# Linear regression with one variable

## Model representation

Machine Learning

**Housing Prices (Portland, OR)**

Price (in 1000s of dollars) vs. Size (feet$^2$)

<u>Supervised Learning</u>

Given the "right answer" for each example in the data.

<u>Regression Problem</u>

Predict real-valued output

Andrew Ng

**Training set of housing prices (Portland, OR)**

| Size in feet² (x) | Price ($) in 1000's (y) |
|---|---|
| 2104 | 460 |
| 1416 | 232 |
| 1534 | 315 |
| 852 | 178 |
| … | … |

Notation:

   **m** = Number of training examples

   **x**'s = "input" variable / features

   **y**'s = "output" variable / "target" variable

Training Set

Learning Algorithm

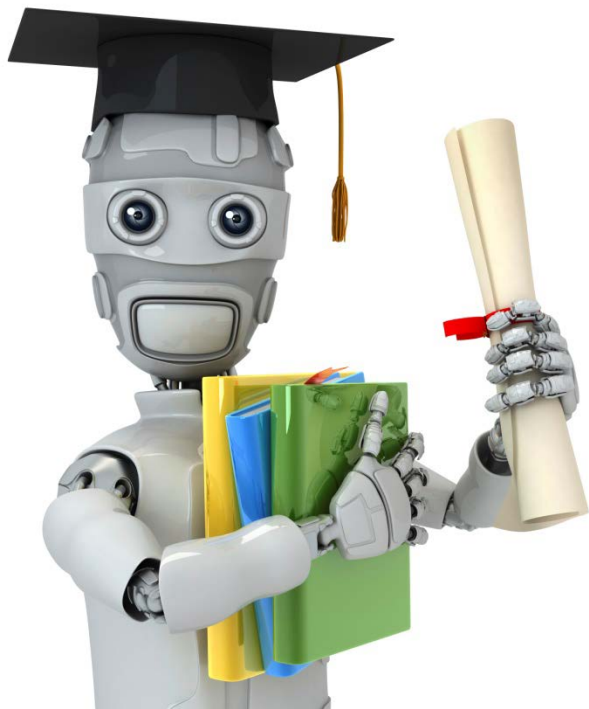Size of house → h → Estimated price

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

Linear regression with one variable.
Univariate linear regression.

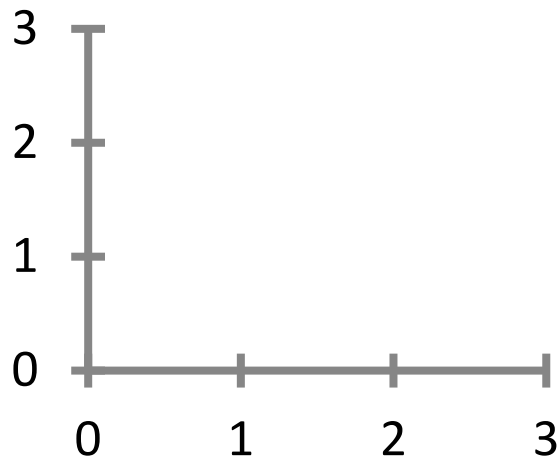Andrew Ng

Linear regression
with one variable

Cost function

Machine Learning

Training Set

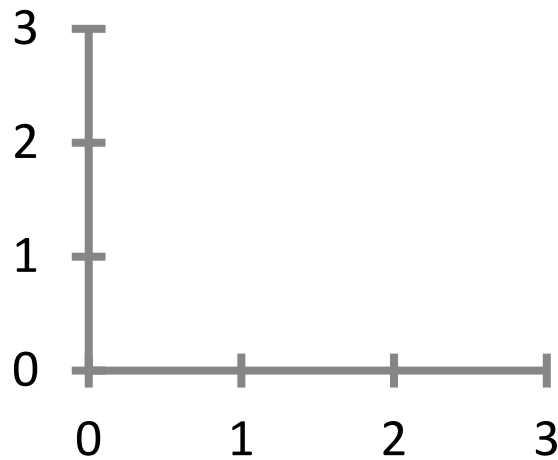| Size in feet² (x) | Price ($) in 1000's (y) |
|---|---|
| 2104 | 460 |
| 1416 | 232 |
| 1534 | 315 |
| 852 | 178 |
| … | … |

Hypothesis: $h_\theta(x) = \theta_0 + \theta_1 x$
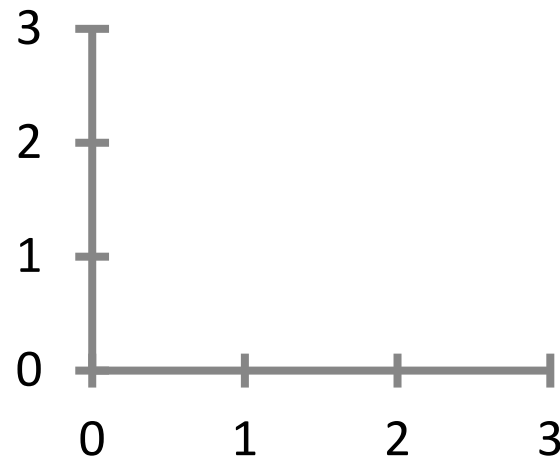
$\theta_i$'s:  Parameters

How to choose $\theta_i$'s ?

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



$\theta_0 = 1.5$
$\theta_1 = 0$

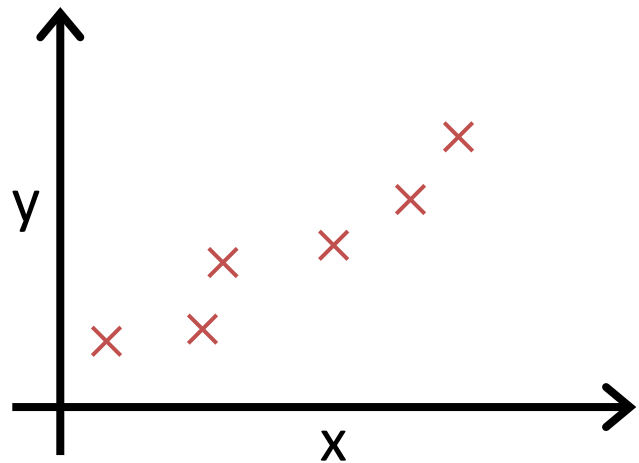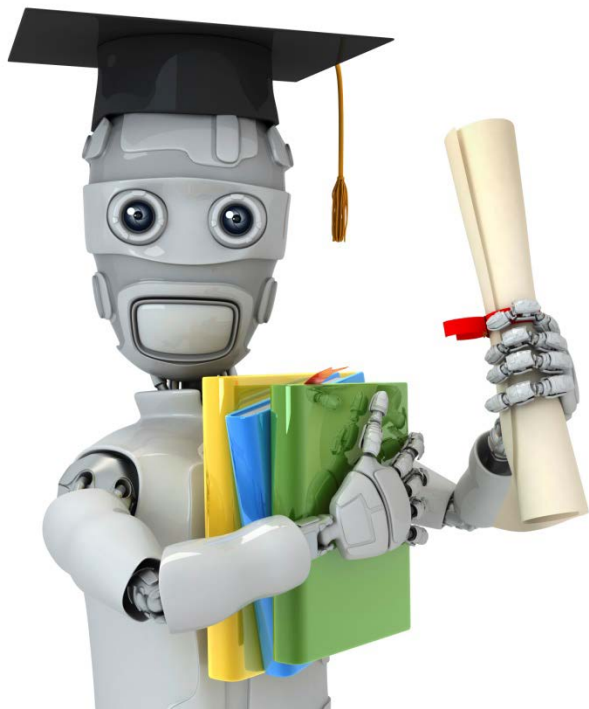$\theta_0 = 0$
$\theta_1 = 0.5$

$\theta_0 = 1$
$\theta_1 = 0.5$

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

Linear regression with one variable

Cost function intuition I

Machine Learning

**Simplified**

Hypothesis:

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

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

Parameters:

$$\theta_0, \theta_1$$

$$\theta_1$$

Cost Function:

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

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

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

$$\underset{\theta_1}{\text{minimize}} \; J(\theta_1)$$

# $h_\theta(x)$

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



$h_\theta(x)$

$\theta_1 = 1$

# $J(\theta_1)$

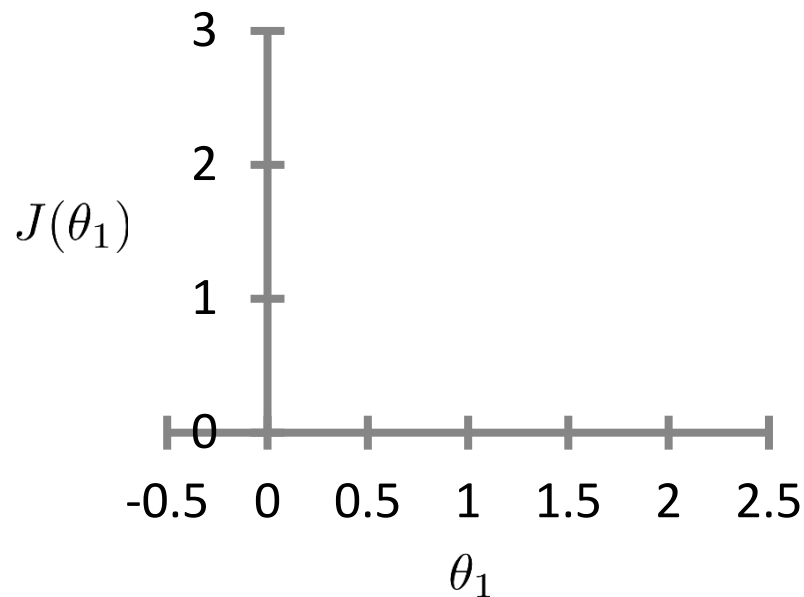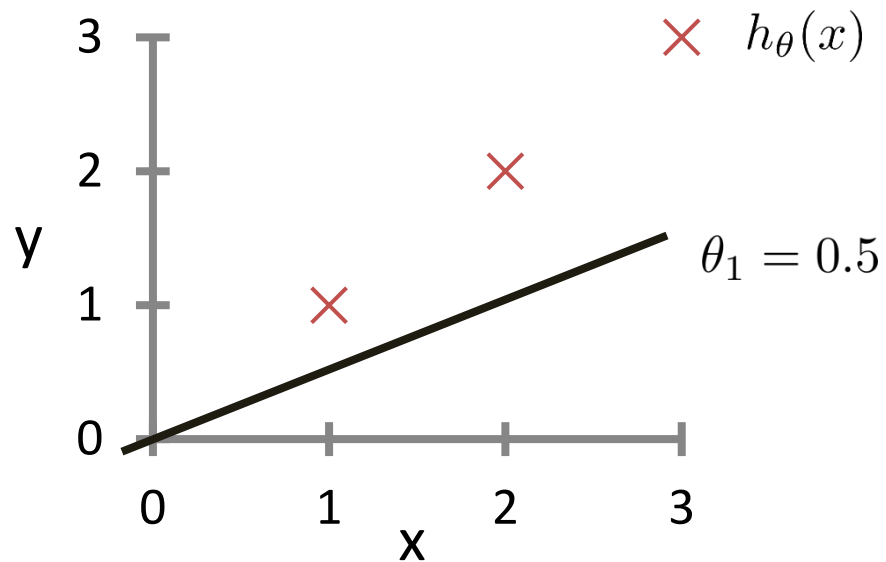(function of the parameter $\theta_1$)



$J(\theta_1)$

# $h_\theta(x)$

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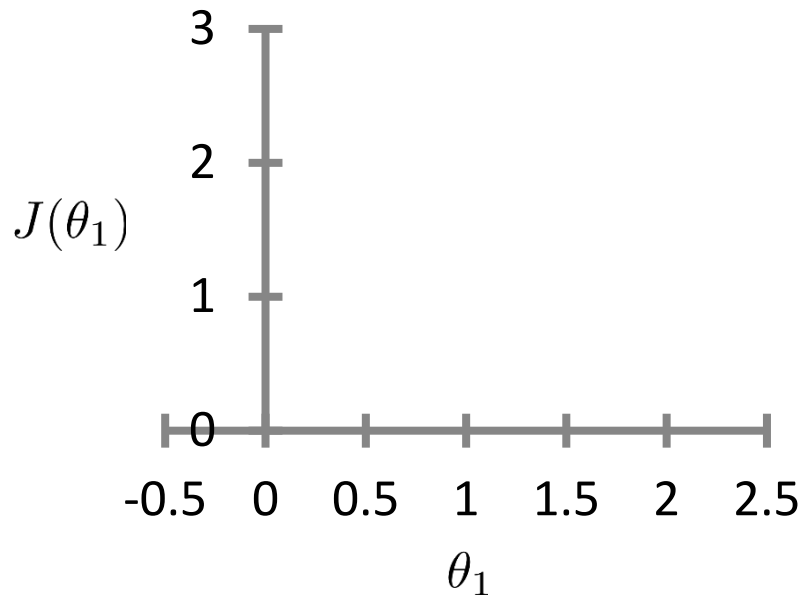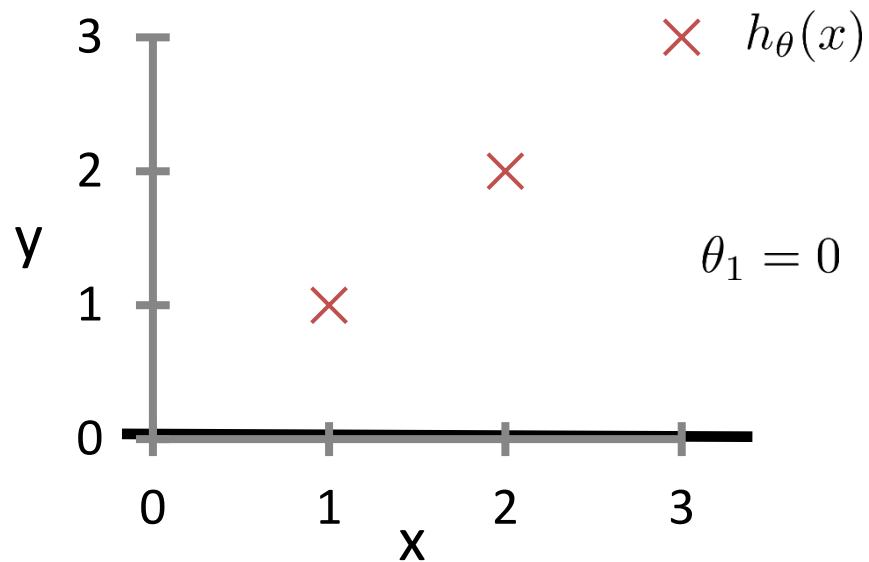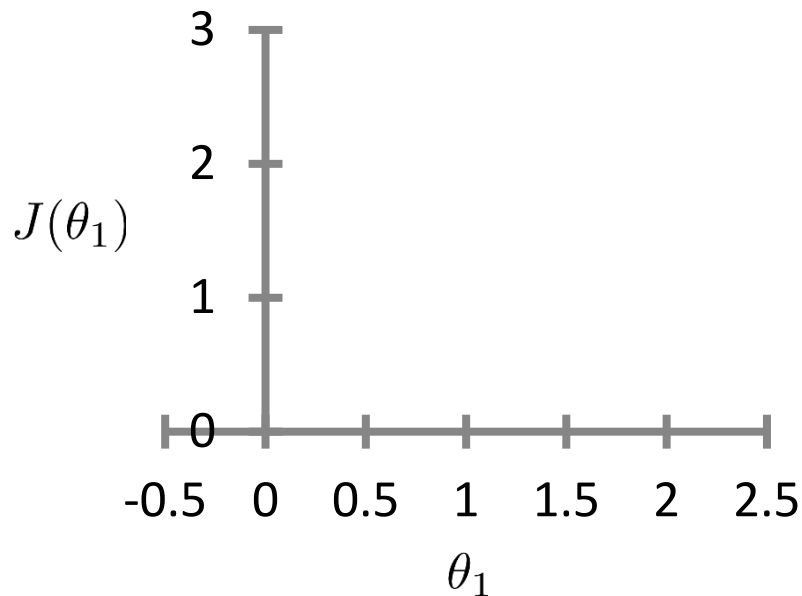
# $J(\theta_1)$

(function of the parameter $\theta_1$)

$\times$  $h_\theta(x)$

y

$\theta_1 = 0.5$

x

$J(\theta_1)$

$\theta_1$

# $h_\theta(x)$

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

# $J(\theta_1)$

(function of the parameter $\theta_1$)

$\times$  $h_\theta(x)$

$\theta_1 = 0$

Linear regression
with one variable

Cost function
intuition II

Machine Learning

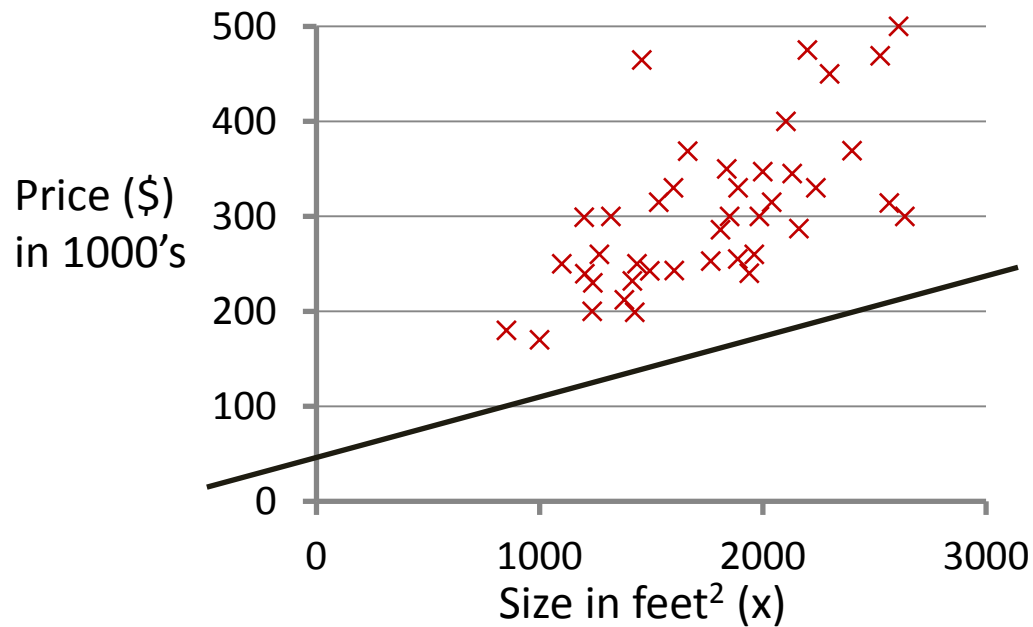Hypothesis: $\qquad h_\theta(x) = \theta_0 + \theta_1 x$

Parameters: $\qquad \theta_0, \theta_1$

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

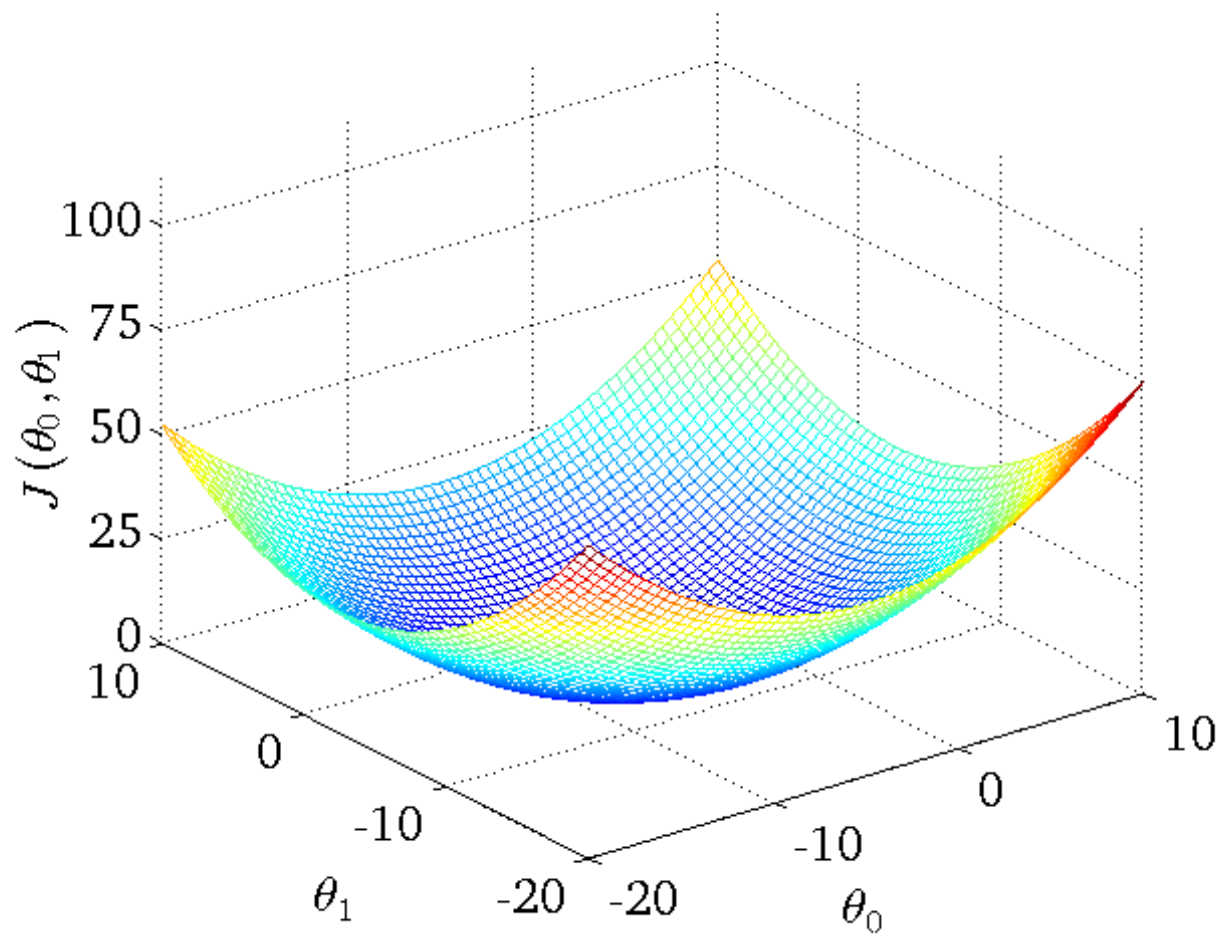Goal: $\qquad\qquad \underset{\theta_0, \theta_1}{\text{minimize}} \; J(\theta_0, \theta_1)$

# $h_\theta(x)$

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

# $J(\theta_0, \theta_1)$

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



$$h_\theta(x) = 50 + 0.06x$$

# $h_\theta(x)$

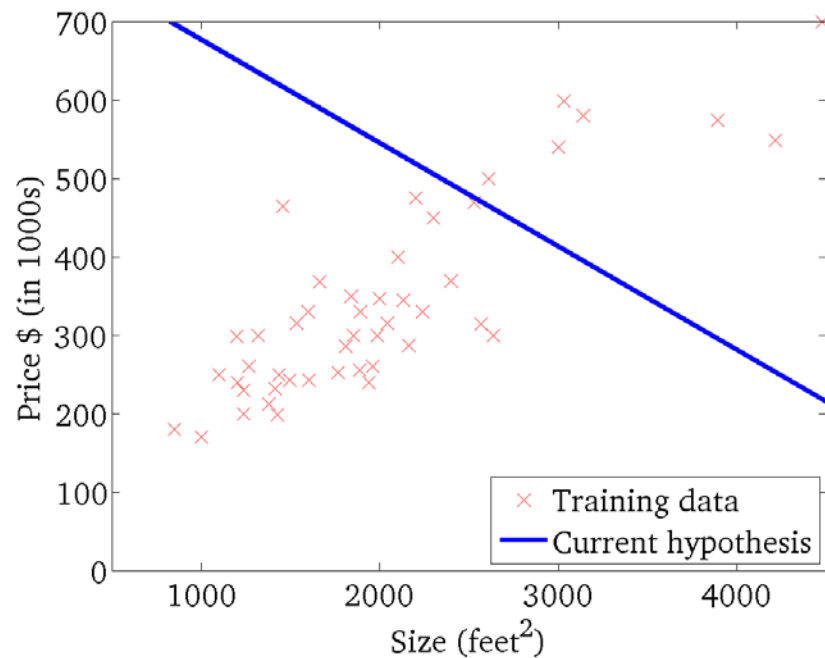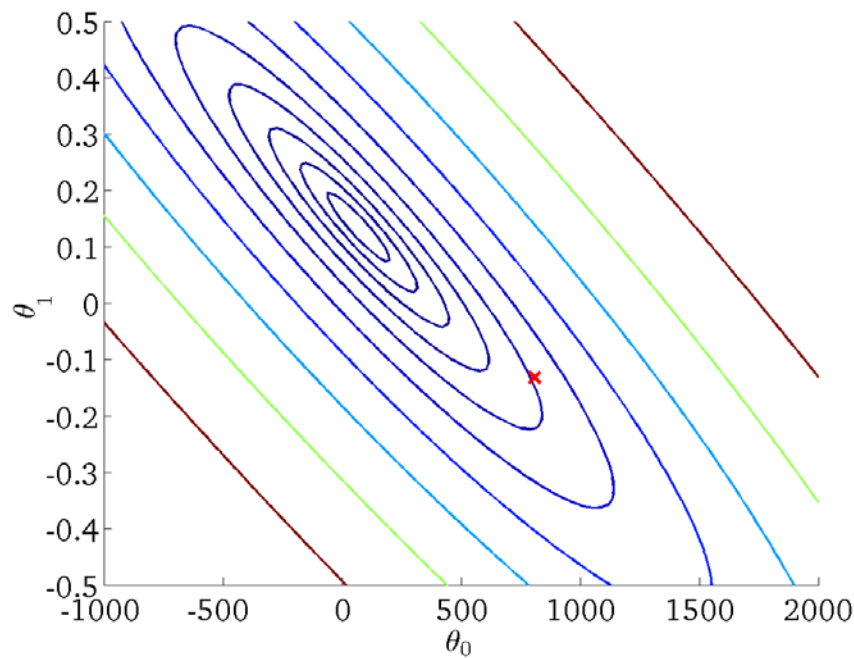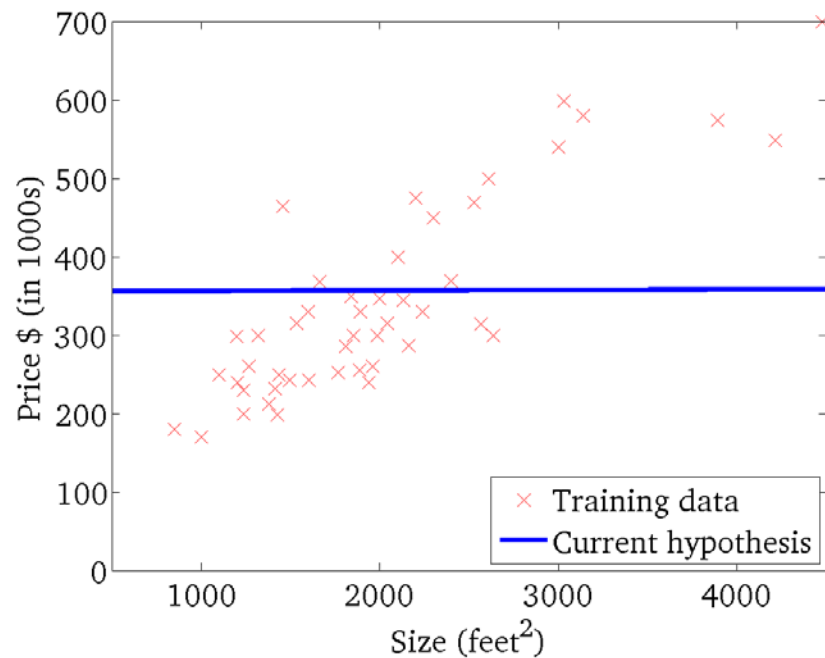(for fixed $\theta_0, \theta_1$, this is a function of x)

# $J(\theta_0, \theta_1)$

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



Andrew Ng
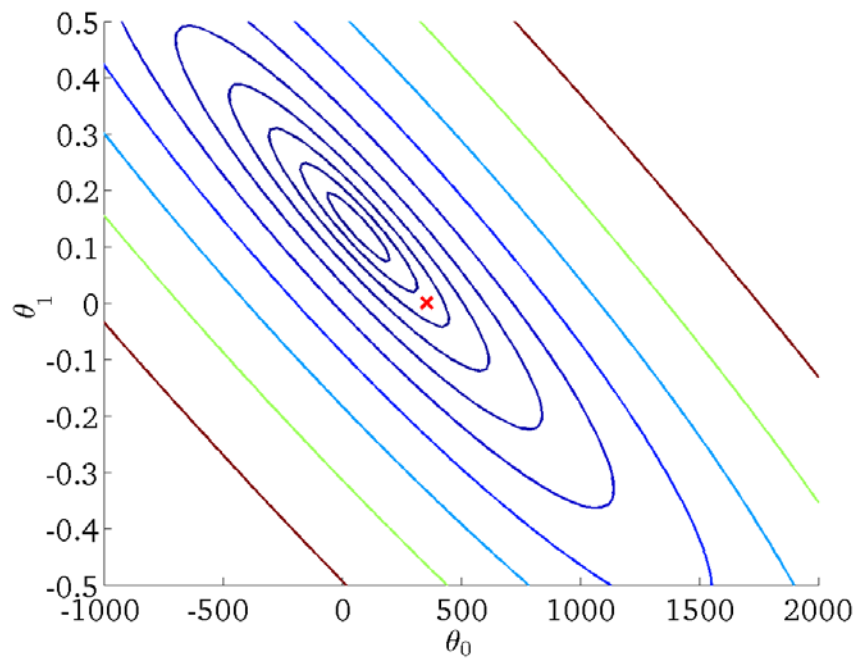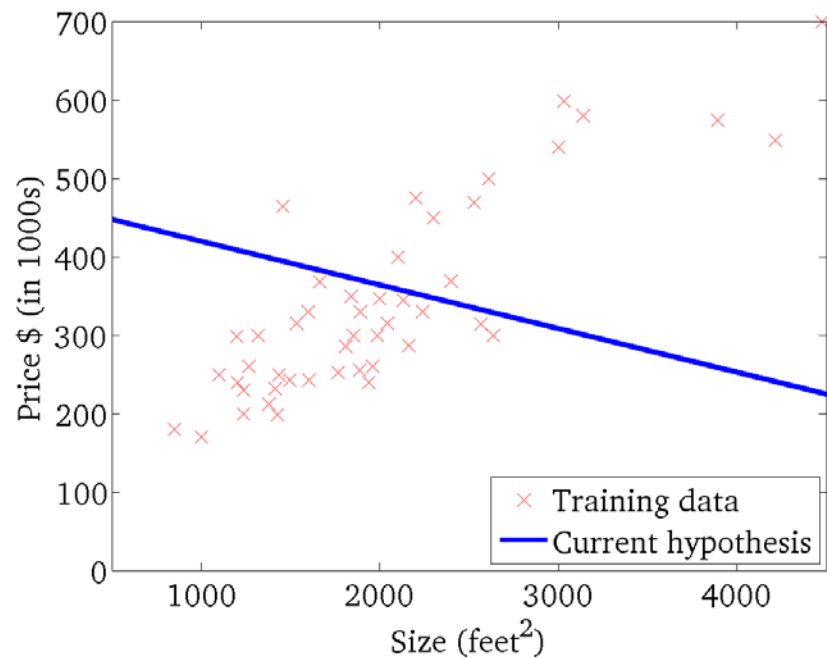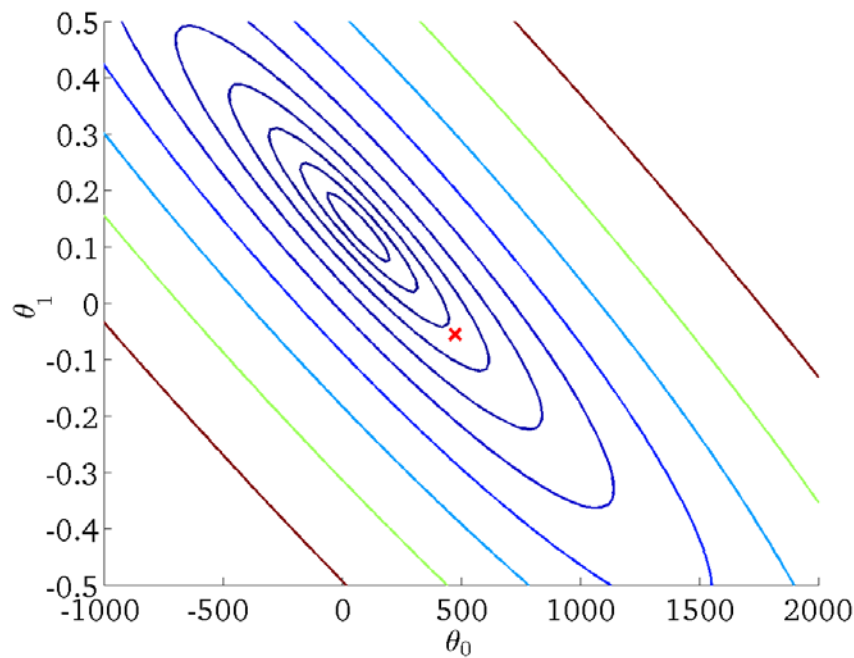
$$h_\theta(x)$$

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

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

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



Andrew Ng

# $h_\theta(x)$

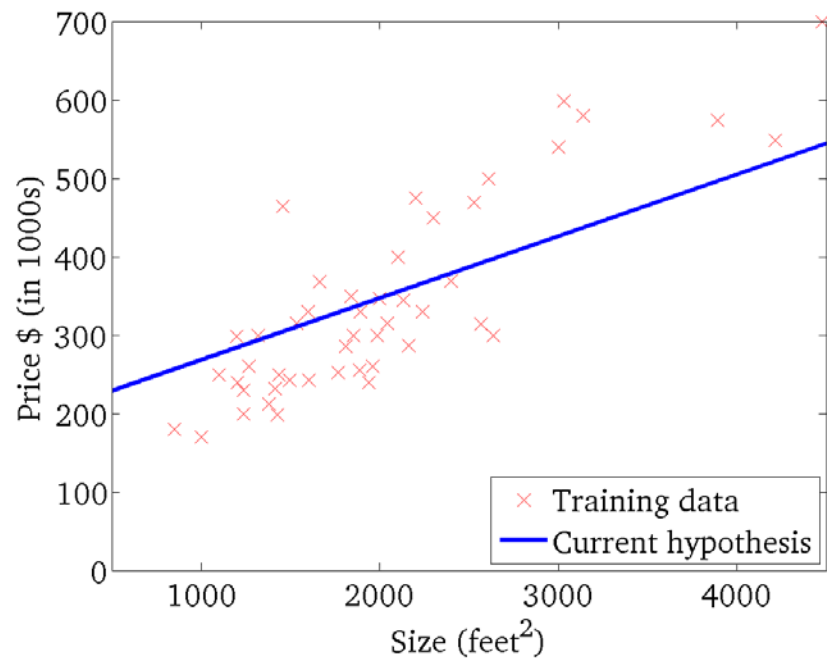## (for fixed $\theta_0, \theta_1$, this is a function of x)

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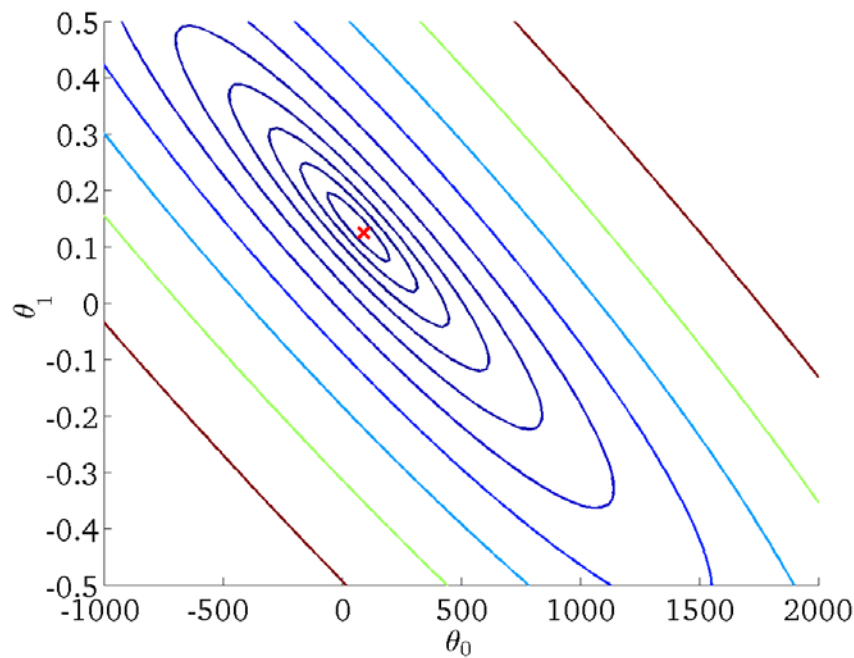
## (function of the parameters $\theta_0, \theta_1$)

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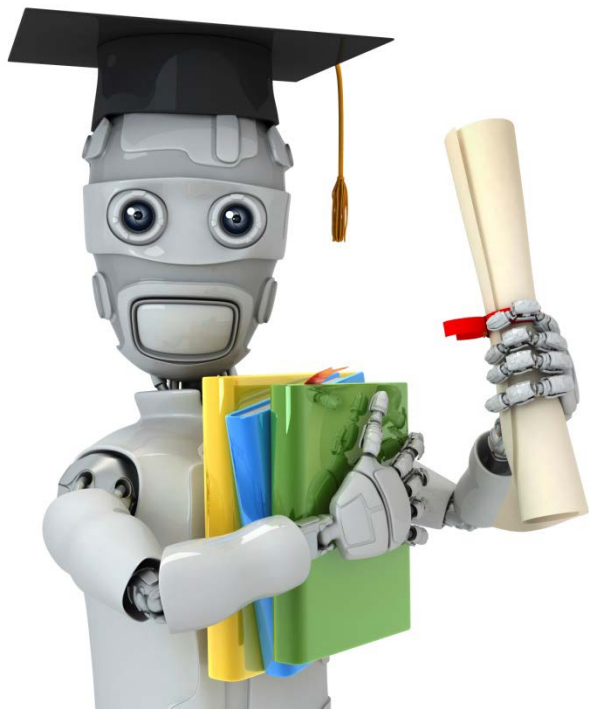
(for fixed $\theta_0, \theta_1$, this is a function of x)

$J(\theta_0, \theta_1)$

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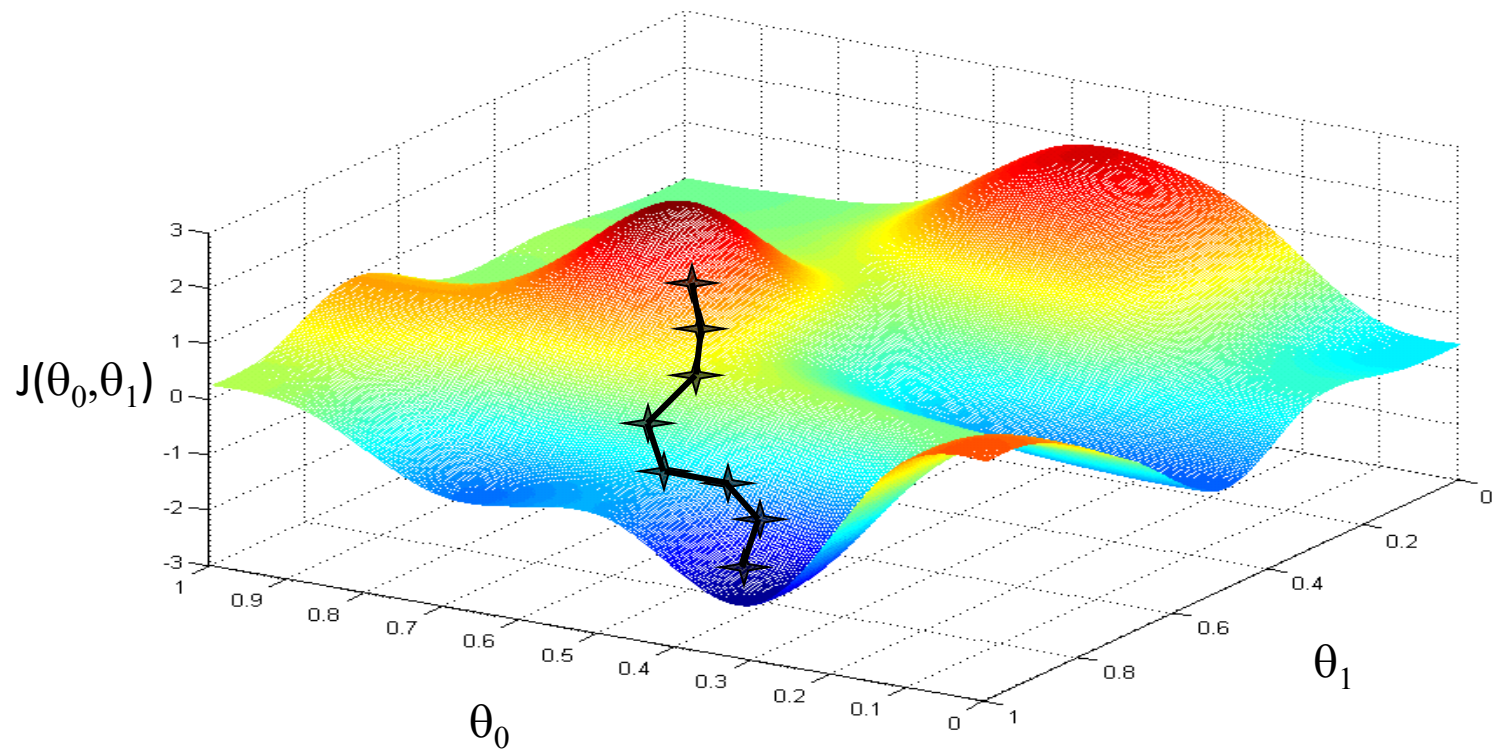
Andrew Ng

Linear regression with one variable

Gradient descent

Machine Learning
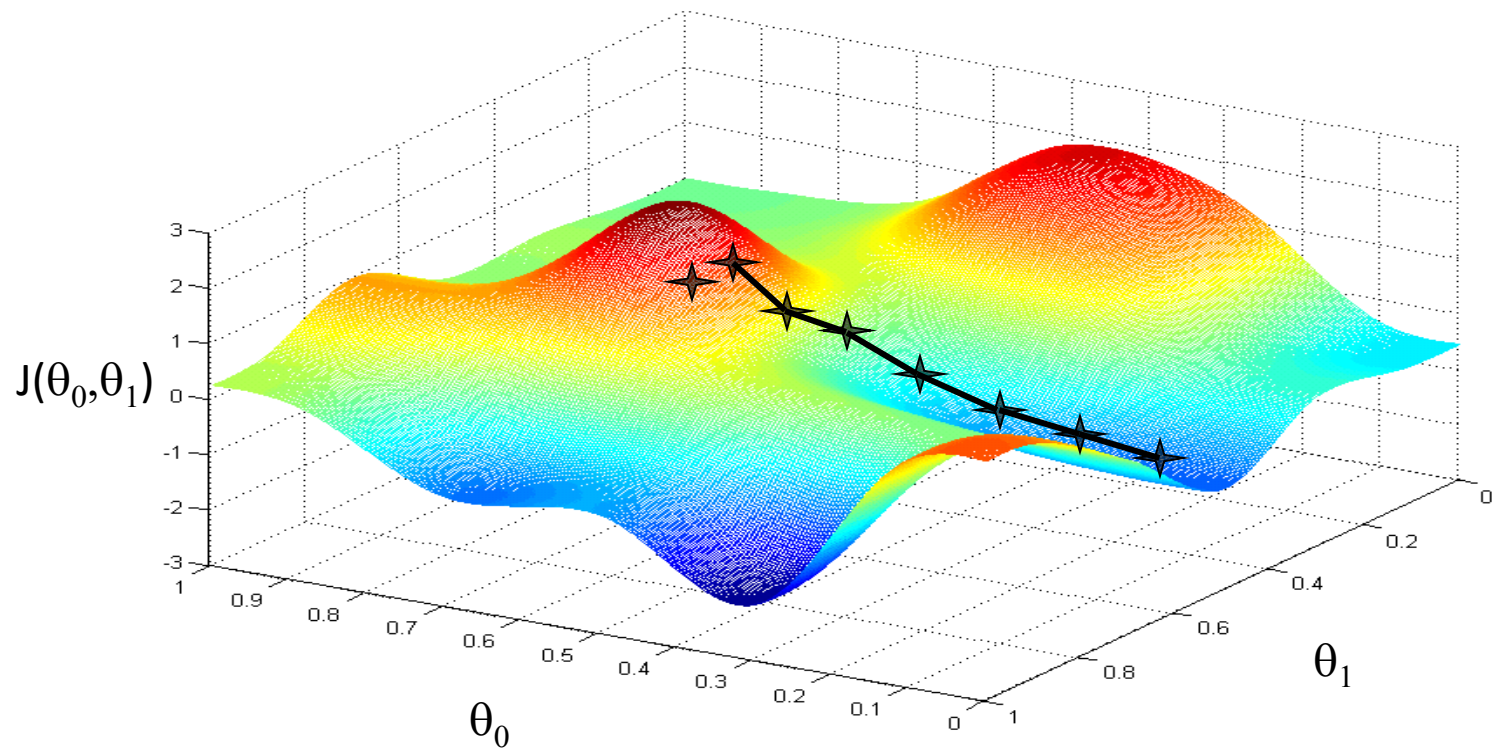
Have some function $J(\theta_0, \theta_1)$

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

**Outline:**

- Start with some $\theta_0, \theta_1$

- Keep changing $\theta_0, \theta_1$ to reduce $J(\theta_0, \theta_1)$

  until we hopefully end up at a minimum

$J(\theta_0, \theta_1)$

$\theta_0$

$\theta_1$

Andrew Ng

$J(\theta_0, \theta_1)$

$\theta_0$

$\theta_1$

Andrew Ng

# Gradient descent algorithm

repeat until convergence {

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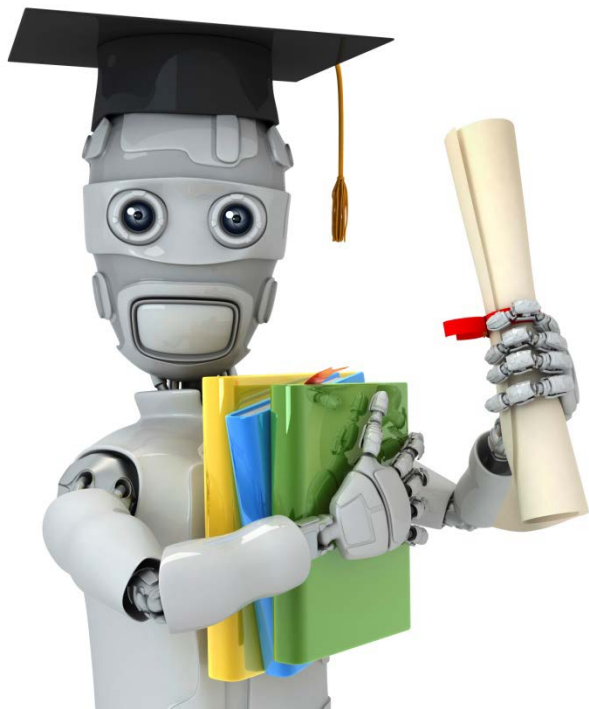
}

---

**Correct: Simultaneous update**

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

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

$\theta_0 := \text{temp0}$

$\theta_1 := \text{temp1}$

**Incorrect:**

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

$\theta_0 := \text{temp0}$

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

$\theta_1 := \text{temp1}$

Andrew Ng

Linear regression with one variable

Gradient descent intuition

Machine Learning
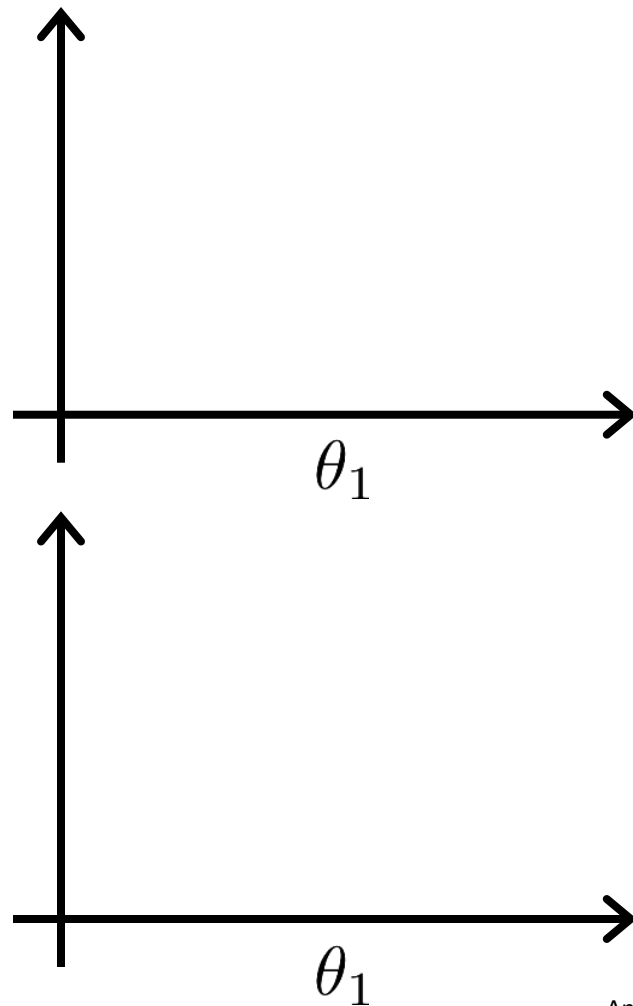
# Gradient descent algorithm

repeat until convergence {

$$\theta_j := \theta_j - \alpha \frac{\partial}{\partial \theta_j} J(\theta_0, \theta_1)$$
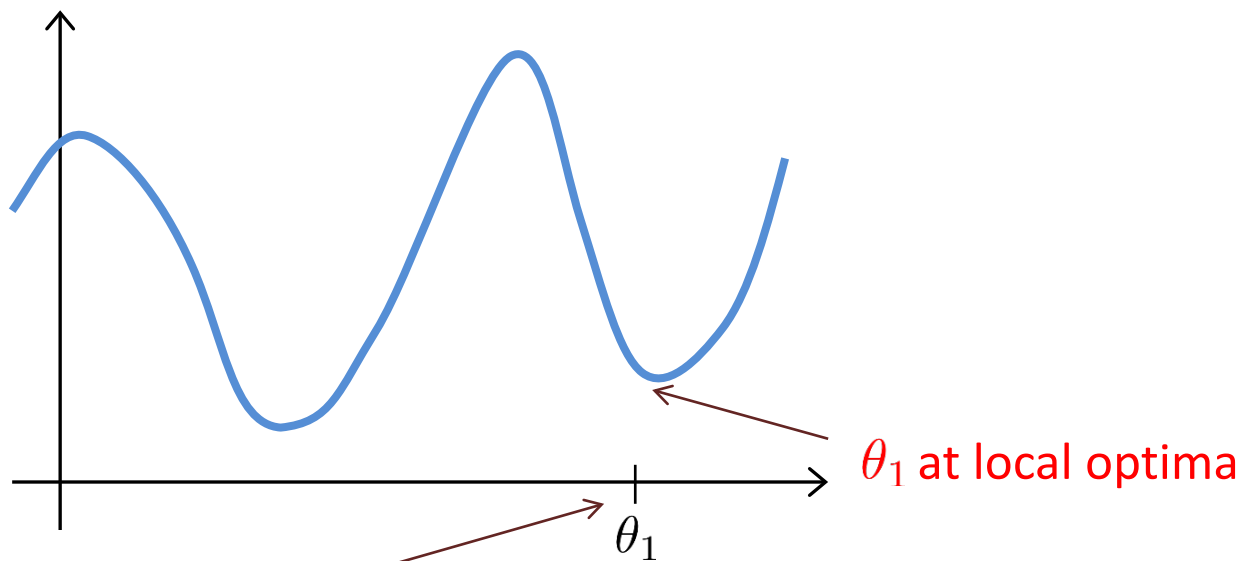
(simultaneously update $j = 0$ and $j = 1$)

}

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

If α is too small, gradient descent can be slow.

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



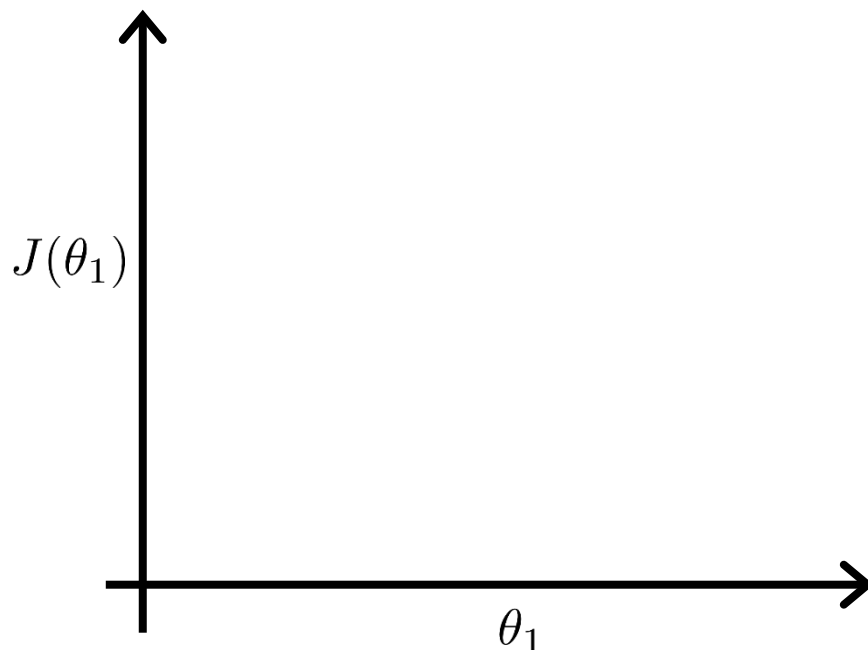$\theta_1$

$\theta_1$

Andrew Ng

$\theta_1$ at local optima

Current value of $\theta_1$

$\theta_1$

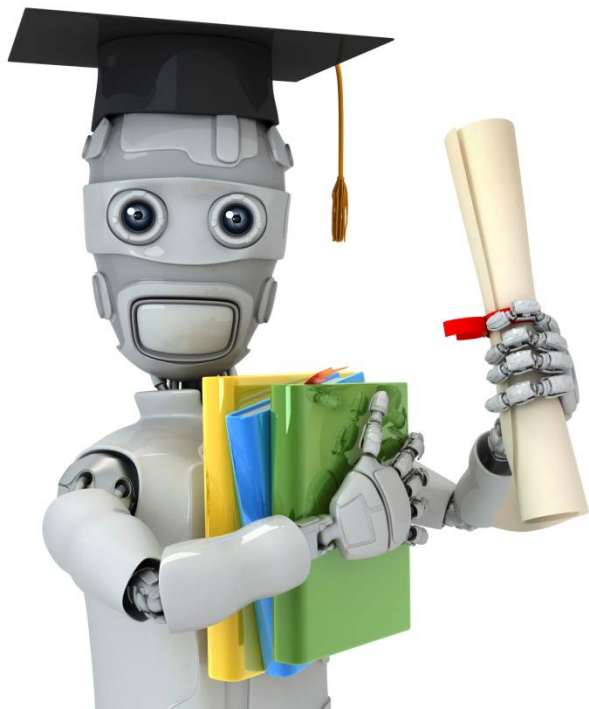$$\theta_1 := \theta_1 - \alpha \frac{d}{d\theta_1} J(\theta_1)$$

Andrew Ng

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.

$J(\theta_1)$

$\theta_1$

Andrew Ng

Linear regression
with one variable

Gradient descent for
linear regression

Machine Learning

# Gradient descent algorithm

repeat until convergence {

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

$$\text{(for } j = 1 \text{ and } j = 0)$$

}

# Linear Regression Model

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

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

Andrew Ng

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

$$j = 0 : \frac{\partial}{\partial \theta_0} J(\theta_0, \theta_1) =$$

$$j = 1 : \frac{\partial}{\partial \theta_1} J(\theta_0, \theta_1) =$$

# Gradient descent algorithm

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

$J(\theta_0, \theta_1)$

$\theta_0$

$\theta_1$

Andrew Ng

$h_\theta(x)$

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

$J(\theta_0, \theta_1)$

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

Andrew Ng

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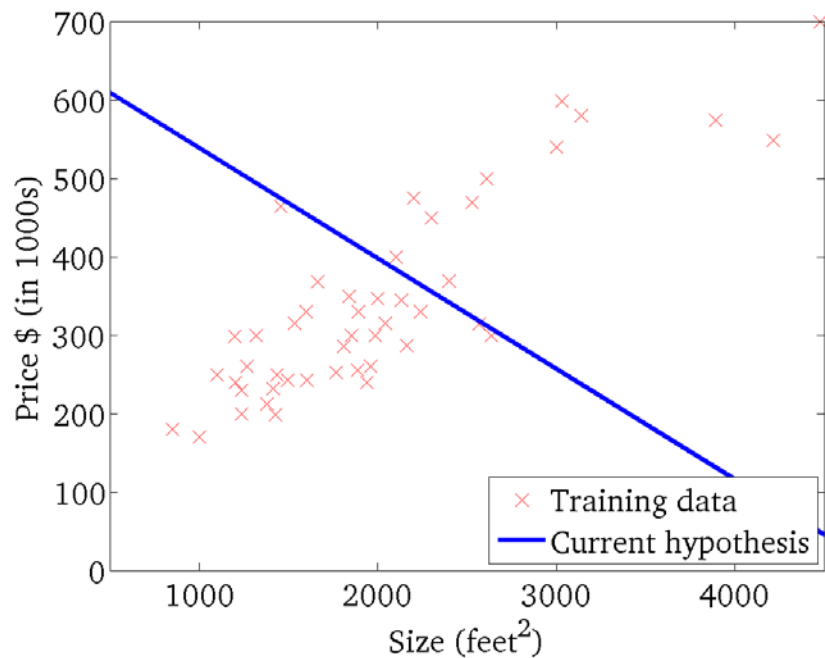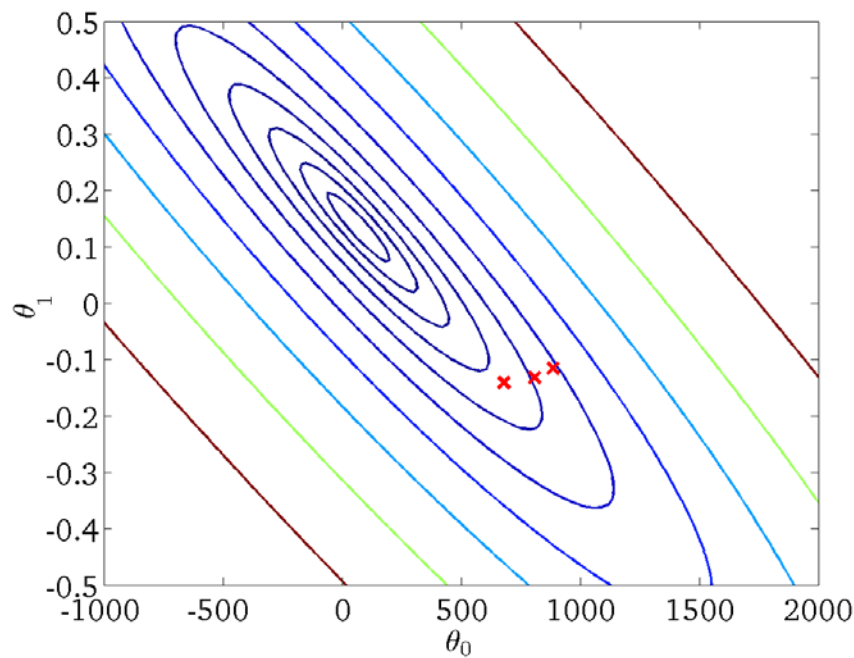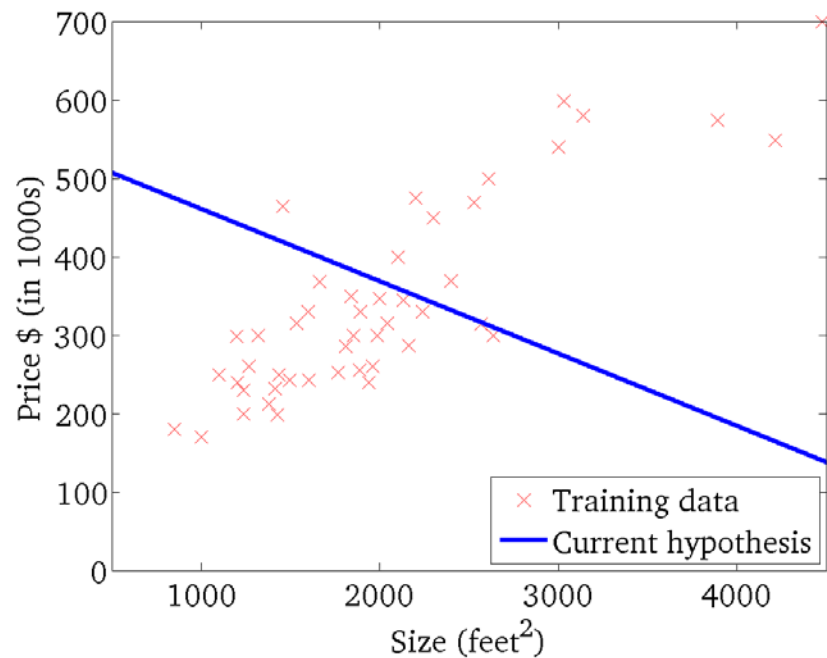
(for fixed $\theta_0, \theta_1$, this is a function of x)

# $J(\theta_0, \theta_1)$

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

$h_\theta(x)$

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

$J(\theta_0, \theta_1)$

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

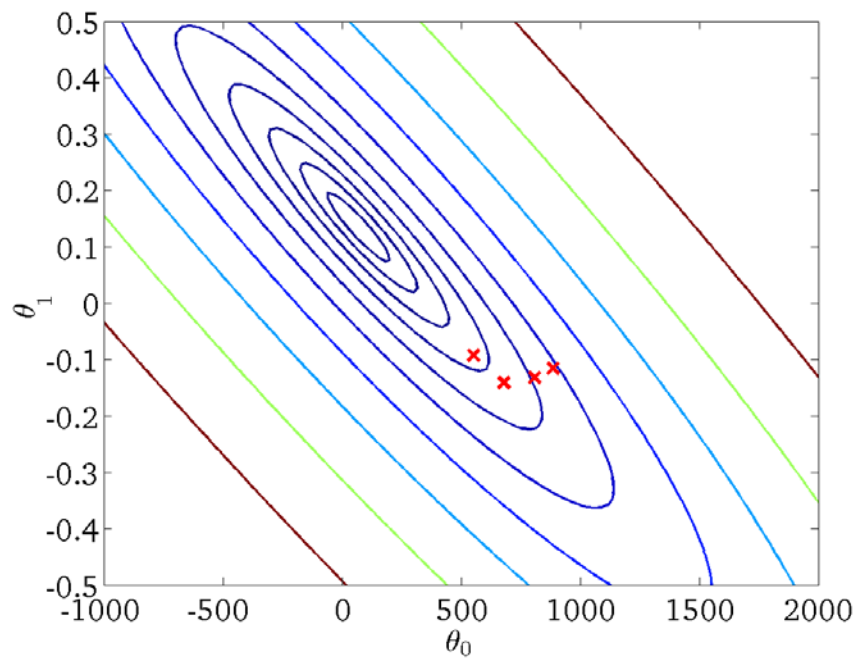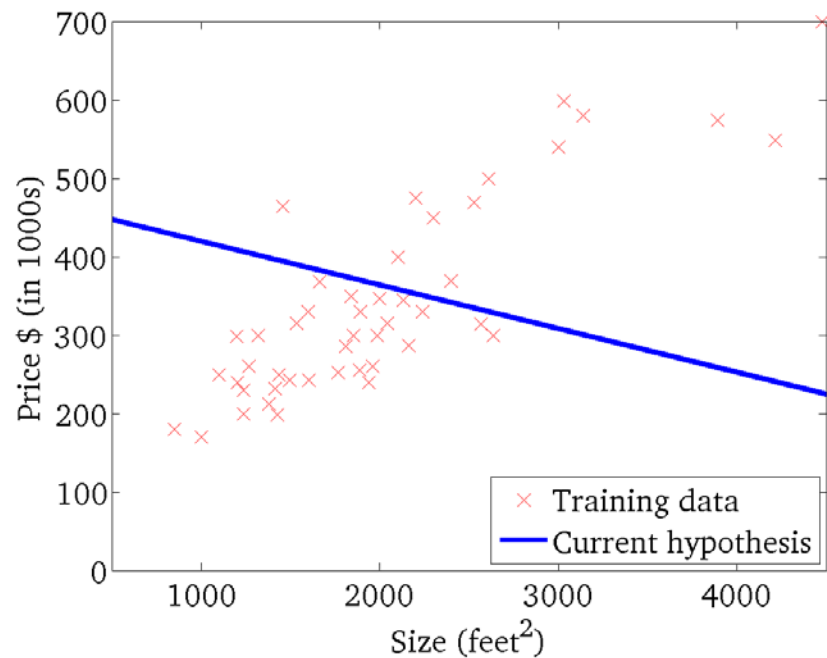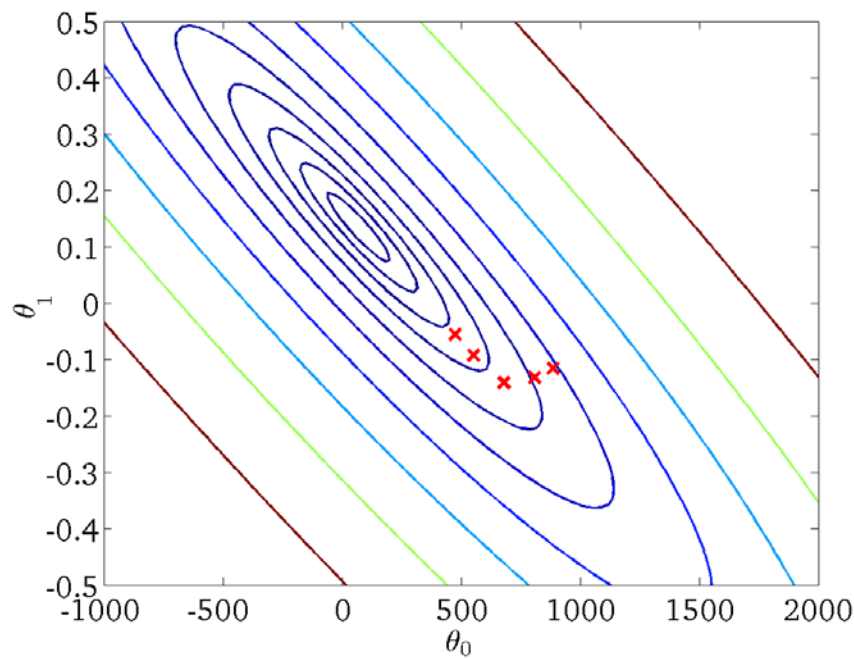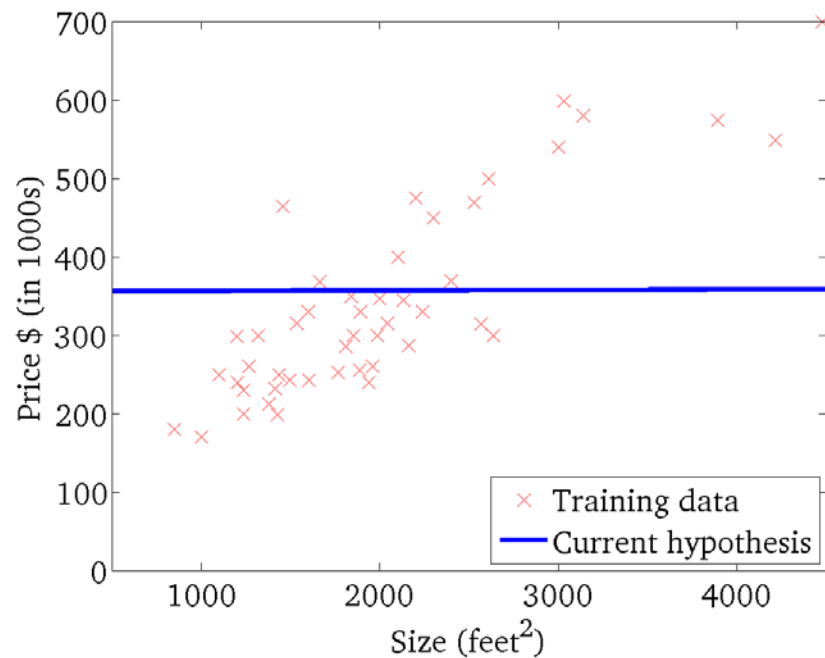# $h_\theta(x)$

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

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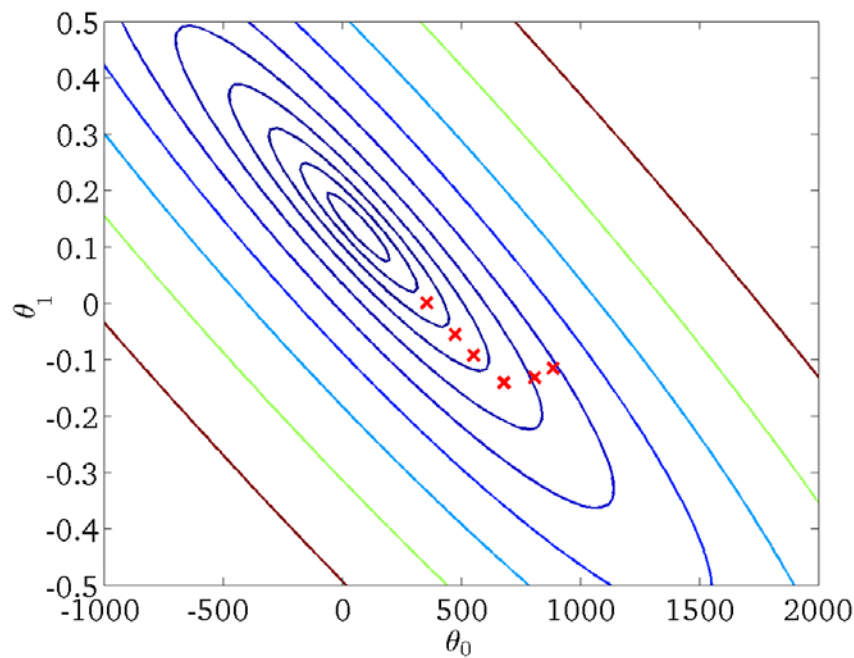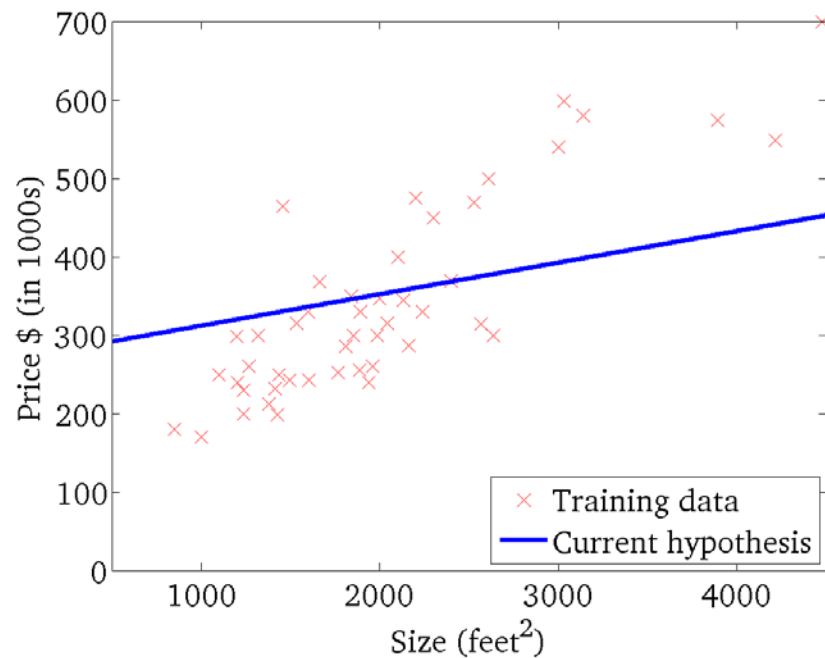
(function of the parameters $\theta_0, \theta_1$)

$h_\theta(x)$

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

$J(\theta_0, \theta_1)$

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



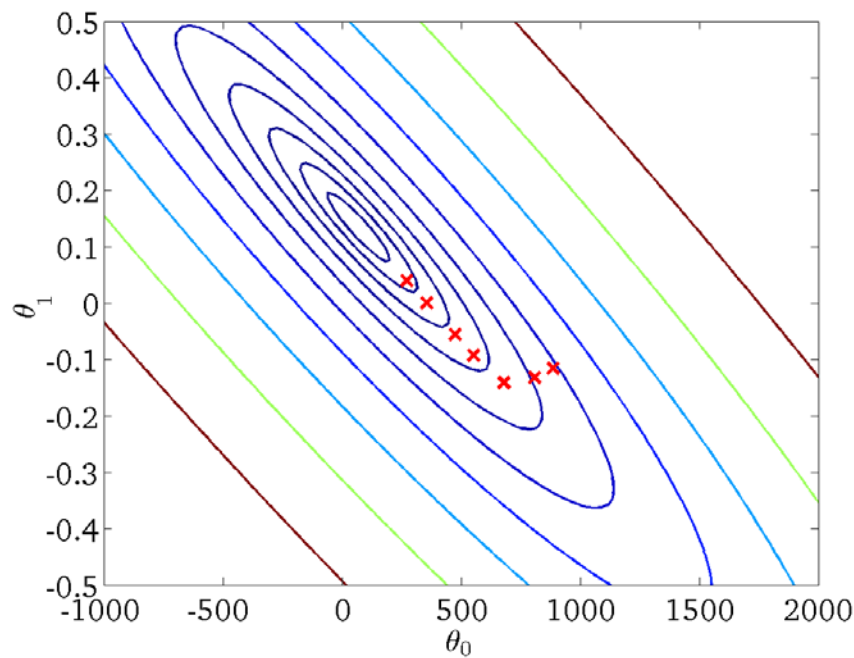Andrew Ng

# $h_\theta(x)$

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

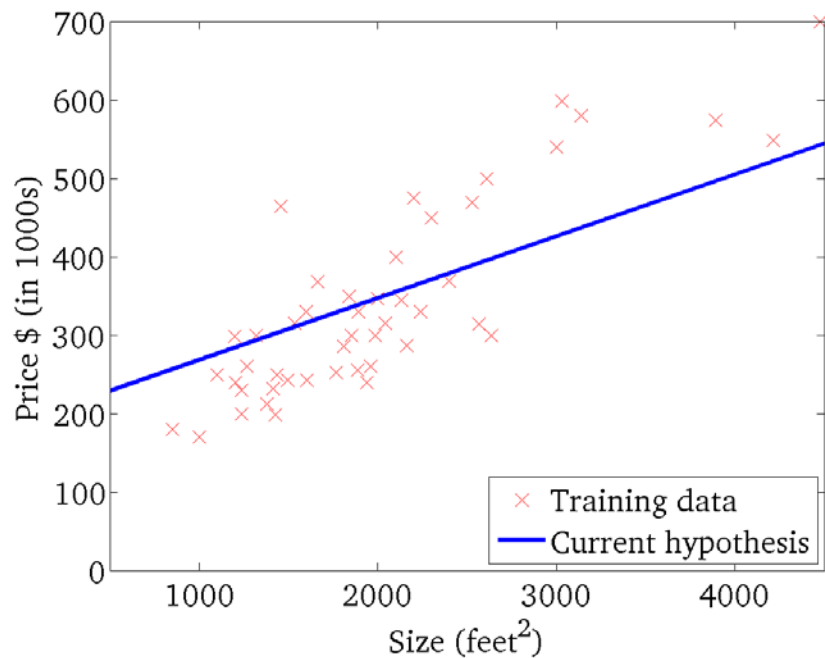# $J(\theta_0, \theta_1)$

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

$h_\theta(x)$

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

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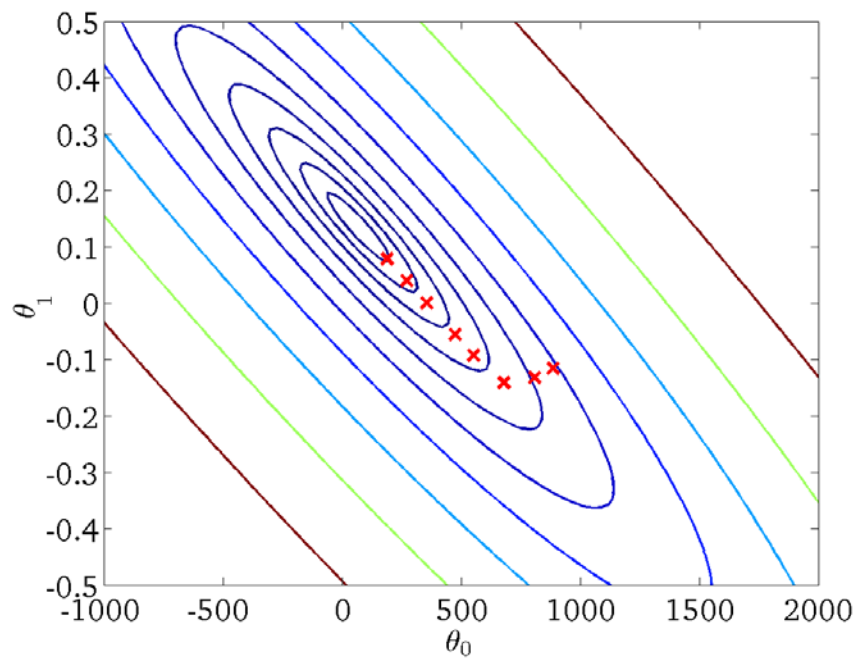
(function of the parameters $\theta_0, \theta_1$)

Andrew Ng

$h_\theta(x)$

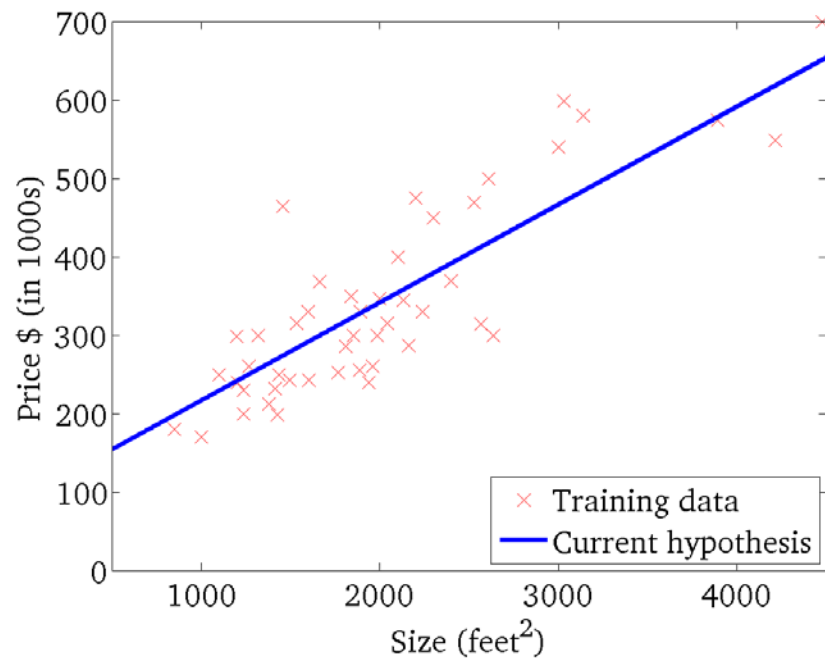(for fixed $\theta_0, \theta_1$, this is a function of x)

$J(\theta_0, \theta_1)$

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

Andrew Ng

$$h_\theta(x)$$

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

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

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



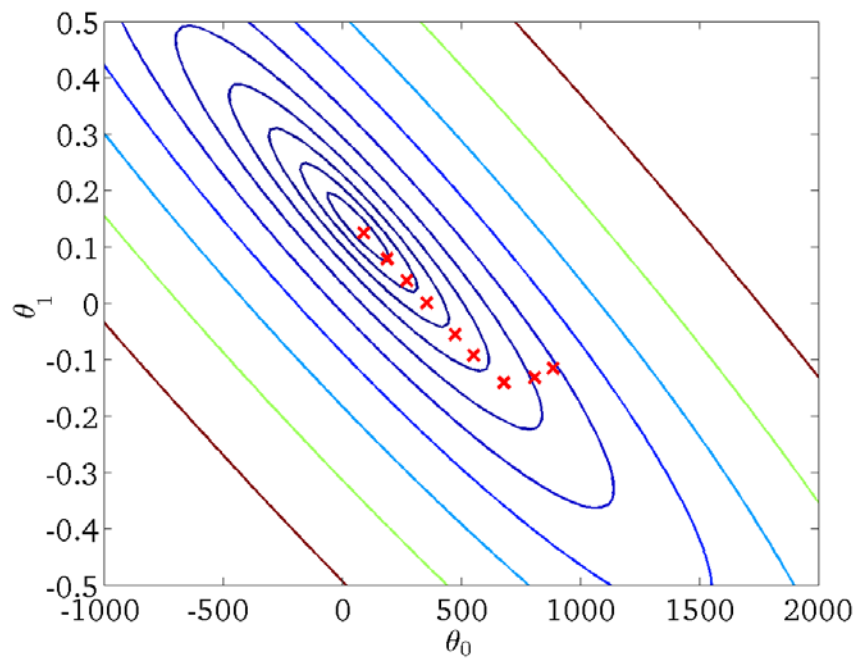Andrew Ng

**"Batch" Gradient Descent**

"Batch": Each step of gradient descent uses all the training examples.