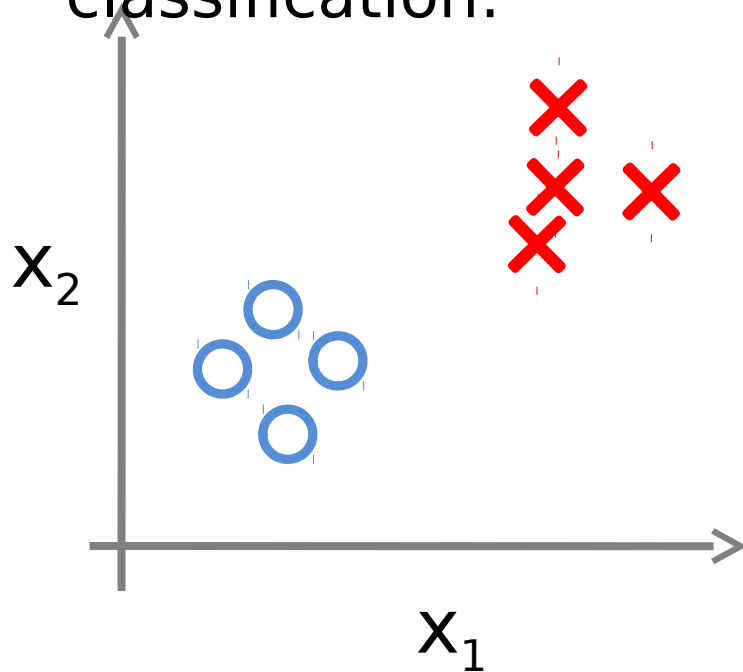## Multi-class classification: One-vs-all

Machine Learning

## Multiclass classification
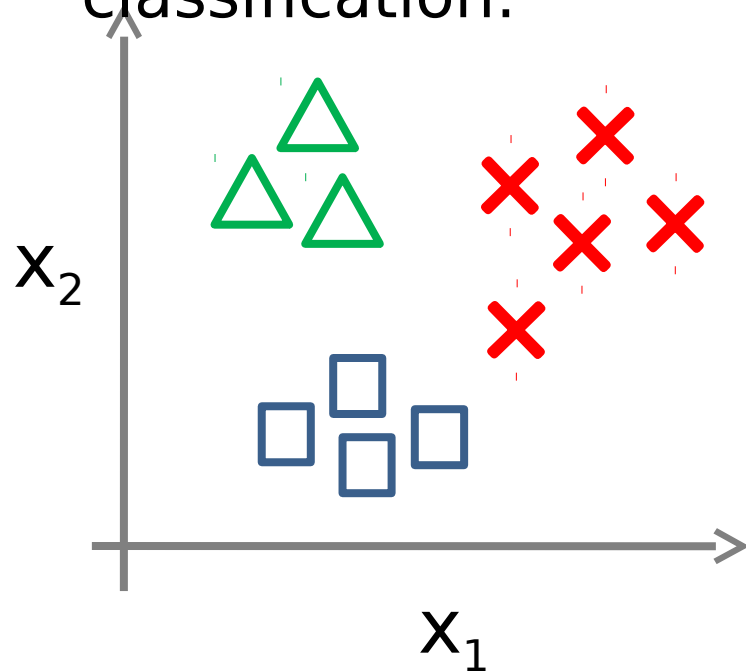
Email foldering/tagging: Work, Friends, Family, Hobby

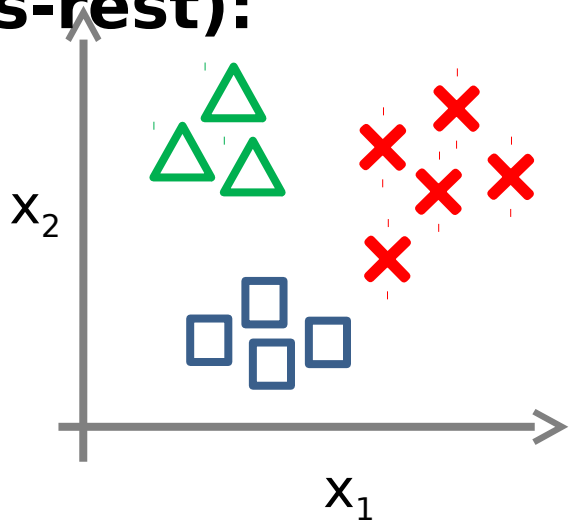Medical diagrams: Not ill, Cold, Flu

Weather: Sunny, Cloudy, Rain, Snow

Binary classification:

$x_2$

$x_1$

Multi-class classification:

$x_2$

$x_1$

Andrew Ng

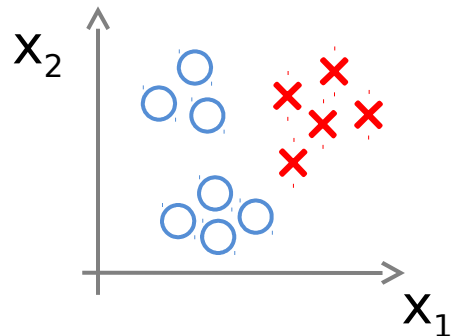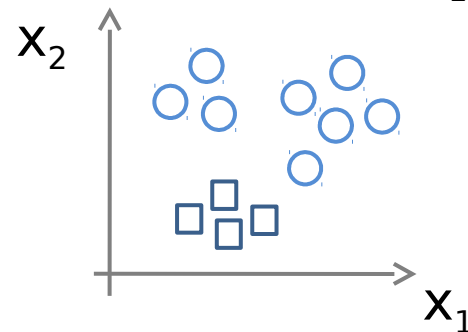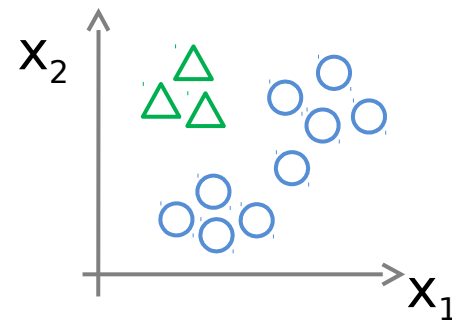# One-vs-all (one-vs-rest):



Class 1: △
Class 2: □
Class 3: ✖

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

Andrew Ng

## 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

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