# Linear regression with one variable

## Model representation

Machine Learning

**Housing Prices (Portland, OR)**



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

220k

1250

**Supervised Learning**

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

**Regression Problem**

Predict real-valued output

Classification: Discrete-valued output
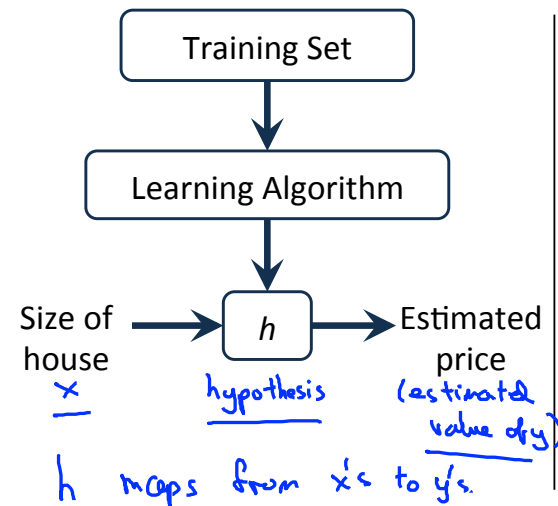
---

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

| Size in feet$^2$ (x) | Price ($) in 1000's (y) |
| --- | --- |
| 2104 | 460 |
| 1416 | 232 |
| 1534 | 315 |
| 852 | 178 |
| ... | ... |

m = 47

Notation:
- m = Number of training examples
- x's = "input" variable / features
- y's = "output" variable / "target" variable

$(x, y)$ – one training example

$(x^{(i)}, y^{(i)})$ – $i^{th}$ training example

$x^{(1)} = 2104$

$x^{(2)} = 1416$

$y^{(1)} = 460$

---

Training Set

↓

Learning Algorithm

↓

Size of house → $h$ → Estimated price

x

hypothesis

(estimated value of y)

$h$ 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.

one variable

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 |
| ... | ... |

$m = 47$
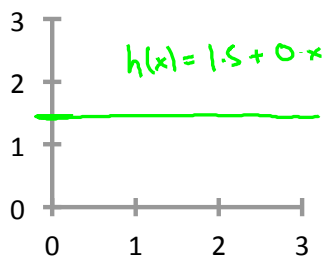
Hypothesis: $h_\theta(x) = \theta_0 + \theta_1 x$

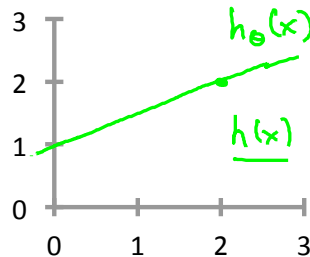$\theta_i$'s:    Parameters

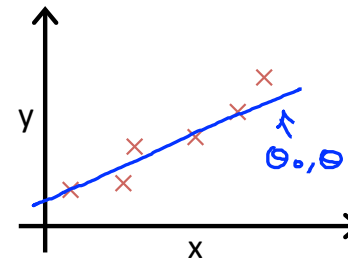How to choose $\theta_i$'s ?

---

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

$h(x) = 1.5 + 0 \cdot x$

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

$h(x) = 0.5x$

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

$h_\theta(x)$

$h(x)$

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

$\theta_0, \theta_1$

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

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

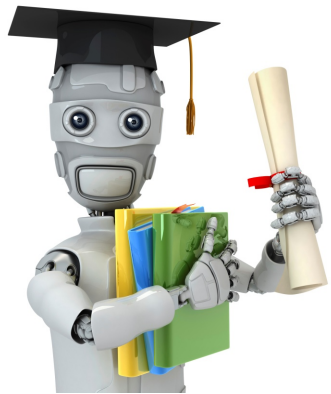$\underset{\theta_0, \theta_1}{\text{minimize}} \; \dfrac{1}{2m} \sum_{i=1}^{m} \left( h_\theta(x^{(i)}) - y^{(i)} \right)^2$

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

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

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

Cost function

Squared error function

# Linear regression with one variable

## Cost function intuition I

Machine Learning

Andrew Ng

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

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

$\theta_0 = 0$

$\theta_1$

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

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

$\theta_1 x^{(i)}$

Andrew Ng

---

$h_\theta(x)$

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

$h_\theta(x)$

$\theta_1 = 1$

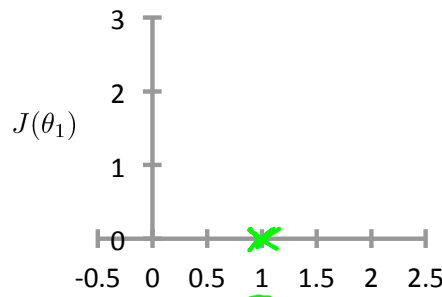$h_\theta(x^{(i)}) = y^{(i)}$

$\theta_1 = 1$

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

$J(\theta_1)$

(function of the parameter $\theta_1$)

$J(\theta_1)$

$\theta_1 = 0.5?$

$\theta_1$

$J(1) = 0$

$h_\theta(x)$

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

$h_\theta(x)$

$\theta_1 = 0.5$
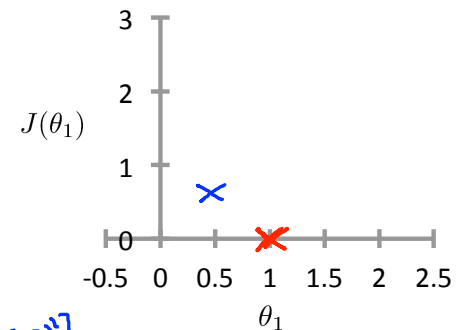
$h_\theta(x^{(i)})$

$J(0.5) = \frac{1}{2m}\left[(0.5-1)^2 + (1-2)^2 + (1.5-3)^2\right]$
$= \frac{1}{2 \times 3}(3.5) = \frac{3.5}{6} \approx 0.58$
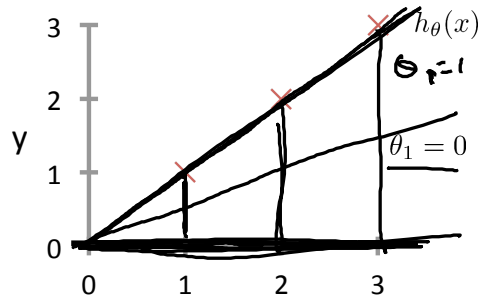
$J(\theta_1)$

(function of the parameter $\theta_1$)

$J(\theta_1)$

$\theta_1$

$\theta_1 = 0?$
$J(0) = ?$

Andrew Ng

## Slide 1

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

$h_\theta(x)$
$\theta_1 = 1$
$\theta_1 = 0$

$J(\theta_1)$
(function of the parameter $\theta_1$)

$J(\theta_1)$

$\theta_1 = 1$

$J(\theta) = \frac{1}{2m}(1^2 + 2^2 + 3^2)$
$= \frac{1}{6}\cdot 14 \approx 2.3$

$h(x) = -0.5x$

minimize $J(\theta_1)$
$\theta_1$

Linear regression
with one variable

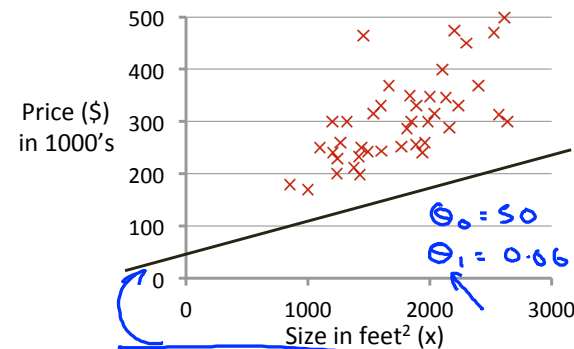Cost function
intuition II

Machine Learning

## Slide 2

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

Parameters:   $\theta_0, \theta_1$

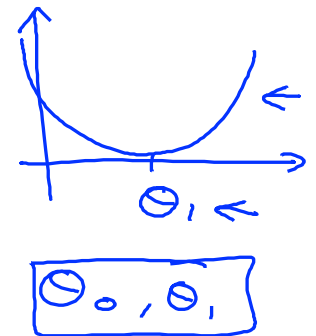Cost Function:   $J(\theta_0,\theta_1) = \frac{1}{2m}\sum_{i=1}^{m}\left(h_\theta(x^{(i)}) - y^{(i)}\right)^2$

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

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

Price ($) in 1000's

Size in feet² (x)

$\theta_0 = 50$
$\theta_1 = 0.06$
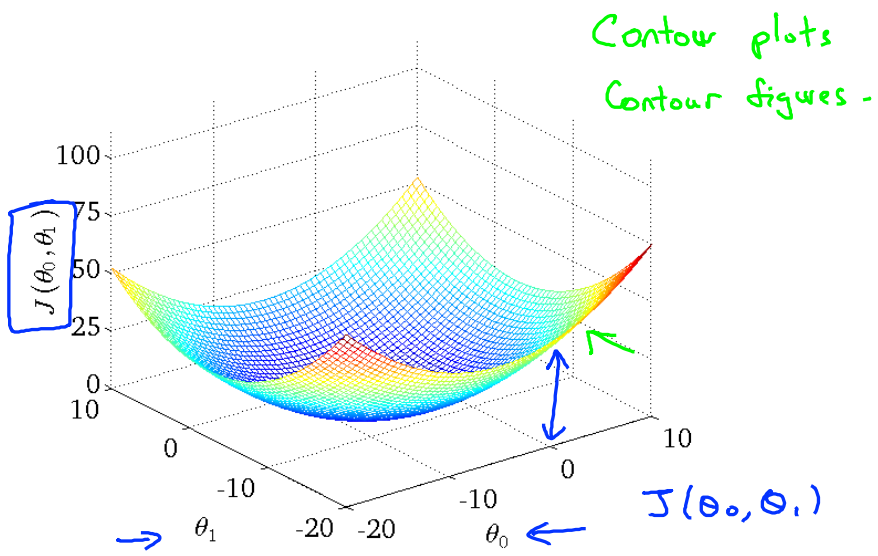
$h_\theta(x) = 50 + 0.06x$

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

$\theta_1$
$\theta_0, \theta_1$

Contour plots
Contour figures -

$J(\theta_0, \theta_1)$



$J(\theta_0, \theta_1)$

$\theta_1$  →

$\theta_0$  ←

$J(\theta_0, \theta_1)$

$h_\theta(x)$

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

Price \$ (in 1000s)

Size (feet$^2$)

× Training data
— Current hypothesis

h(x)

$J(\theta_0, \theta_1)$

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

$\theta_1$

$\theta_0$

$J(\theta_0, \theta_1)$

$\theta_0, \theta_1$
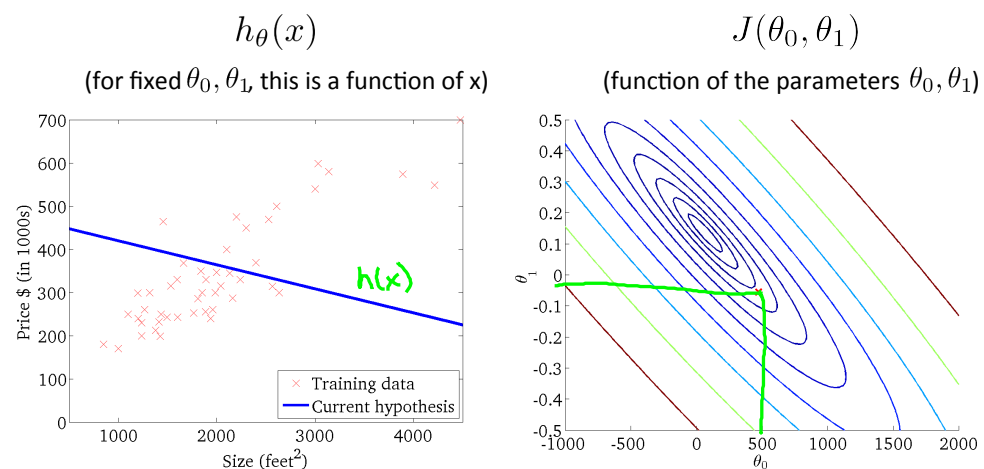
Andrew Ng

Andrew Ng

$h_\theta(x)$

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

Price \$ (in 1000s)

Size (feet$^2$)

h(x): 360 + 0·x

× Training data
— Current hypothesis

$\theta_0 = 360$
$\theta_1 = 0$

$J(\theta_0, \theta_1)$

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

$\theta_1$

$\theta_0$

$h_\theta(x)$

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

Price \$ (in 1000s)

Size (feet$^2$)

× Training data
— Current hypothesis

h(x)

$J(\theta_0, \theta_1)$

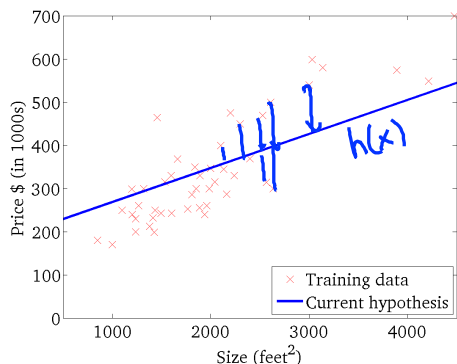(function of the parameters $\theta_0, \theta_1$)

$\theta_1$

$\theta_0$

Andrew Ng
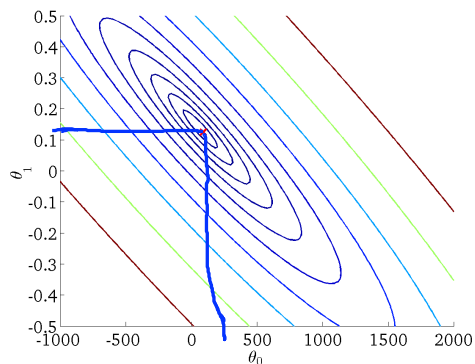
$h_\theta(x)$

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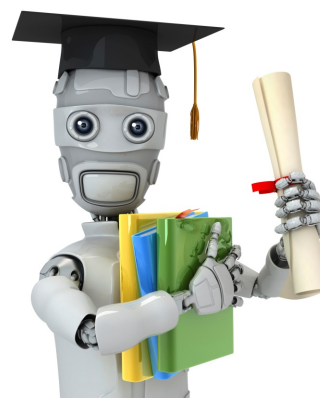
$J(\theta_0, \theta_1)$

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



Linear regression with one variable

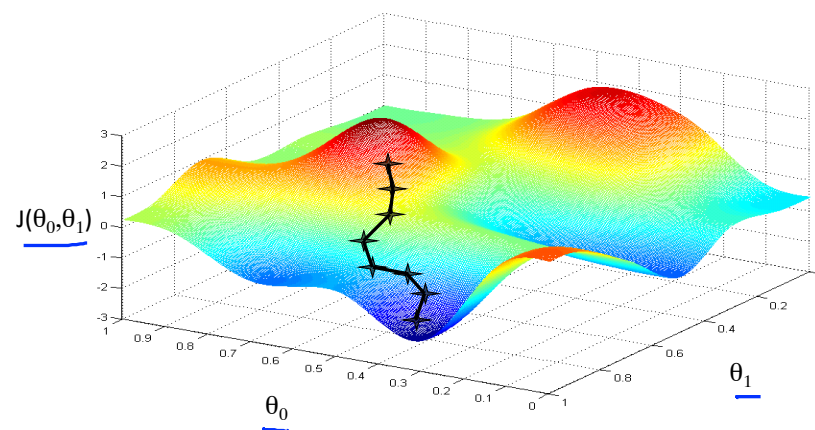Gradient descent

Machine Learning

Have some function $J(\theta_0, \theta_1)$   $J(\Theta_0, \Theta_1, \Theta_2, \ldots, \Theta_n)$

Want $\min_{\theta_0, \theta_1} J(\theta_0, \theta_1)$   $\min_{\Theta_0 \ldots \Theta_n} J(\Theta_0, \ldots, \Theta_n)$
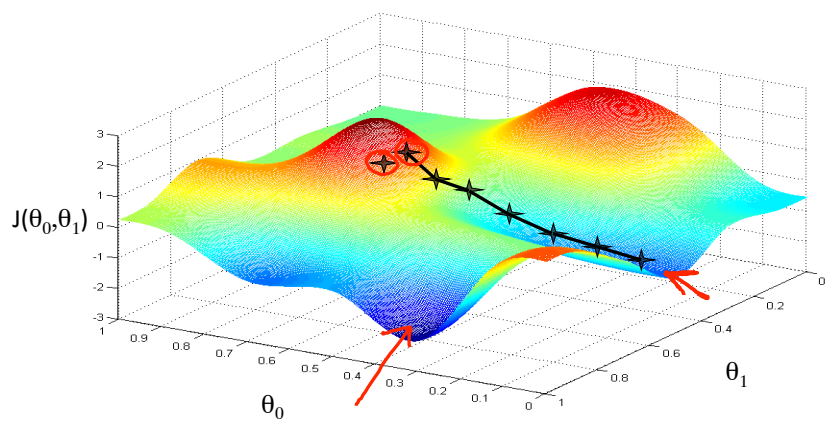
**Outline:**

- Start with some $\theta_0, \theta_1$   $(Say \quad \Theta_0 = 0, \Theta_1 = 0)$

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

   until we hopefully end up at a minimum



$J(\theta_0, \theta_1)$

$\theta_1$

$\theta_0$

**Gradient descent algorithm**

Assignment
$a := b$

Truth assertion
$a = b$

$a := a+1$

$a = a+1$ ✗

$\theta_0, \theta_1$

repeat until convergence {

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

}

learning rate

Simultaneously update $\theta_0$ and $\theta_1$

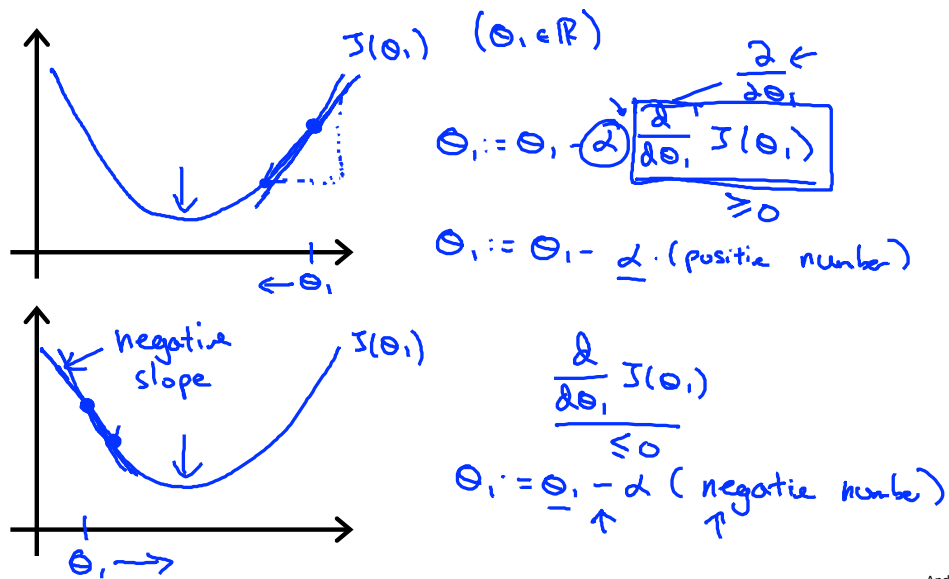| Correct: Simultaneous update | Incorrect: |
|---|---|
| $\text{temp0} := \theta_0 - \alpha \frac{\partial}{\partial \theta_0} J(\theta_0, \theta_1)$ | $\text{temp0} := \theta_0 - \alpha \frac{\partial}{\partial \theta_0} J(\theta_0, \theta_1)$ |
| $\text{temp1} := \theta_1 - \alpha \frac{\partial}{\partial \theta_1} J(\theta_0, \theta_1)$ | $\theta_0 := \text{temp0}$ |
| $\theta_0 := \text{temp0}$ | $\text{temp1} := \theta_1 - \alpha \frac{\partial}{\partial \theta_1} J(\theta_0, \theta_1)$ |
| $\theta_1 := \text{temp1}$ | $\theta_1 := \text{temp1}$ |

Andrew Ng

Linear regression with one variable

Gradient descent intuition

Machine Learning

**Gradient descent algorithm**

repeat until convergence {

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

}

learning rate

Derivative

$\min_{\theta_1} J(\theta_1)$

$\theta_1 \in \mathbb{R}.$
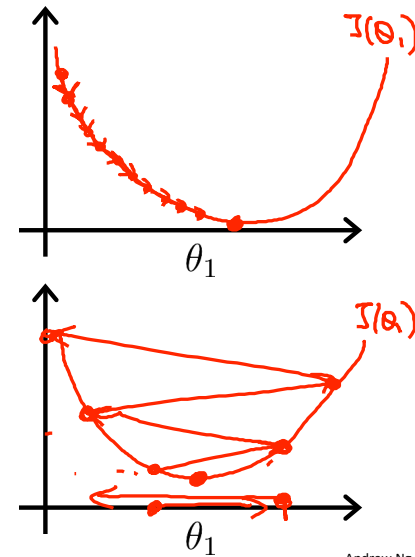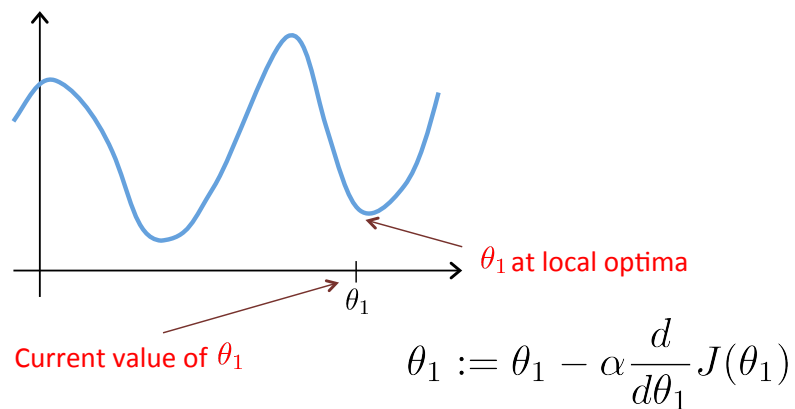
Andrew Ng

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

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

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

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

negative slope  $J(\theta_1)$

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

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

$\theta_1 \longrightarrow$

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

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

$J(\theta_1)$

$\theta_1$

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

$J(\theta_1)$

$\theta_1$

$\theta_1$ at local optima

Current value of $\theta_1$

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

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.

$J(\theta_1)$

$J(\theta_1)$

$\theta_1$

# Linear regression with one variable

## Gradient descent for linear regression

### Machine Learning

**Gradient descent algorithm**

repeat until convergence {

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

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

}

**Linear Regression Model**

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

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

---

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

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

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

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

**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

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

$J(\theta_0,\theta_1)$

$\theta_0$

$\theta_1$

Andrew Ng



$J(\theta_0,\theta_1)$

$\theta_0$

$\theta_1$

Andrew Ng



"Convex function"

Bowl-shaped

$J(\theta_0,\theta_1)$

$\theta_1$

$\theta_0$

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

Price $ (in 1000s)

Size (feet$^2$)

$h(x) = -900 - 0.1x$

Training data
Current hypothesis

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

$\theta_1$

$\theta_0$

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 $\theta_0, \theta_1$)

$h_\theta(x)$

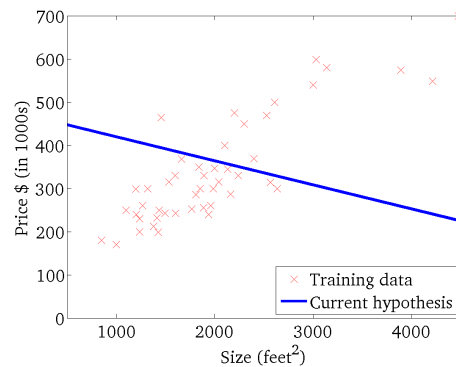(for fixed $\theta_0, \theta_1$, this is a function of x)

$J(\theta_0, \theta_1)$

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

Andrew Ng

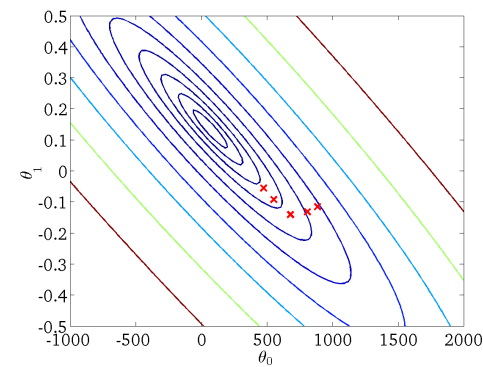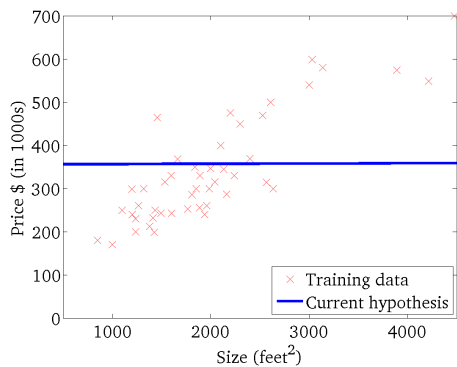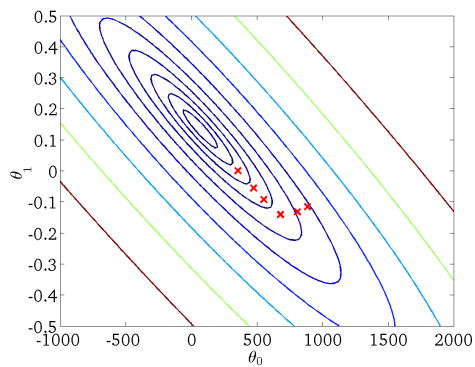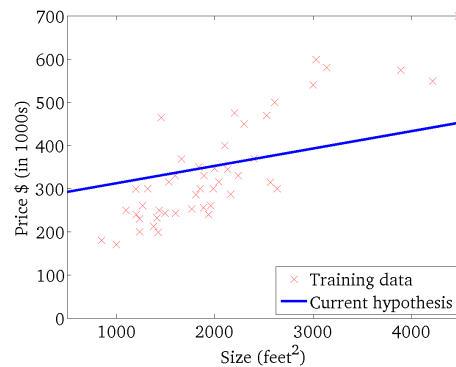$h_\theta(x)$

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

$J(\theta_0, \theta_1)$

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

$h_\theta(x)$

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

$J(\theta_0, \theta_1)$

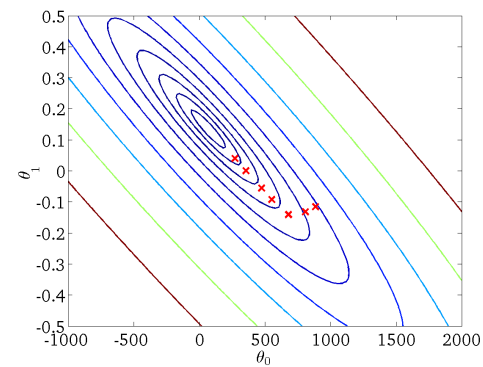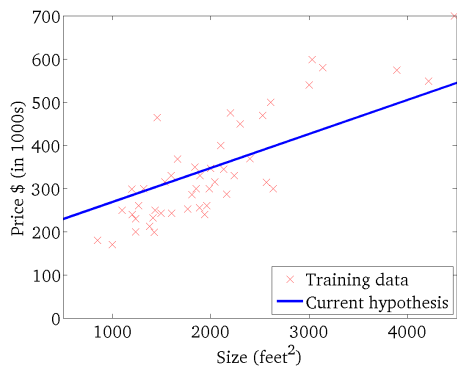(function of the parameters $\theta_0, \theta_1$)

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)

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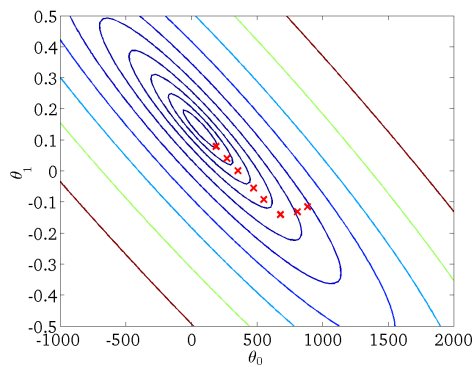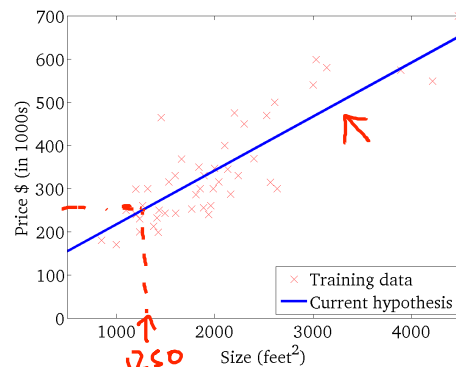(function of the parameters $\theta_0, \theta_1$)

Andrew Ng

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

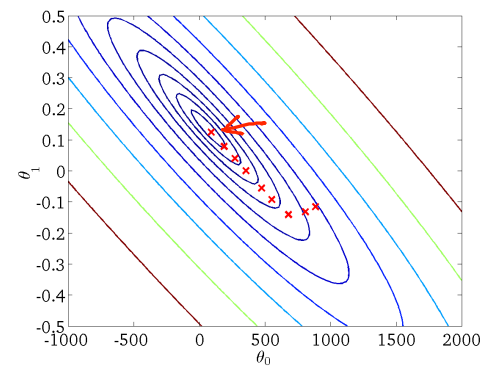$J(\theta_0, \theta_1)$
(function of the parameters $\theta_0, \theta_1$)

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

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

1250

Andrew Ng

**"Batch" Gradient Descent**

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

$$\rightarrow \sum_{i=1}^{m} \left( h_\theta(x^{(i)}) - y^{(i)} \right)$$