# Linear Regression

# Simple Linear Regression

Training Set

Learning Algorithm

Size of house $\rightarrow$ $h$ $\rightarrow$ Estimated price

$x$ (estimated value of $y$)

hypothesis

$h$ maps from $x$'s to $y$'s.

**How do we represent $h$ ?**

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

Shorthand: $h(x)$

$h(x) = \theta_0 + \theta_1 x$

Linear regression with one variable.
Univariate linear regression.

Andrew Ng

# Cost Function



y

x

$(x^{(i)}, y^{(i)})$

Idea: Choose $\theta_0, \theta_1$ so that $h_\theta(x)$ is close to $y$ for our training examples $(x, y)$

x, y

#training examples

$$\text{minimize}_{\theta_0, \theta_1} \frac{1}{2m} \sum_{i=1}^{m} \left( h_\theta(x^{(i)}) - y^{(i)} \right)^2$$

$$h_\theta(x^{(i)}) = \theta_0 + \theta_1 x^{(i)}$$

$$J(\theta_0, \theta_1) = \frac{1}{2m} \sum_{i=1}^{m} \left( h_\theta(x^{(i)}) - y^{(i)} \right)^2$$

minimize $J(\theta_0, \theta_1)$
$\theta_0, \theta_1$

Cost function

Squared error function

$\theta_0, \theta_1$

Andrew Ng

**Used for linear regression**

# Cost Function - Intuition I

**Hypothesis:**

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

**Parameters:**

$$\theta_0, \theta_1$$

**Cost Function:**

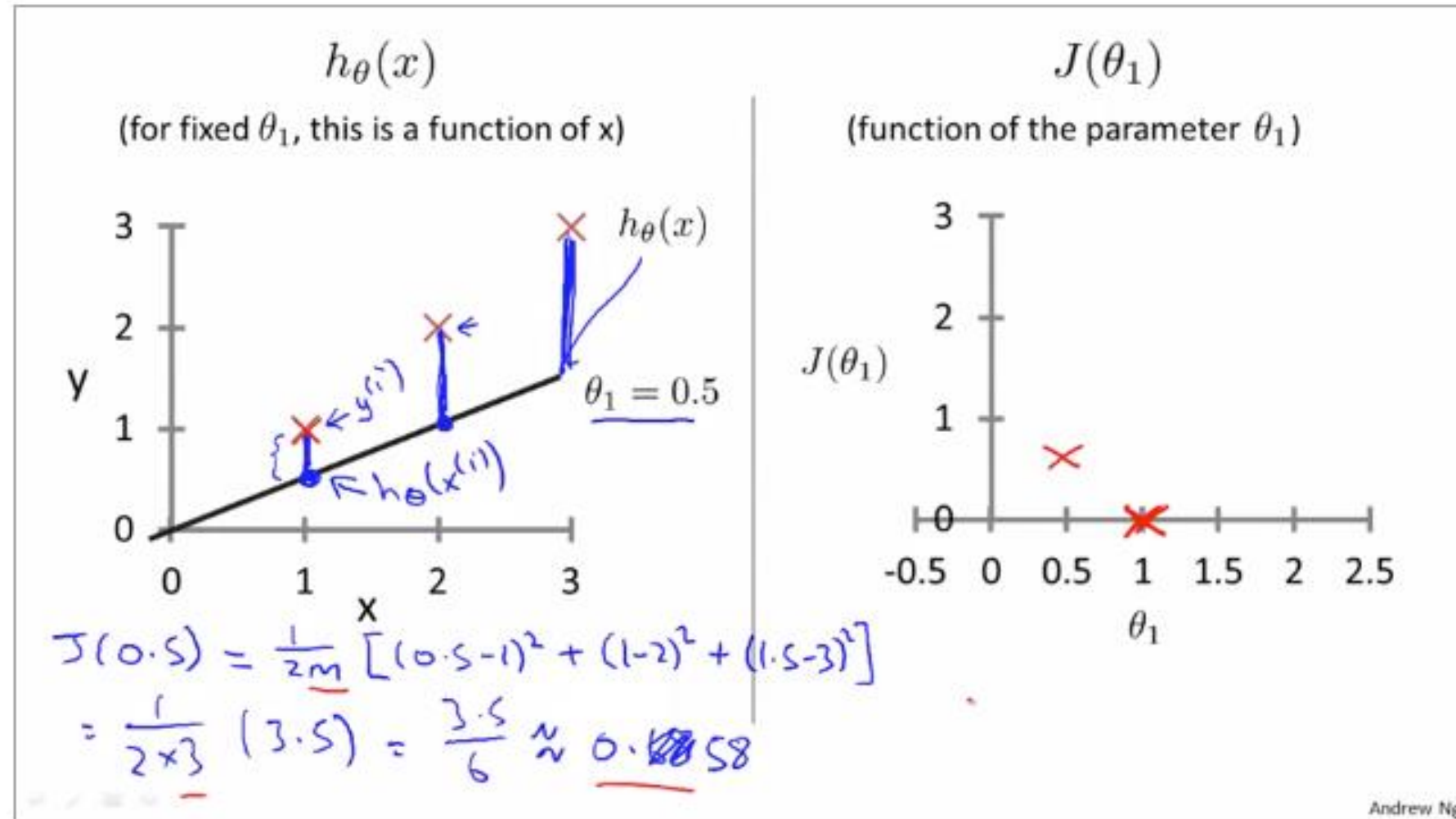$$J(\theta_0, \theta_1) = \frac{1}{2m} \sum_{i=1}^{m} \left( h_\theta(x^{(i)}) - y^{(i)} \right)^2$$

**Goal:** $\underset{\theta_0, \theta_1}{\text{minimize}} \; J(\theta_0, \theta_1)$

Andrew Ng

# Cost Function - Intuition I

Consider Theta0 = 0

$\rightarrow h_\theta(x)$

(for fixed $\theta_1$, this is a function of x)

$\rightarrow J(\theta_1)$

(function of the parameter $\theta_1$)



$\theta_1 = 1$

$h_\theta(x^{(i)}) = y^{(i)}$

$J(\theta_1) = \frac{1}{2m} \sum_{i=1}^{m} (h_\theta(x^{(i)}) - y^{(i)})^2$

$= \frac{1}{2m} \sum_{i=1}^{m} (\theta_1 x^{(i)} - y^{(i)})^2 = \frac{1}{2m}(0^2 + 0^2 + 0^2) = 0^2$

$J(1) = 0$

# Cost Function - Intuition I



$$h_\theta(x)$$

(for fixed $\theta_1$, this is a function of x)

$$J(\theta_1)$$

(function of the parameter $\theta_1$)

$$J(0.5) = \frac{1}{2m}\left[(0.5-1)^2 + (1-2)^2 + (1.5-3)^2\right]$$

$$= \frac{1}{2\times3}(3.5) = \frac{3.5}{6} \approx 0.58$$

# Cost Function - Intuition I



$h_\theta(x)$

(for fixed $\theta_1$, this is a function of x)

$J(\theta_1)$

(function of the parameter $\theta_1$)

$\theta_1 = 0$

$J(\theta_1)$

$J(0) = \frac{1}{2m}(1^2 + 2^2 + 3^2)$

$= \frac{1}{6} \cdot 14 \approx 2.3$

$h(x) = -0.5x$

Andrew Ng

# Cost Function - Intuition II

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

Parameters: $\theta_0, \theta_1$

Cost Function: $J(\theta_0, \theta_1) = \frac{1}{2m} \sum_{i=1}^{m} \left( h_\theta(x^{(i)}) - y^{(i)} \right)^2$

Goal: $\underset{\theta_0, \theta_1}{\text{minimize}} \ J(\theta_0, \theta_1)$

Andrew Ng

Considering both parameters

# Cost Function - Intuition II

$$h_\theta(x)$$

(for fixed $\theta_0, \theta_1$, this is a function of x)

$$J(\theta_0, \theta_1)$$

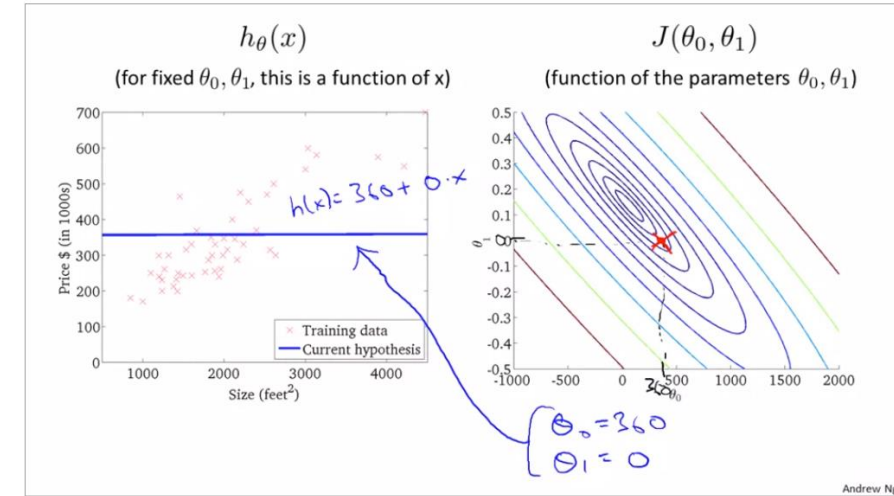(function of the parameters $\theta_0, \theta_1$)

# Cost Function - Intuition II

Any point on the same eclipses will give us the same value of error function J.

contour plot



Cost Function - Intuition II

$h_\theta(x)$
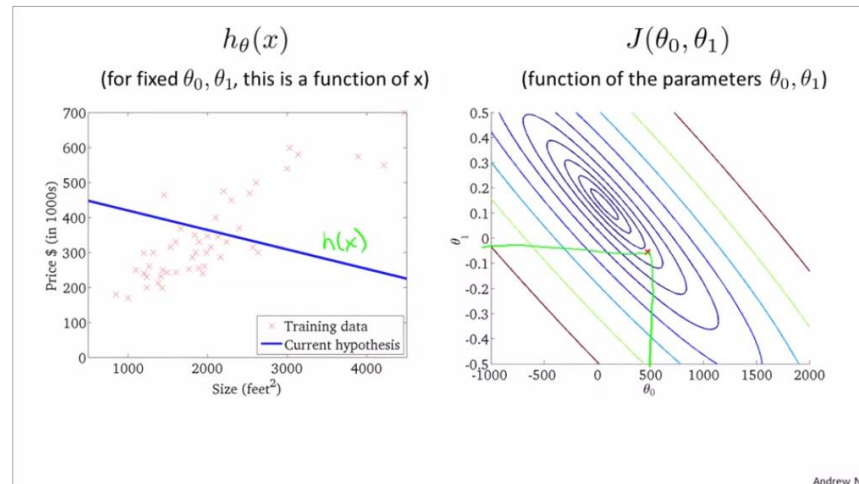(for fixed $\theta_0, \theta_1$, this is a function of x)

$J(\theta_0, \theta_1)$
(function of the parameters $\theta_0, \theta_1$)

$J(\theta_0, \theta_1)$



Cost Function - Intuition II

$h_\theta(x)$
(for fixed $\theta_0, \theta_1$, this is a function of x)

$J(\theta_0, \theta_1)$
(function of the parameters $\theta_0, \theta_1$)

$h(x) = 360 + 0 \cdot x$
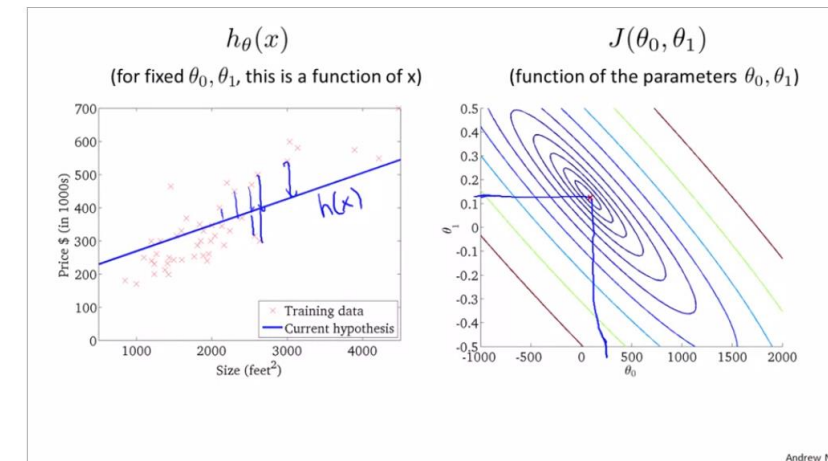
$\begin{pmatrix} \theta_0 = 360 \\ \theta_1 = 0 \end{pmatrix}$

We want algorithm which Automatically finding the value of Theta0 and Theta1 and minimize J (Cost function)

Cost Function - Intuition II



$h_\theta(x)$
(for fixed $\theta_0, \theta_1$, this is a function of x)

$J(\theta_0, \theta_1)$
(function of the parameters $\theta_0, \theta_1$)

$h(x)$

Cost Function - Intuition II



$h_\theta(x)$
(for fixed $\theta_0, \theta_1$, this is a function of x)

$J(\theta_0, \theta_1)$
(function of the parameters $\theta_0, \theta_1$)

$h(x)$

# Gradient Descent

- Automatically finding the value of Theta0 and Theta1 and minimize J (Cost function) for Linear Regression.
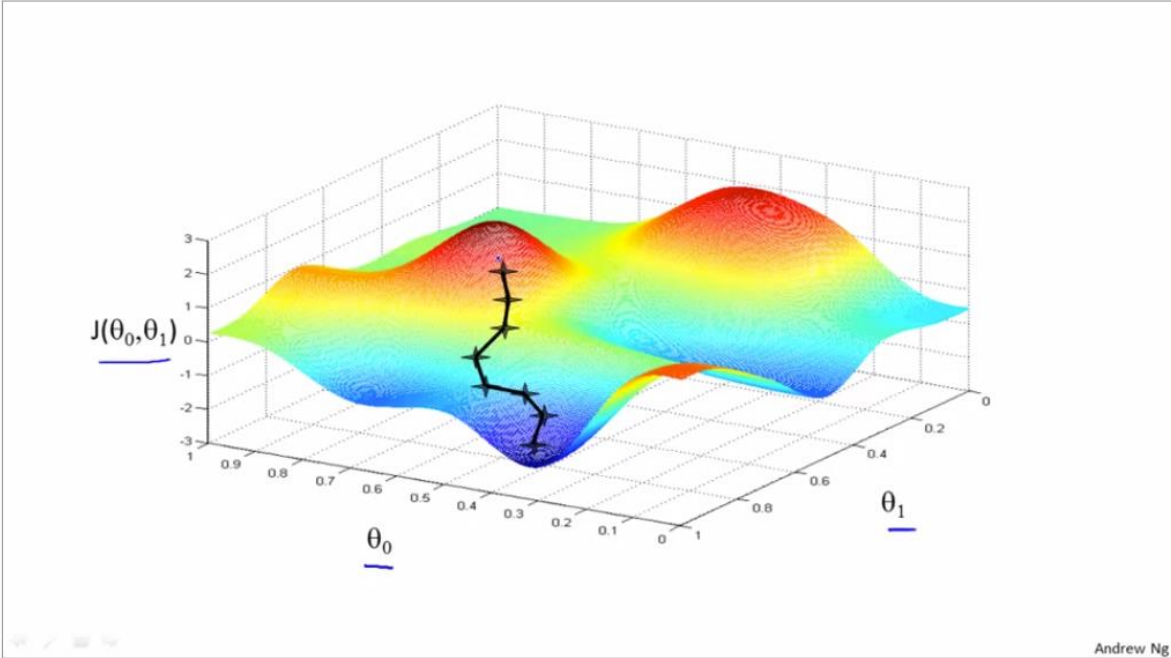
- Used for other algorithm as well.

Have some function $J(\theta_0, \theta_1)$ $\quad J(\theta_0, \theta_1, \theta_2, \ldots, \theta_n)$

Want $\min\limits_{\theta_0, \theta_1} J(\theta_0, \theta_1)$ $\qquad \min\limits_{\theta_0 \ldots \theta_n} J(\theta_0, \ldots, \theta_n)$

**Outline:**

- Start with some $\theta_0, \theta_1$ $\quad ($ Say $\theta_0 = 0, \theta_1 = 0)$
- Keep changing $\theta_0, \theta_1$ to reduce $J(\theta_0, \theta_1)$

  until we hopefully end up at a minimum

Andrew Ng

## Gradient Descent



## Gradient Descent

**Gradient descent algorithm**

$\theta_0, \theta_1$

repeat until convergence {

$$\theta_j := \theta_j - \alpha \frac{\partial}{\partial \theta_j} J(\theta_0, \theta_1)$$

(for $j = 0$ and $j = 1$)

}

learning rate

Simultaneously update $\theta_0$ and $\theta_1$

Assignment
$a := b$
$a := a+1$

Truth assertion
$a = b$
$a = a+1$ ✗

---

**Correct: Simultaneous update**

$$temp0 := \theta_0 - \alpha \frac{\partial}{\partial \theta_0} J(\theta_0, \theta_1)$$
$$temp1 := \theta_1 - \alpha \frac{\partial}{\partial \theta_1} J(\theta_0, \theta_1)$$
$$\theta_0 := temp0$$
$$\theta_1 := temp1$$

**Incorrect:**

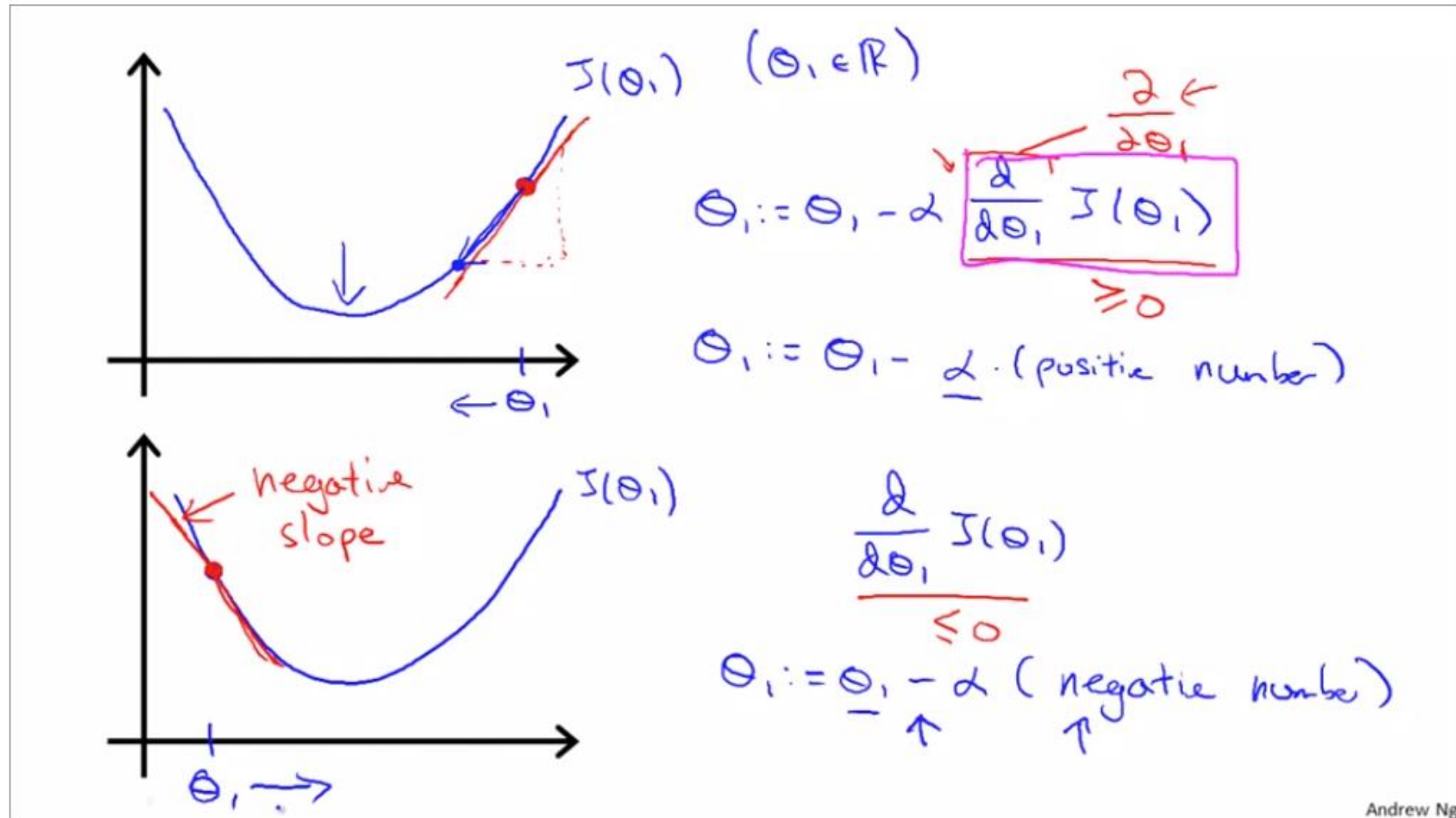$$temp0 := \theta_0 - \alpha \frac{\partial}{\partial \theta_0} J(\theta_0, \theta_1)$$
$$\theta_0 := temp0$$
$$temp1 := \theta_1 - \alpha \frac{\partial}{\partial \theta_1} J(\theta_0, \theta_1)$$
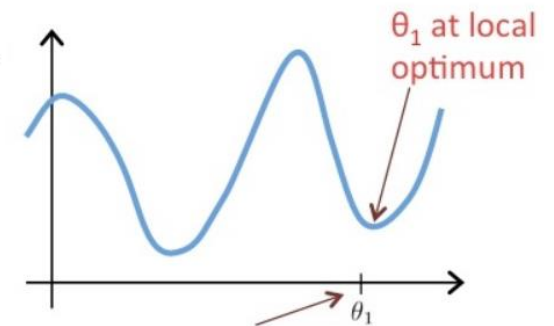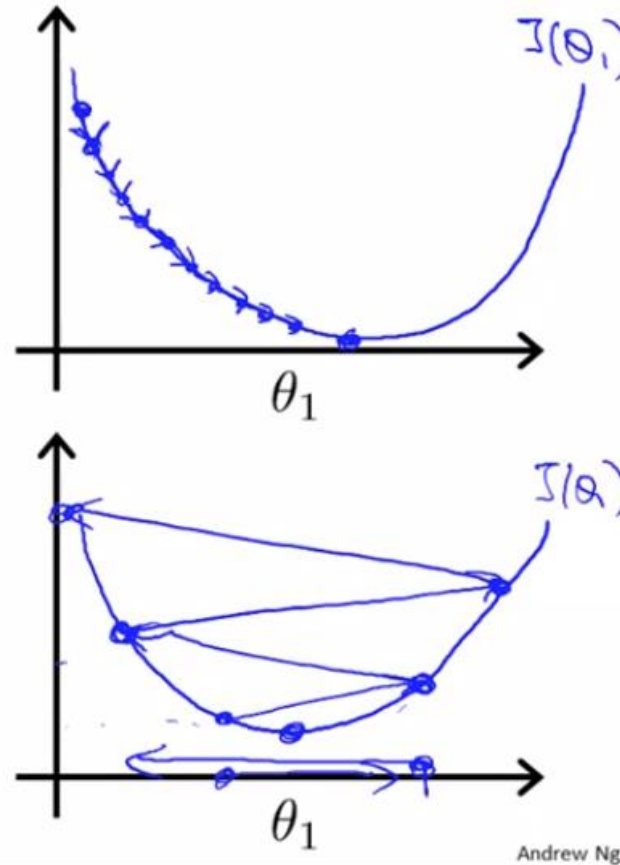$$\theta_1 := temp1$$

Andrew Ng

# Gradient Descent Intuition



$J(\theta_1)$ $(\theta_1 \in \mathbb{R})$

$$\theta_1 := \theta_1 - \alpha \frac{d}{d\theta_1} J(\theta_1)$$

$$\frac{d}{d\theta_1} \geq 0$$

$$\theta_1 := \theta_1 - \alpha \cdot (\text{positive number})$$

negative slope

$J(\theta_1)$

$$\frac{\frac{d}{d\theta_1} J(\theta_1)}{\leq 0}$$

$$\theta_1 := \theta_1 - \alpha (\text{negative number})$$

All about slope: https://www.occc.edu/wp-content/legacy/sem/mathhandouts/All%20About%20Slopes.pdf

# Gradient Descent Intuition

$$\theta_1 := \theta_1 - \alpha \frac{\partial}{\partial \theta_1} J(\theta_1)$$
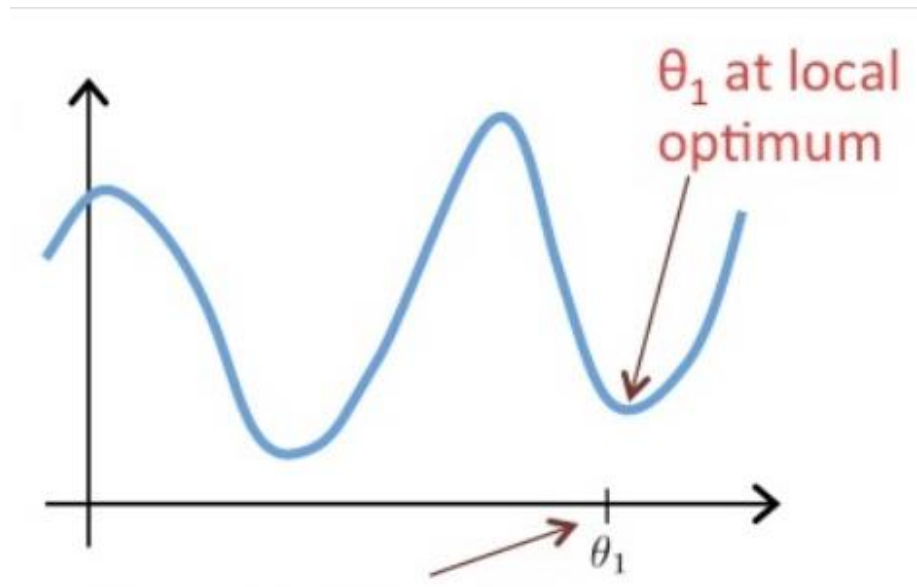
If $\alpha$ is too small, gradient descent can be slow.

If $\alpha$ is too large, gradient descent can overshoot the minimum. It may fail to converge, or even diverge.

$J(\theta_1)$

$J(\theta_1)$

$\theta_1$

$\theta_1$

Andrew Ng

Question: Suppose theta_1 is at a local optimum of J(\theta_1), such as shown in the figure. What will one step of gradient descent do?

$\theta_1$ at local optimum

$\theta_1$

Question: Suppose theta_1 is at a local optimum of J(\theta_1), such as shown in the figure. What will one step of gradient descent do?
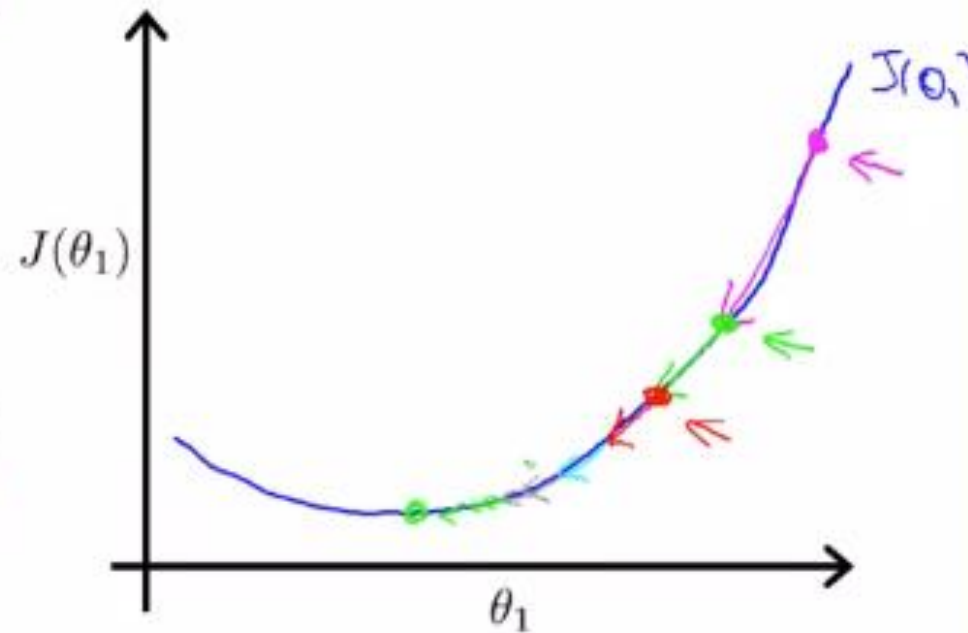
# Gradient Descent Intuition

Gradient descent can converge to a local minimum, even with the learning rate α fixed.

$$\theta_1 := \theta_1 - \alpha \frac{d}{d\theta_1} J(\theta_1)$$

As we approach a local minimum, gradient descent will automatically take smaller steps. So, no need to decrease α over time.



Andrew Ng

Linear regression
with one variable

Gradient descent for
linear regression

Machine Learning

## Gradient Descent For Linear Regression

**Gradient descent algorithm**

repeat until convergence {
$$\theta_j := \theta_j - \alpha \frac{\partial}{\partial \theta_j} J(\theta_0, \theta_1)$$
(for $j = 1$ and $j = 0$)
}

**Linear Regression Model**

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

$$J(\theta_0, \theta_1) = \frac{1}{2m} \sum_{i=1}^{m} \left( h_\theta(x^{(i)}) - y^{(i)} \right)^2$$

- $\underset{\theta_0, \theta_1}{\text{Min}} \ J(\theta_0, \theta_1)$

Andrew Ng

# Gradient Descent For Linear Regression

$$\frac{\partial}{\partial \theta_j} J(\theta_0, \theta_1) = \frac{\partial}{\partial \theta_j} \cdot \frac{1}{2m} \sum_{i=1}^{m} \left( h_\theta(x^{(i)}) - y^{(i)} \right)^2$$

$$= \frac{\partial}{\partial \theta_j} \frac{1}{2m} \sum_{i=1}^{m} \left( \theta_0 + \theta_1 x^{(i)} - y^{(i)} \right)^2$$

$$\theta_0 \, j = 0 : \frac{\partial}{\partial \theta_0} J(\theta_0, \theta_1) = \frac{1}{m} \sum_{i=1}^{m} \left( h_\theta(x^{(i)}) - y^{(i)} \right)$$

$$\theta_1 \, j = 1 : \frac{\partial}{\partial \theta_1} J(\theta_0, \theta_1) = \frac{1}{m} \sum_{i=1}^{m} \left( h_\theta(x^{(i)}) - y^{(i)} \right) \cdot x^{(i)}$$

Andrew Ng

# Gradient descent algorithm

$$\frac{\partial}{\partial \theta_0} J(\theta_0, \theta_1)$$

repeat until convergence {

$$\theta_0 := \theta_0 - \alpha \frac{1}{m} \sum_{i=1}^{m} \left( h_\theta(x^{(i)}) - y^{(i)} \right)$$

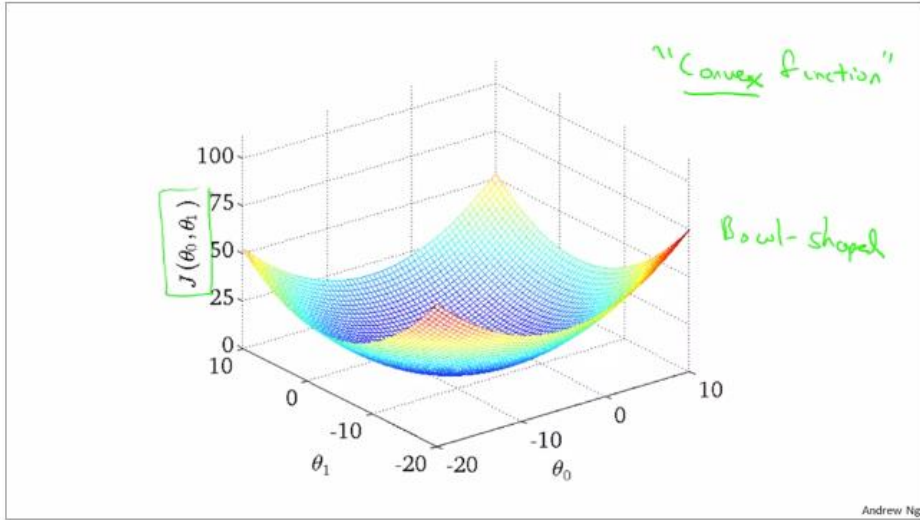$$\theta_1 := \theta_1 - \alpha \frac{1}{m} \sum_{i=1}^{m} \left( h_\theta(x^{(i)}) - y^{(i)} \right) \cdot x^{(i)}$$
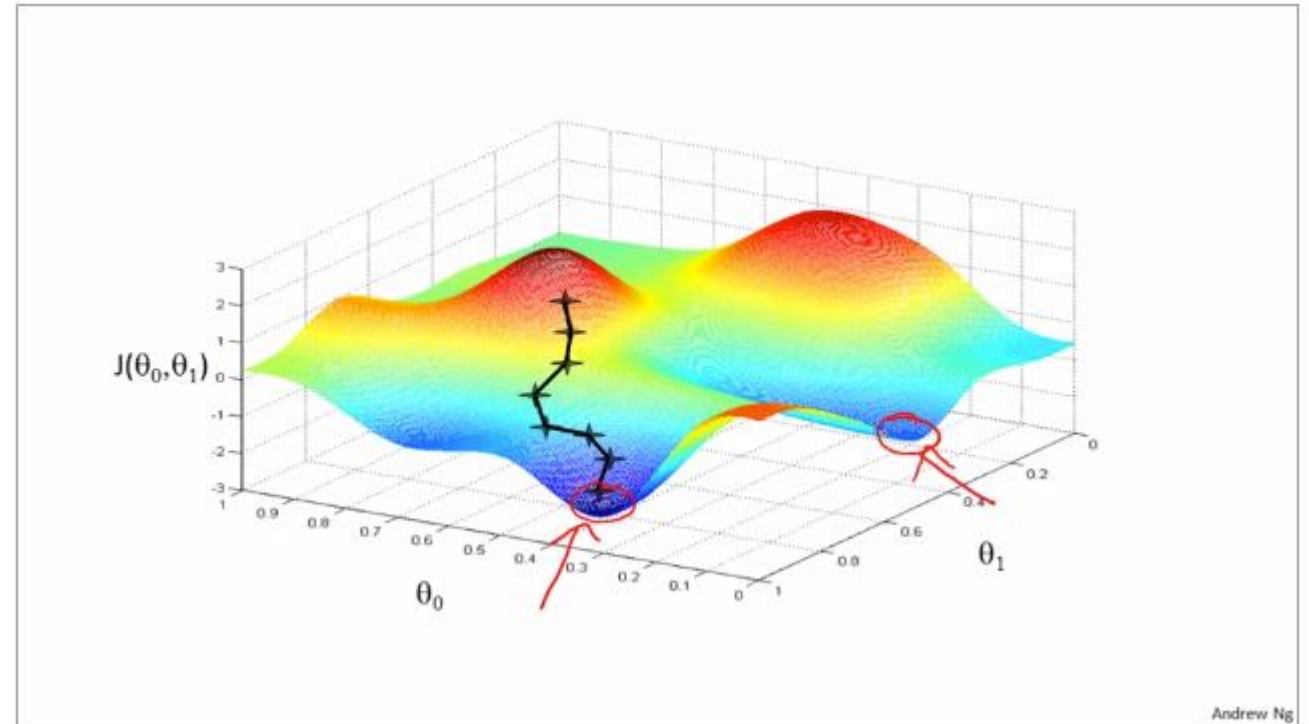
}

update $\theta_0$ and $\theta_1$ simultaneously

$$\frac{\partial}{\partial \theta_1} J(\theta_0, \theta_1)$$

Same local and global minima



Multiple local minima

# Multiple Features



Linear Regression with multiple variables

Multiple features

Machine Learning

# Multiple Features

**Multiple features (variables).**

| Size (feet²) | Number of bedrooms | Number of floors | Age of home (years) | Price ($1000) |
|---|---|---|---|---|
| $x_1$ | $x_2$ | $x_3$ | $x_4$ | $y$ |
| 2104 | 5 | 1 | 45 | 460 |
| 1416 | 3 | 2 | 40 | 232 |
| 1534 | 3 | 2 | 30 | 315 |
| 852 | 2 | 1 | 36 | 178 |
| ... | ... | ... | ... | ... |

$m = 47$

$$x^{(2)} = \begin{bmatrix} 1416 \\ 3 \\ 2 \\ 40 \end{bmatrix}$$

Notation:

$n$ = number of features    $n = 4$

$x^{(i)}$ = input (features) of $i^{th}$ training example.

$x_j^{(i)}$ = value of feature $j$ in $i^{th}$ training example.

Andrew Ng

# Multiple Features

Hypothesis:

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

$$h_\theta(x) = \theta_0 + \theta_1 x_1 + \theta_2 x_2 + \theta_3 x_3 + \theta_4 x_4$$

E.g. $h_\theta(x) = 80 + 0.1 x_1 + 0.01 x_2 + 3 x_3 - 2 x_4$

# Multiple Features

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

For convenience of notation, define $x_0 = 1$. $\quad \left( x_0^{(i)} = 1 \right)$

$$x = \begin{bmatrix} x_0 \\ x_1 \\ x_2 \\ \vdots \\ x_n \end{bmatrix} \in \mathbb{R}^{n+1} \qquad \theta = \begin{bmatrix} \theta_0 \\ \theta_1 \\ \theta_2 \\ \vdots \\ \theta_n \end{bmatrix} \in \mathbb{R}^{n+1}$$

$$[\theta_0 \ \theta_1 \cdots \theta_n] \begin{bmatrix} x_0 \\ x_1 \\ \vdots \\ x_n \end{bmatrix}$$

$$\underbrace{\qquad}_{\theta^T}$$

$(n+1) \times 1$ matrix

$x$

$, \theta^T x$

$$h_\theta(x) = \theta_0 x_0 + \theta_1 x_1 + \cdots + \theta_n x_n \quad \leftarrow \quad (\ell = 1)$$

$$= \boxed{\theta^T x.}$$

Multivariate linear regression. $\leftarrow$

Andrew Ng

# Gradient Descent for Multiple Variables

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

$\rightarrow x_0 = 1$

**Parameters:** $\theta_0, \theta_1, \ldots, \theta_n$ $\quad \theta \quad$ n+1 - dimensional vector

**Cost function:**
$$J(\theta_0, \theta_1, \ldots, \theta_n) = \frac{1}{2m} \sum_{i=1}^{m} (h_\theta(x^{(i)}) - y^{(i)})^2$$

$J(\theta)$

**Gradient descent:**

Repeat {

$\rightarrow \quad \theta_j := \theta_j - \alpha \frac{\partial}{\partial \theta_j} J(\theta_0, \ldots, \theta_n) \quad J(\theta)$

}
$\quad\quad$ (simultaneously update for every $j = 0, \ldots, n$)

Andrew Ng

# Gradient Descent for Multiple Variables

**Gradient Descent**

Previously (n=1):

Repeat {

$$\theta_0 := \theta_0 - \alpha \frac{1}{m} \sum_{i=1}^{m} (h_\theta(x^{(i)}) - y^{(i)})$$

$$\boxed{\frac{\partial}{\partial \theta_0} J(\theta)}$$

$$\theta_1 := \theta_1 - \alpha \frac{1}{m} \sum_{i=1}^{m} (h_\theta(x^{(i)}) - y^{(i)})x^{(i)}$$

$$x_1^{(i)}$$

(simultaneously update $\theta_0, \theta_1$)

}

New algorithm $(n \geq 1)$:

Repeat {

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

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

(simultaneously update $\theta_j$ for $j = 0, \ldots, n$)

$$x_0^{(i)} = 1$$

}

$$\theta_0 := \theta_0 - \alpha \frac{1}{m} \sum_{i=1}^{m} (h_\theta(x^{(i)}) - y^{(i)})x_0^{(i)}$$

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

$$\theta_2 := \theta_2 - \alpha \frac{1}{m} \sum_{i=1}^{m} (h_\theta(x^{(i)}) - y^{(i)})x_2^{(i)}$$

. . .

Andrew Ng
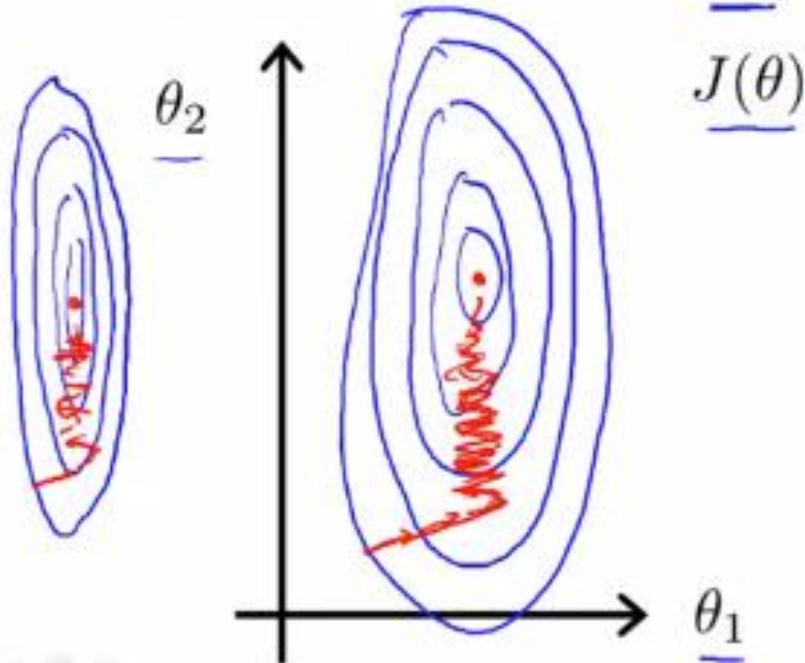
# Gradient Descent in Practice I - Feature Scaling

Make gradient Descent run much faster and take fewer steps

**Feature Scaling**

Idea: Make sure features are on a similar scale.
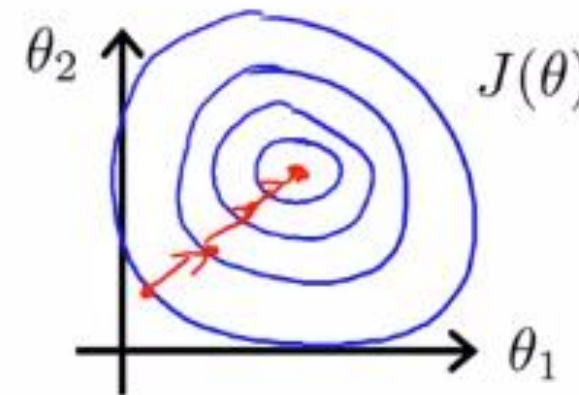
E.g. $x_1$ = size (0-2000 feet²) ←

$x_2$ = number of bedrooms (1-5) ←

$\rightarrow x_1 = \dfrac{\text{size (feet}^2)}{2000}$

$\rightarrow x_2 = \dfrac{\text{number of bedrooms}}{5}$

$0 \leq x_1 \leq 1 \qquad 0 \leq x_2 \leq 1$

$J(\theta)$

$\theta_2$

$\theta_1$

$J(\theta)$

$\theta_2$

$\theta_1$

Andrew Ng

# Gradient Descent in Practice I - Feature Scaling

**Feature Scaling**

Get every feature into approximately a $-1 \leq x_i \leq 1$ range.

$$x_0 = 1$$

$$0 \leq x_1 \leq 3 \quad \checkmark$$

$$-2 \leq x_2 \leq 0.5 \quad \checkmark$$

$$-100 \leq x_3 \boxed{100} \quad \times$$

$$-0.0001 \leq x_4 \leq \boxed{0.0001} \quad \times$$

Andrew Ng

# Gradient Descent in Practice I - Feature Scaling

**Mean normalization**

Replace $x_i$ with $x_i - \mu_i$ to make features have approximately zero mean (Do not apply to $x_0 = 1$).

E.g. $\rightarrow$ $x_1 = \dfrac{size - 1000}{2000}$

$x_2 = \dfrac{\#bedrooms - 2}{5}$

$-0.5 \leq x_1 \leq 0.5, \; -0.5 \leq x_2 \leq 0.5$

Average size = 1000

1-5 bedrooms

$x_1 \leftarrow \dfrac{x_1 - \boxed{\mu_1}}{\boxed{S_1}}$   ← avg value of $x_1$ in training set

$\quad$ range (max - min)
(or standard deviation)

$x_2 \leftarrow \dfrac{x_2 - \mu_1}{S_2}$
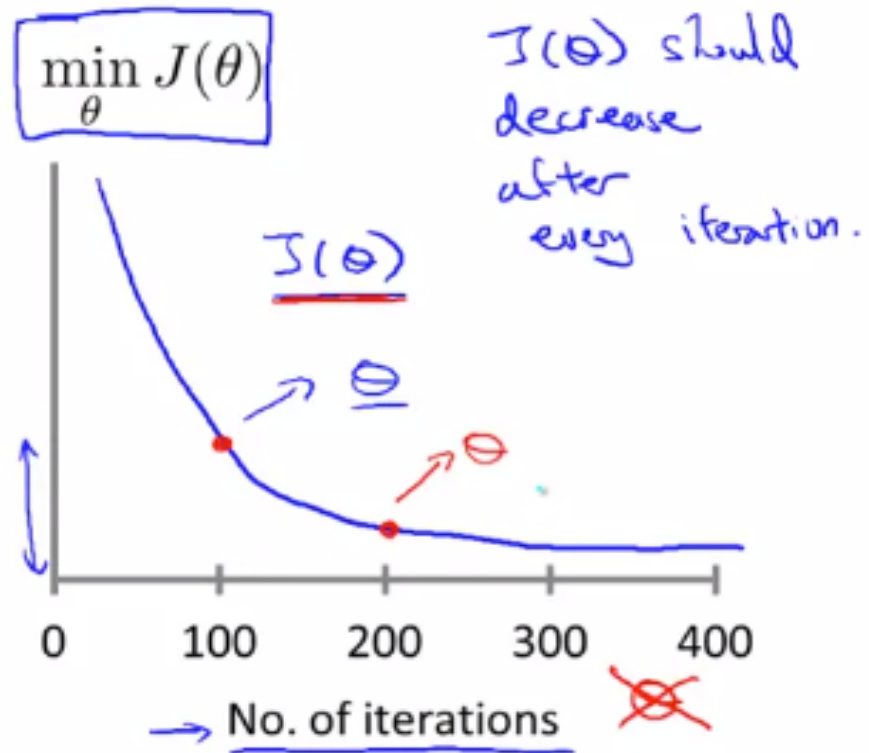
# Gradient Descent in Practice II - Learning Rate

**Gradient descent**

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

- "Debugging": How to make sure gradient descent is working correctly.

- How to choose learning rate $\alpha$.

Andrew Ng

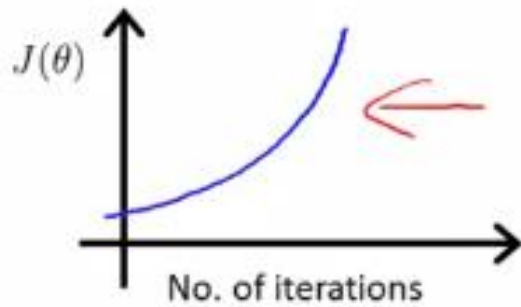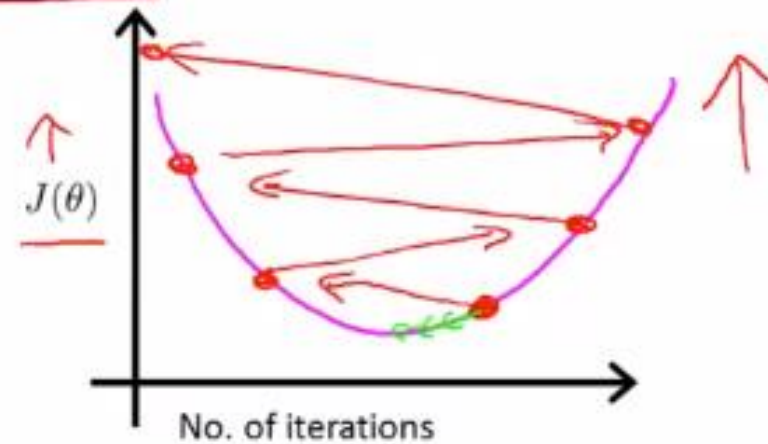# Gradient Descent in Practice II - Learning Rate



**Making sure gradient descent is working correctly.**

$$\min_{\theta} J(\theta)$$

$J(\theta)$

$J(\theta)$ should decrease after every iteration.

0    100    200    300    400

No. of iterations

Andrew Ng

# Gradient Descent in Practice II - Learning Rate



**Making sure gradient descent is working correctly.**

$J(\theta)$

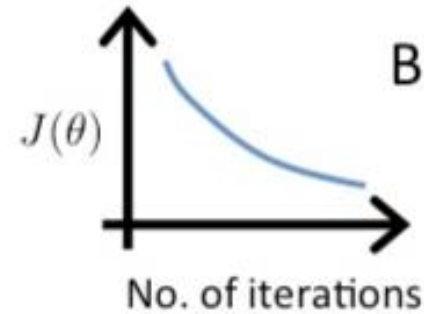No. of iterations

Gradient descent not working.

Use smaller $\alpha$.

$J(\theta)$

No. of iterations

Suppose a friend ran gradient descent three times, with $\alpha = 0.01$, $\alpha = 0.1$, and $\alpha = 1$, and got the following three plots (labeled A, B, and C):
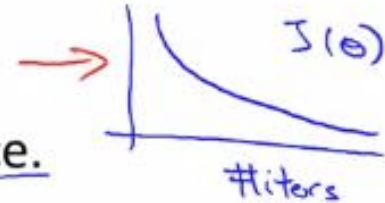


$J(\theta)$ — A — No. of iterations

$J(\theta)$ — B — No. of iterations

$J(\theta)$ — C — No. of iterations

Which plots corresponds to which values of $\alpha$?

○ A is $\alpha = 0.1$, B is $\alpha = 0.01$, C is $\alpha = 1$.

**Summary:**

- If $\alpha$ is too small: slow convergence.
- If $\alpha$ is too large: $J(\theta)$ may not decrease on every iteration; may not converge. (Slow converge also possible.)

$J(\theta)$

#iters

To choose $\alpha$, try

$\ldots, 0.001, \quad , 0.01, \quad , 0.1, \quad , 1, \ldots$

# Housing prices prediction

$$h_\theta(x) = \theta_0 + \theta_1 \times frontage + \theta_2 \times depth$$

$$\underbrace{frontage}_{x_1} \qquad \underbrace{depth}_{x_2}$$

Area

$$x = frontage * depth$$

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

↖ land area

$m$ **training examples,** $n$ **features.**

| Gradient Descent | Normal Equation |
|---|---|
| • Need to choose $\alpha$. | • No need to choose $\alpha$. |
| • Needs many iterations. | • Don't need to iterate. |
| • Works well even when $n$ is large. | • Need to compute $(X^T X)^{-1}$  $n \times n$  $O(n^3)$ |
|  | • Slow if $n$ is very large. |

$n = 10^6$

$n = 100$
$n = 1000$
$n = 10000$

Andrew Ng

# Logistic Regression

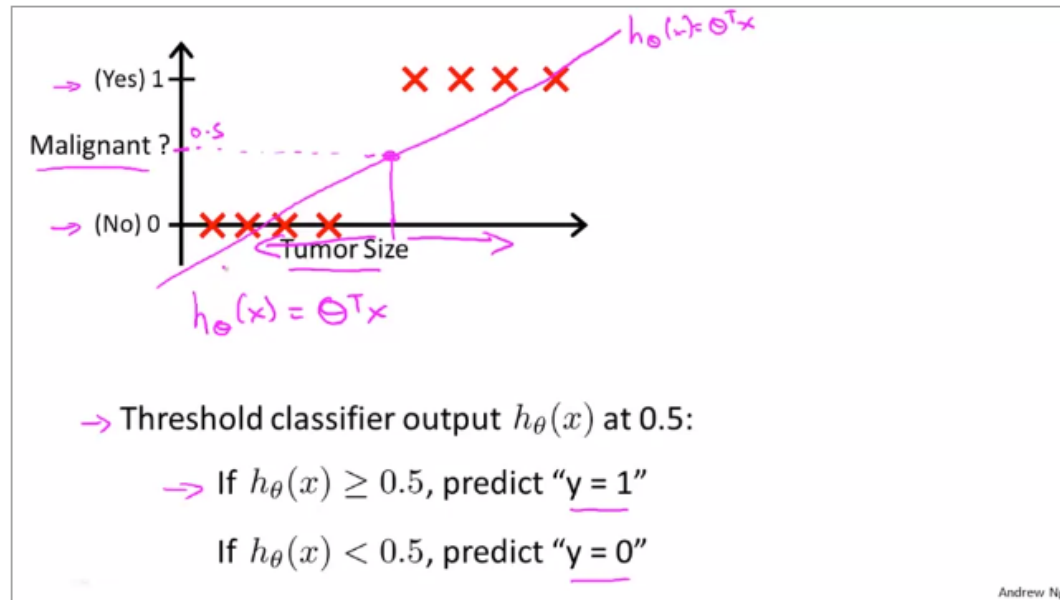# Logistic Regression

## Classification

Classification

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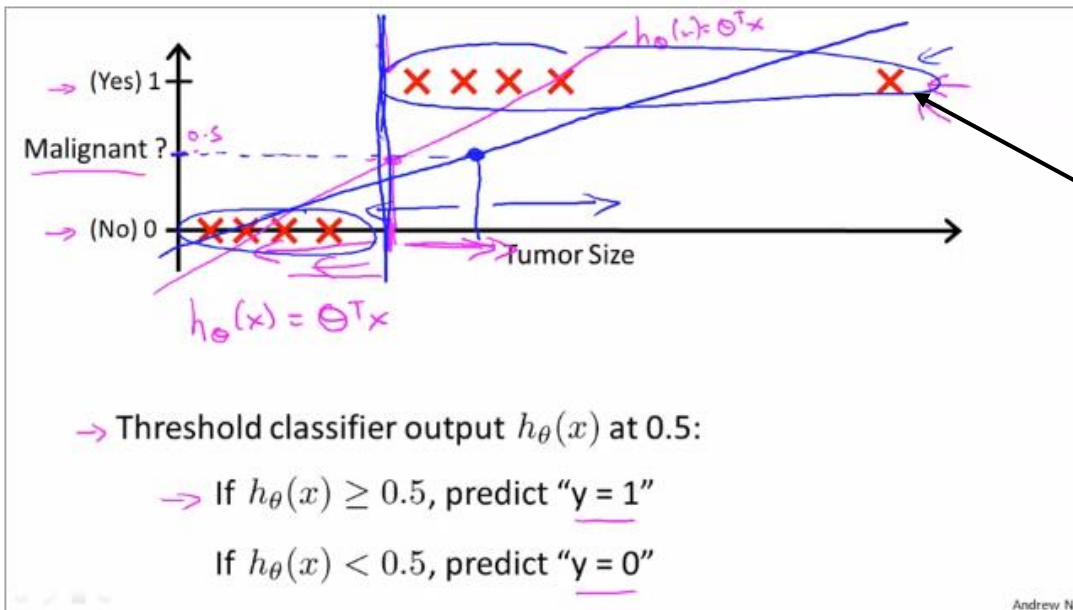→ Tumor: Malignant / Benign ?

$$y \in \{0, 1\}$$

0: "Negative Class" (e.g., benign tumor)
1: "Positive Class" (e.g., malignant tumor)

## Classification



Linear regression is doing good in this case



Linear regression is doing bad after adding one data in this case

# Hypothesis Representation

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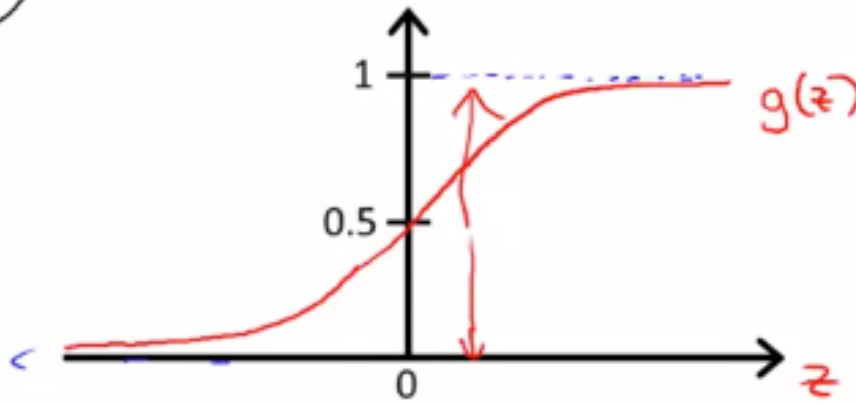
Want $0 \leq h_\theta(x) \leq 1$

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

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

$\theta^T x$

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

→ Sigmoid function
↳ Logistic function



Parameters $\theta$.

# Hypothesis Representation

$h_\theta(x)$

**Interpretation of Hypothesis Output**

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

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

$h_\theta(x) = 0.7$

$y = 1$

Tell patient that 70% chance of tumor being malignant

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

$y = 0$  or  $1$

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

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

Andrew Ng

# Decision Boundary

**Logistic regression**

$\rightarrow h_\theta(x) = g(\theta^T x) = P(y=1 | x; \theta)$
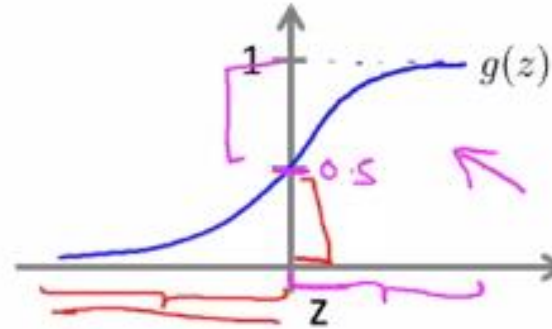
$\rightarrow g(z) = \frac{1}{1+e^{-z}}$

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

$\theta^T x \geq 0$

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

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

$\theta^T x < 0$

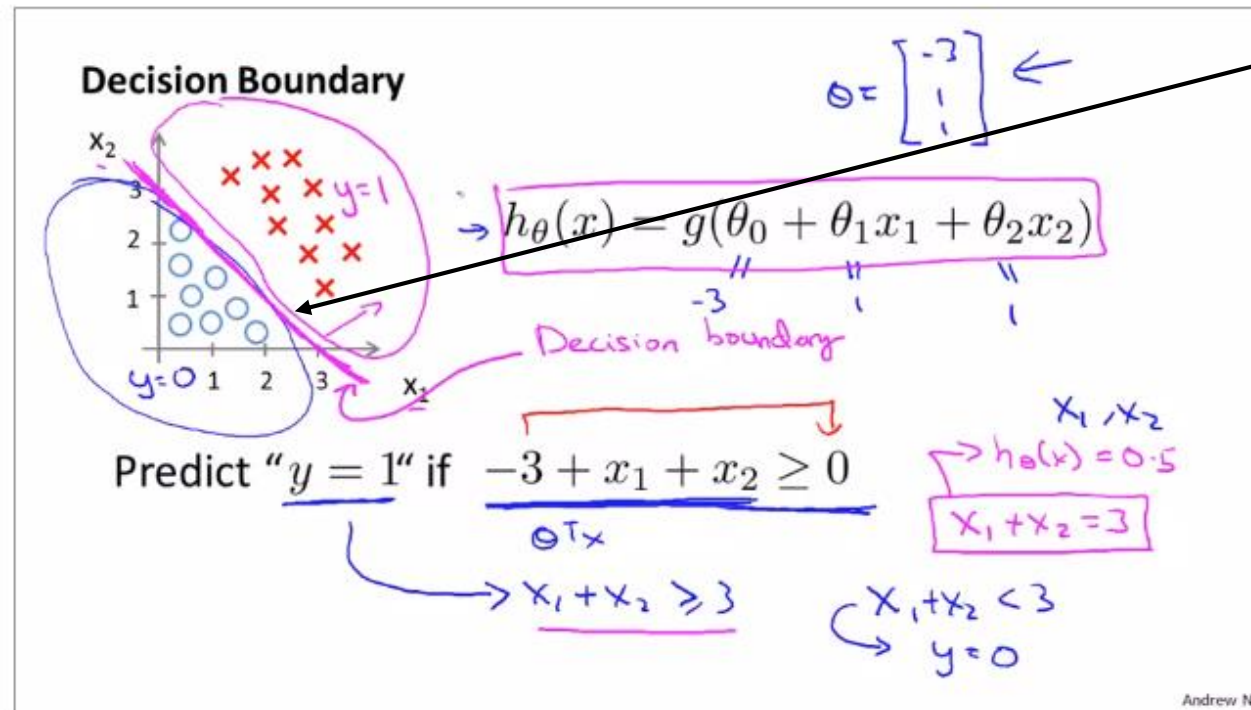$g(z) < 0.5$

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

$h_\theta(x) = g(\theta^T x) \geq 0.5$
wherever $\theta^T x \geq 0$

$z$

Andrew Ng

# Decision Boundary



Decision Boundary

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

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

Decision boundary

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

$\theta^T x$

$\rightarrow x_1 + x_2 \geq 3$

$\rightarrow h_\theta(x) = 0.5$

$x_1 + x_2 = 3$

$x_1 + x_2 < 3$

$\rightarrow y = 0$
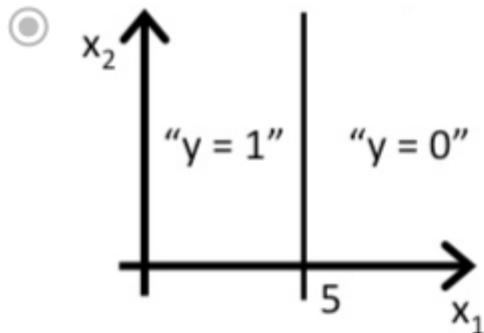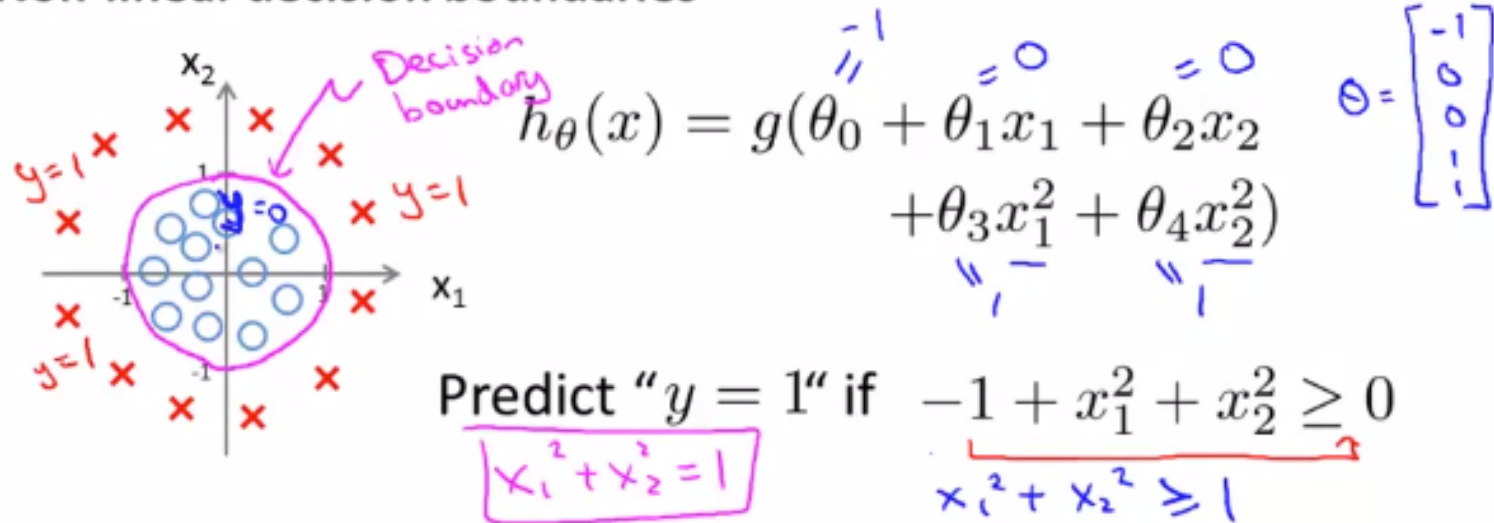
Andrew Ng

Decision boundary is property of hypothesis and Parameters but not of Data.

Consider logistic regression with two features $x_1$ and $x_2$. Suppose $\theta_0 = 5$, $\theta_1 = -1$, $\theta_2 = 0$, so that $h_\theta(x) = g(5 - x_1)$. Which of these shows the decision boundary of $h_\theta(x)$?



"y = 1"    "y = 0"

5

# Decision Boundary

## Non-linear decision boundaries



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

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

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

$$x_1^2 + x_2^2 = 1$$

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

# Cost Function

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

m examples $\qquad x \in \begin{bmatrix} x_0 \\ x_1 \\ \cdots \\ x_n \end{bmatrix} \mathbb{R}^{n+1} \qquad x_0 = 1, y \in \{0, 1\}$

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

How to choose parameters $\theta$ ?

# Cost Function

**Cost function**

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

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

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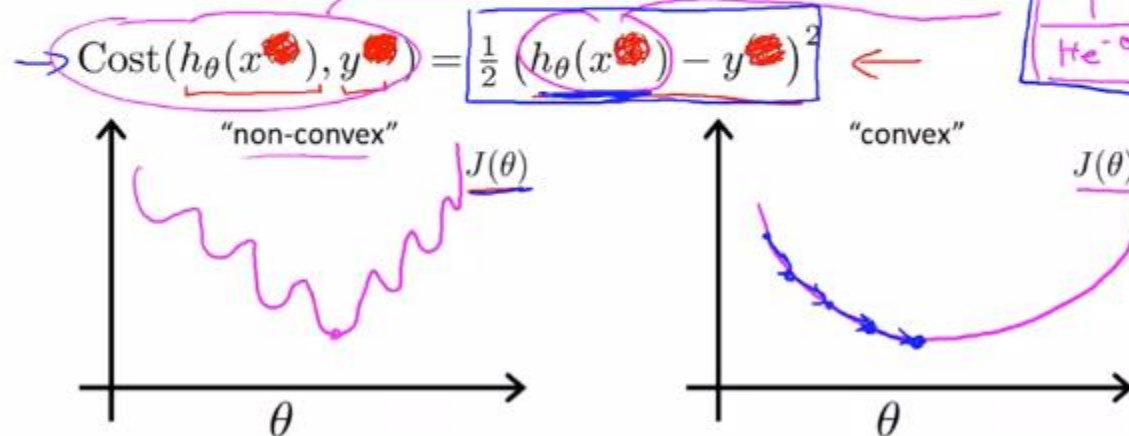
# Cost Function

**Cost function**

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

Logistic

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

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

$\frac{1}{1+e^{-\theta^T y}}$
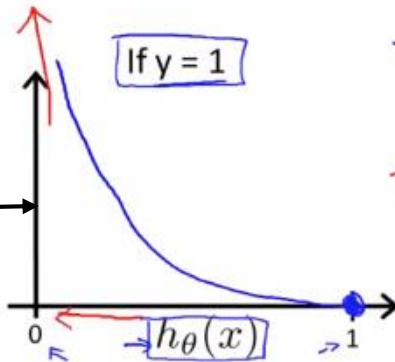
"non-convex"  $J(\theta)$

"convex"  $J(\theta)$

$\theta$

$\theta$

Due to sigmoid, this cost function will be non convex And that's why we do not use squared error loss in logistic regression

Andrew Ng

## Cost Function

### Logistic regression cost function

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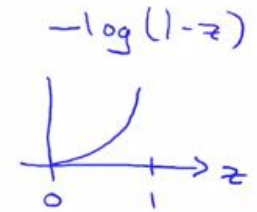
If y = 1

Cost

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

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.

$h_\theta(x)$   0   1

## Cost Function

### Logistic regression cost function

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

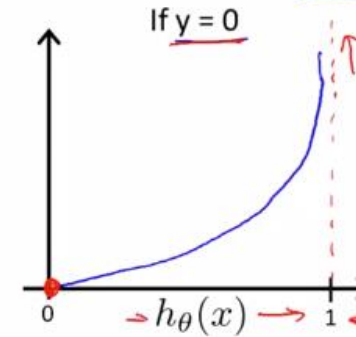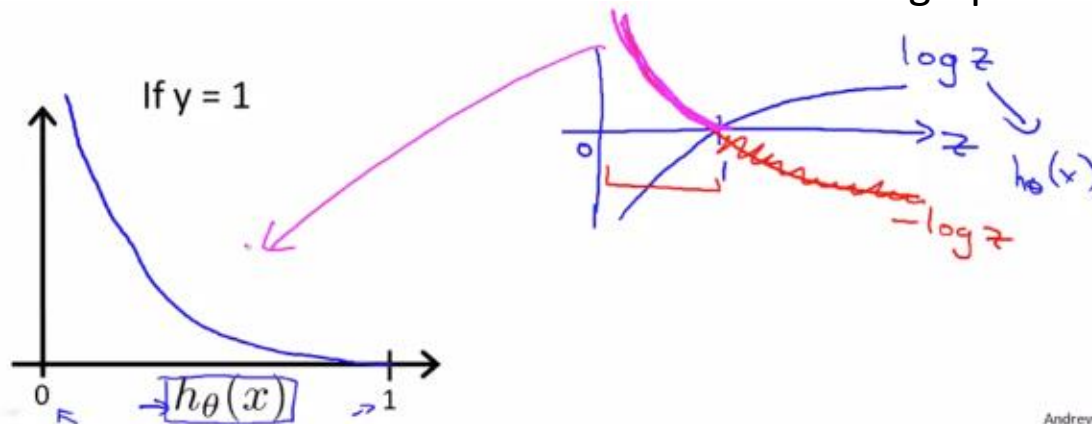If y = 0

$-\log(1 - z)$

$h_\theta(x)$   0   1

Andrew Ng

With proper choice of cost function we can have convex graph of it and it local optimum free

If y = 1

$\log z$

$-\log z$

$h_\theta(x)$

$h_\theta(x)$   0   1

Andrew N

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:

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

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

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

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

☐ If $y = 0$, then $\text{cost}(h_\theta(x), y) \rightarrow \infty$ as $h_\theta(x) \rightarrow 0$.

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

# 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

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

**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

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

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

Repeat {

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

(simultaneously update all $\theta_j$)

}

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

## Gradient Descent

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

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

$$\theta = \begin{bmatrix} \theta_0 \\ \theta_1 \\ \theta_2 \\ \vdots \\ \theta_n \end{bmatrix}$$

Repeat {

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

(simultaneously update all $\theta_j$)

}

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

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

Algorithm looks identical to linear regression!

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

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

$$\frac{d(\sigma(x))}{dx} = \frac{0 * (1 + e^{-x}) - (1) * (e^{-x} * (-1))}{(1 + e^{-x})^2}$$

$$\frac{d(\sigma(x))}{dx} = \frac{(e^{-x})}{(1 + e^{-x})^2} = \frac{1\text{-}1 + (e^{-x})}{(1 + e^{-x})^2} = \frac{1 + e^{-x}}{(1 + e^{-x})^2} - \frac{1}{(1 + e^{-x})^2}$$

$$\frac{d(\sigma(x))}{dx} = \frac{1}{1 + e^{-x}} * \left(1 - \frac{1}{1 + e^{-x}}\right) = \sigma(x)(1 - \sigma(x))$$

# Step 1:

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

Applying Chain rule and writing in terms of partial derivatives.

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

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

$$\frac{\partial(J(\theta))}{\partial(\theta_j)} = -\frac{1}{m} * (\sum_{i=1}^{m} \left[ y^{(i)} * \frac{1}{h_\theta(x^{(i)})} * \sigma(z)(1 - \sigma(z)) * \frac{\partial(\theta^T x)}{\partial(\theta_j)} \right]$$

$$+ \sum_{i=1}^{m} \left[ (1 - y^{(i)}) * \frac{1}{\left(1 - h_\theta(x^{(i)})\right)} * (-\sigma(z)(1 - \sigma(z)) * \frac{\partial(\theta^T x)}{\partial(\theta_j)} \right])$$

# Step 2:

Evaluating the partial derivative using the pattern of the derivative of the sigmoid function.

$$\frac{\partial(J(\theta))}{\partial(\theta j)} = -\frac{1}{m} * (\sum_{i=1}^{m} \left[ y^{(i)} * \frac{1}{h_\theta(x^{(i)})} * \sigma(z)(1-\sigma(z)) * \frac{\partial(\theta^T x)}{\partial(\theta j)} \right]$$

$$+ \sum_{i=1}^{m} \left[ (1-y^{(i)}) * \frac{1}{\left(1 - h_\theta(x^{(i)})\right)} * (\text{-}\sigma(z)(1-\sigma(z))* \frac{\partial(\theta^T x)}{\partial(\theta j)} \right])$$

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

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

# Step 3:

Simplifying the terms by multiplication

$$\frac{\partial(J(\theta))}{\partial(\theta j)} = -\frac{1}{m} * \left(\sum_{i=1}^{m} \left[y^{(i)} * \left(1 - h_\theta(x^{(i)})\right) * x_j^i - (1 - y^{(i)}) * h_\theta(x^{(i)}) * * x_j^i\right]\right)$$

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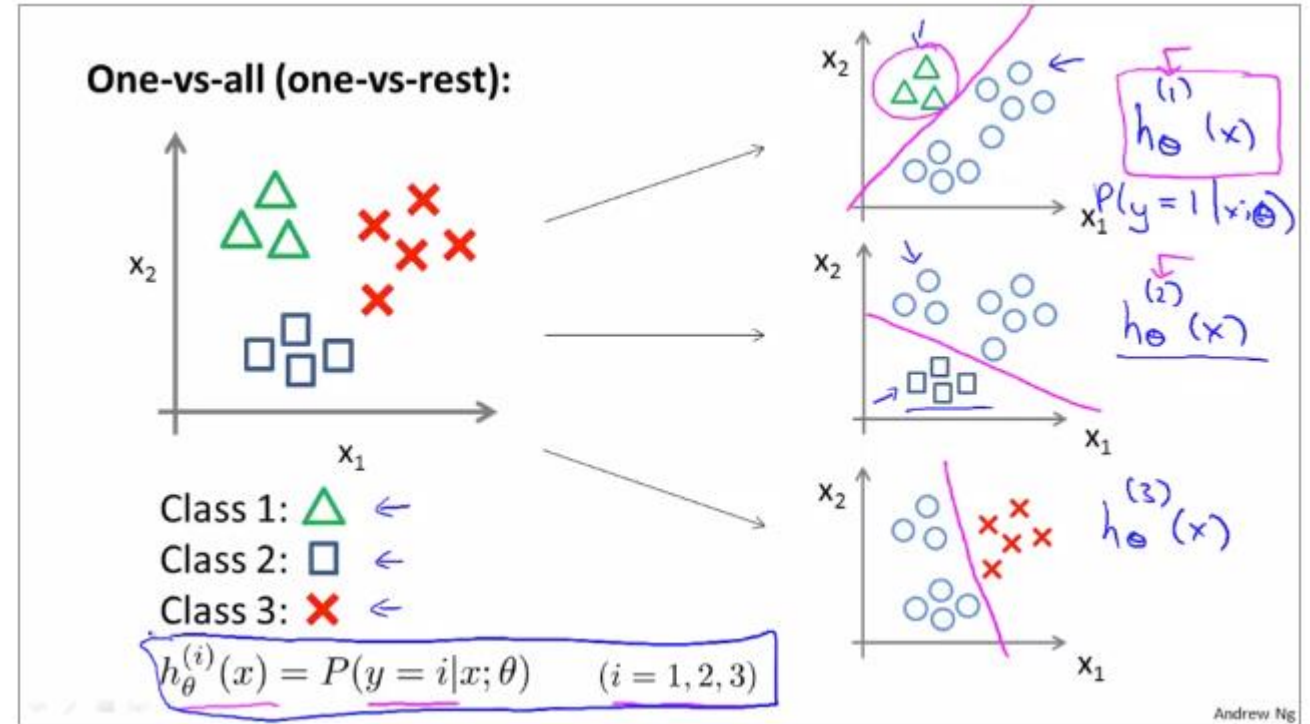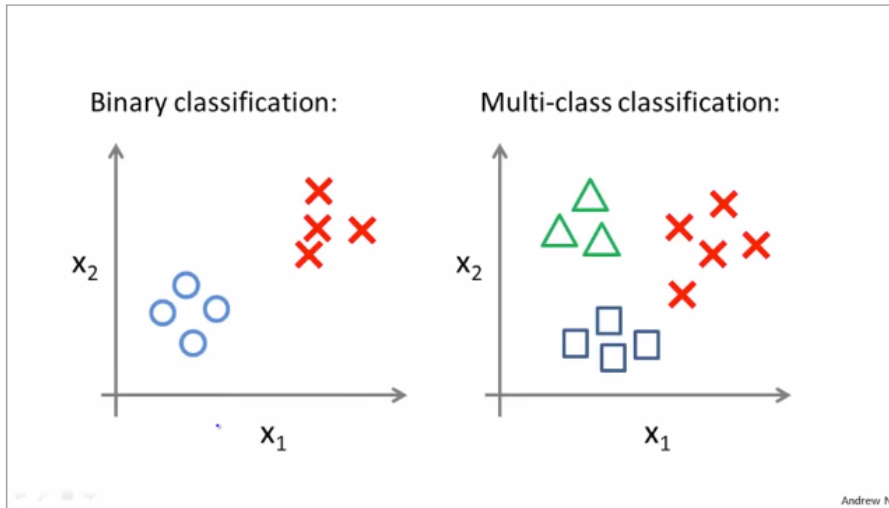
$$\frac{\partial(J(\theta))}{\partial(\theta j)} = -\frac{1}{m} * \left(\sum_{i=1}^{m} \left[y^{(i)} - h_\theta(x^{(i)})\right] * x_j^i\right)$$

# Step 4:

Removing the summation term by converting it into a matrix form for the gradient with respect to all the weights including the bias term.

$$\frac{\partial(J(\theta))}{\partial(\theta)} = \frac{1}{m} X^T [h_\theta(x) - y]$$

## Multiclass Classification: One-vs-all



Binary classification:

Multi-class classification:

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

Suppose you have a multi-class classification problem with $k$ classes (so $y \in \{1, 2, \ldots, k\}$). Using the 1-vs.-all method, how many different logistic regression classifiers will you end up training?
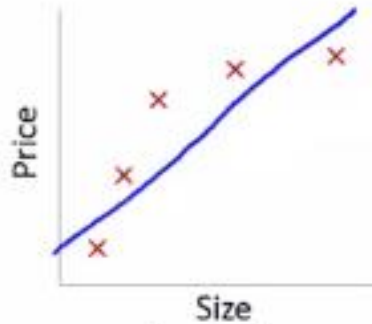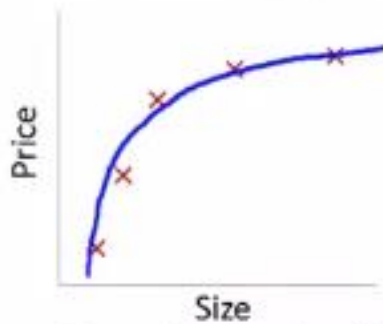
⦿ $k$

# Regularization

← Technique to deal with overfitting
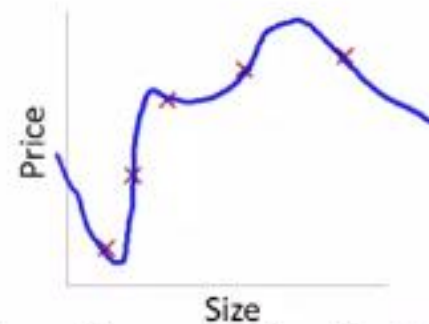
## The problem of overfitting

Example: Linear regression (housing prices)



$\rightarrow \theta_0 + \theta_1 x$

"Underfit"  "High bias"

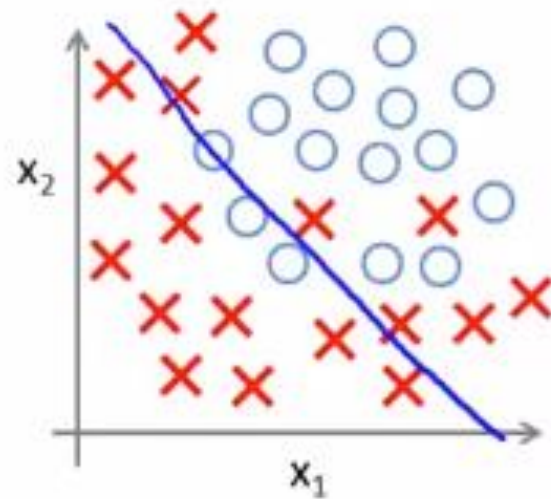$\rightarrow \theta_0 + \theta_1 x + \theta_2 x^2$

"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).
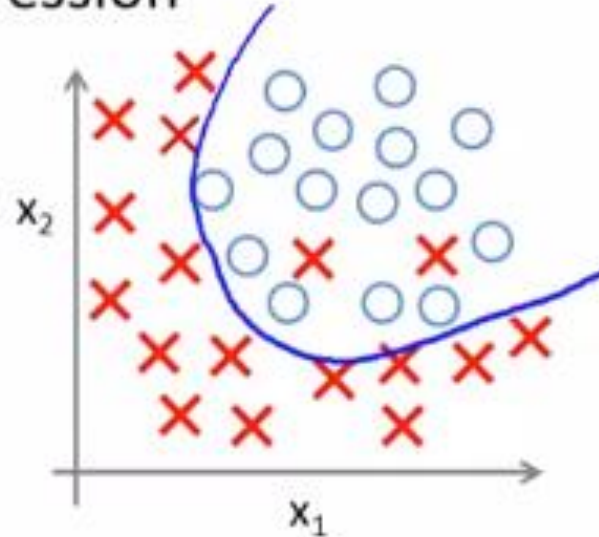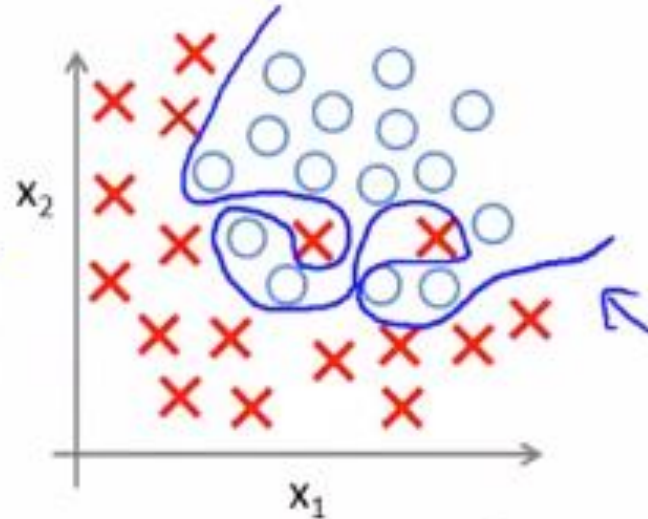
Andrew Ng

Example: Logistic regression

$x_2$ ... $x_1$

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

( $g$ = 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"

Andrew Ng

Consider the medical diagnosis problem of classifying tumors as malignant or benign. If a hypothesis $h_\theta(x)$ has overfit the training set, it means that:

It makes accurate predictions for examples in the training set and generalizes well to make accurate predictions on new, previously unseen examples.

It does not make accurate predictions for examples in the training set, but it does generalize well to make accurate predictions on new, previously unseen examples.

It makes accurate predictions for examples in the training set, but it does not generalize well to make accurate predictions on new, previously unseen examples.

It does not make accurate predictions for examples in the training set and does not generalize well to make accurate predictions on new, previously unseen examples.
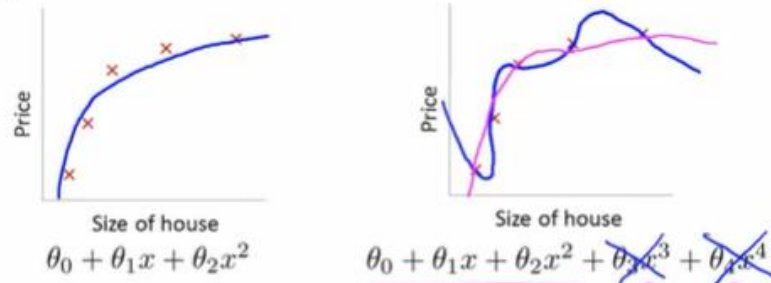
# Addressing Overfitting

1. Collect more Data
2. Select Features
3. Reduce size of parameter (Regularization)

Regularization.
— Keep all the features, but reduce magnitude/values of parameters $\theta_j$.
— Works well when we have a lot of features, each of which contributes a bit to predicting $y$.

# Cost Function

**Intuition**

Price / Size of house

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

Price / Size of house

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

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

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

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

**Regularization.**

Small values for parameters $\theta_0, \theta_1, \ldots, \theta_n$
- — "Simpler" hypothesis
- — Less prone to overfitting

$\rightarrow \theta_3, \theta_4 \approx 0$

Housing:
- — Features: $x_1, x_2, \ldots, x_{100}$
- — Parameters: $\theta_0, \theta_1, \theta_2, \ldots, \theta_{100}$

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

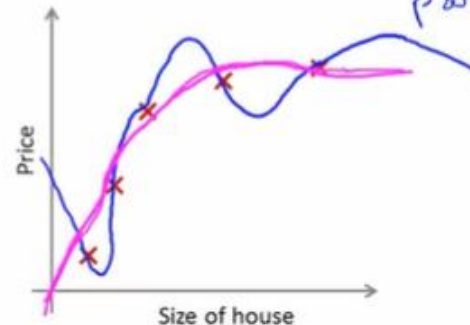$\theta_1, \theta_2, \theta_3, \ldots, \theta_{100}$

**Regularization.**

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

regularization parameter

$$\min_{\theta} J(\theta)$$

Price / Size of house

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

$$J(\theta) = \frac{1}{2m}\left[\sum_{i=1}^{m}(h_\theta(x^{(i)}) - y^{(i)})^2 + \lambda\sum_{j=1}^{n}\theta_j^2\right]$$
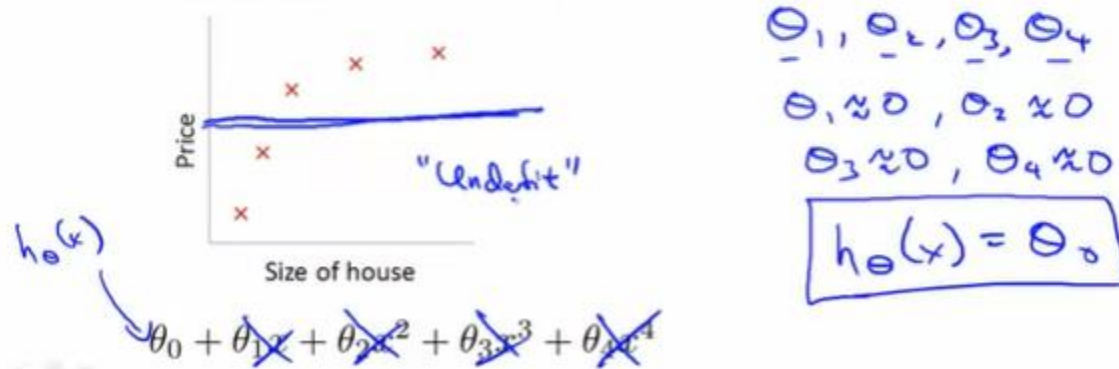
What if $\lambda$ is set to an extremely large value (perhaps too large for our problem, say $\lambda = 10^{10}$)?

○ Algorithm works fine; setting $\lambda$ to be very large can't hurt it.

○ Algorithm fails to eliminate overfitting.

◉ Algorithm results in underfitting (fails to fit even the training set).

○ Gradient descent will fail to converge.

**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 $\lambda = 10^{10}$)?**



Andrew Ng

# Regularized Linear Regression

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

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

**Gradient descent**

$\Theta_0$    $\Theta_1, \Theta_2 \ldots, \Theta_n$

Repeat {

$$\theta_0 := \theta_0 - \alpha \frac{1}{m} \sum_{i=1}^{m} (h_\theta(x^{(i)}) - y^{(i)}) x_0^{(i)}$$

We don't penalize theta0

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

$$(j = \cancel{0}, 1, 2, 3, \ldots, n)$$

}

---

**Gradient descent**

$\Theta_0$    $\Theta_1, \Theta_2 \ldots, \Theta_n$

$\frac{\partial}{\partial \Theta_0} J(\Theta)$

Repeat {

$$\theta_0 := \theta_0 - \alpha \frac{1}{m} \sum_{i=1}^{m} (h_\theta(x^{(i)}) - y^{(i)}) x_0^{(i)}$$

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

$$(j = \cancel{0}, 1, 2, 3, \ldots, n)$$

}

$$\theta_j := \theta_j (1 - \alpha \frac{\lambda}{m}) - \alpha \frac{1}{m} \sum_{i=1}^{m} (h_\theta(x^{(i)}) - y^{(i)}) x_j^{(i)}$$
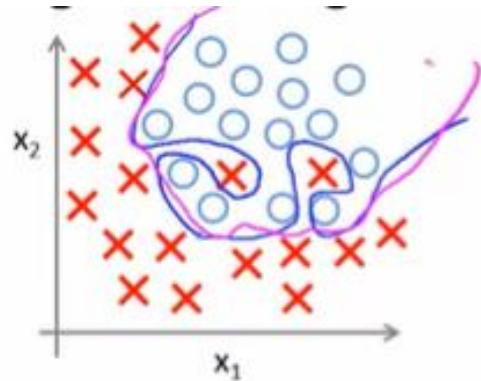
$$1 - \alpha \frac{\lambda}{m} < 1 \qquad 0.99 \qquad \Theta_j \times 0.99 \qquad \Theta_j^2$$

# 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:

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

Andrew Ng

**Reasons for Underfitting:**

1. High bias and low variance

2. The size of the training dataset used is not enough.

3. The model is too simple.

4. Training data is not cleaned and also contains noise in it.

**Techniques to reduce underfitting:**

1. Increase model complexity

2. Increase the number of features, performing feature engineering

3. Remove noise from the data.

4. Increase the number of epochs or increase the duration of training to get better results.