Linear regression with one variable

Model representation

Machine Learning

# Housing Prices (Portland, OR)

Price (in 1000s of dollars)

220K

1250

Size (feet²)

## Supervised Learning

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

## Regression Problem

Predict real-valued output

Classification: Discrete-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 |
| ... | ... |

$m = 47$

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

$h$ maps

in the early days people use this name "hypothesis" and now they keep using it ~

(estimated value of y)

to y's.

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

$$h_\theta(x) = \frac{\theta_0 + \theta_1 x}{\text{Shorthand}: h(x)}$$

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

Linear regression with one variable. $(x)$
Univariate linear regression.
└ one variable

Andrew Ng

Linear regression with one variable

Cost function

Machine Learning

| Size in feet$^2$ (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 ?

Andrew Ng

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



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

$h(x) = 0.5x$

$h_\theta(x)$

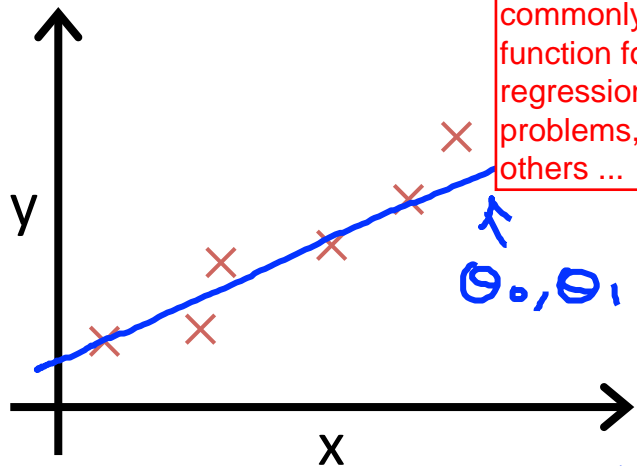$h(x)$

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

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

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

using square error is the most commonly used function for regression problems, there are others ...

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

\#training examples

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

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

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

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

$x, y$

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

square error cost function

Cost function

Squared error function

Andrew Ng

Linear regression
with one variable

Cost function
intuition I

Machine Learning

Hypothesis:

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

$$h_\theta(x) = \underline{\theta_1 x}$$

$$\theta_0 = 0$$

Parameters:

$$\theta_0, \theta_1$$

$$\theta_1$$

h(x)

Cost Function:

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

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

$$\theta, x^{(i)}$$

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

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

The cost function is (1/2m)*(sum of squares of errors). Since we have m points, there will be m error values, one corresponding to each point. The 1/m averages the errors and the 1/2 is included so that when we differentiate the function, it cancels out

Andrew Ng

$\rightarrow h_\theta(x)$

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



$h_\theta(x)$

$\theta_1 = 1$

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

$\theta_1 = 1$

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

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

$\rightarrow J(\theta_1)$

(function of the parameter $\theta_1$)



$J(\theta_1)$

$\theta_1$

$\theta_1 = 0.5$?

$J(1) = 0$

Andrew Ng

# $h_\theta(x)$

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



$h_\theta(x)$

$\{ \, \leftarrow y^{(i)}$

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

$\theta_1 = 0.5$

# $J(\theta_1)$

(function of the parameter $\theta_1$)



$J(\theta_1)$

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

$$= \frac{1}{2\times 3}(3.5) = \frac{3.5}{6} \approx 0.58$$

$\theta_1 = 0?$

$J(0) = ?$

Andrew Ng

$h_\theta(x)$

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

$J(\theta_1)$

(function of the parameter $\theta_1$)

$$J(0) = \frac{1}{2m}(1^2 + 2^2 + 3^2)$$
$$= \frac{1}{6} \cdot 14 \approx 2.3$$

$h(x) = -0.5x$

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

Andrew Ng

Linear regression
with one variable

Cost function
intuition II

Machine Learning

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: $\displaystyle\minimize_{\theta_0, \theta_1} J(\theta_0, \theta_1)$

$$h_\theta(x)$$

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

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

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



Price ($) in 1000's

$\theta_0 = 50$

$\theta_1 = 0.06$

Size in feet$^2$ (x)

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

$\theta_1$

$\theta_0, \theta_1$

Contour plots
Contour figures.

$J(\theta_0, \theta_1)$

$J(\theta_0, \theta_1)$

$\theta_1$

$\theta_0$

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 $\theta_0, \theta_1$)



same value of J on the curve

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 $\theta_0, \theta_1$)



$h(x) = 360 + 0 \cdot x$

$\theta_0 = 360$

$\theta_1 = 0$

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

Linear regression
with one variable

Gradient
descent

Machine Learning

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

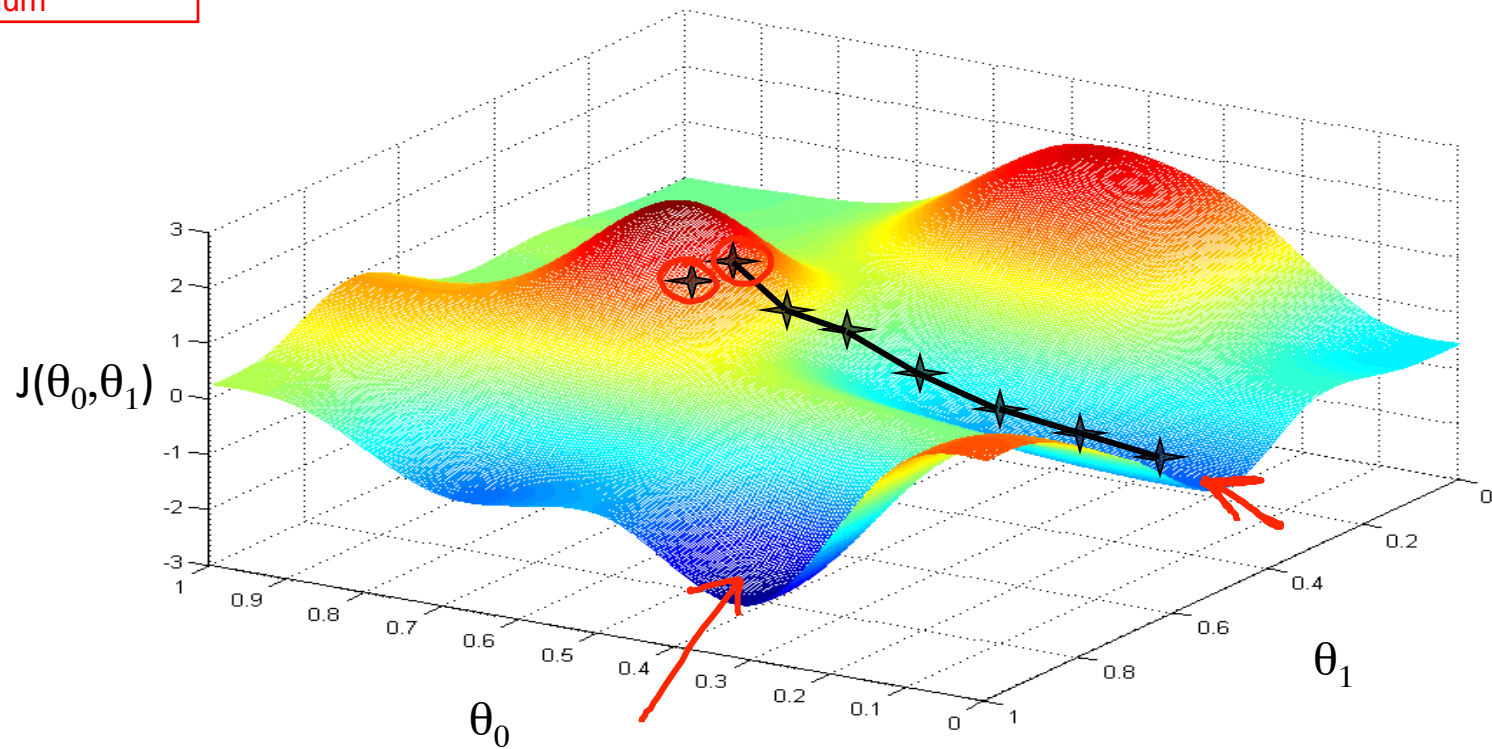Want $\min\limits_{\theta_0, \theta_1} J(\theta_0, \theta_1)$     $\min\limits_{\theta_0 \ldots \theta_n} J(\theta_0, \ldots, \theta_n)$

**Outline:**

- Start with some $\theta_0, \theta_1$     $(\text{say } \theta_0 = 0, \theta_1 = 0)$

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

  until we hopefully end up at a minimum

$J(\theta_0, \theta_1)$

$\theta_0$

$\theta_1$

Andrew Ng

different starting point directs you to different local minimum

$J(\theta_0,\theta_1)$

$\theta_0$

$\theta_1$

Andrew Ng

# Gradient descent algorithm

$\theta_0, \theta_1$

Assignment

$a := b$

$a := a+1$

Truth assertion

$a = b$

$a = a+1$  ✗

repeat until convergence {

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

}

learning rate

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

Simultaneously update $\theta_0$ and $\theta_1$

---

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 {

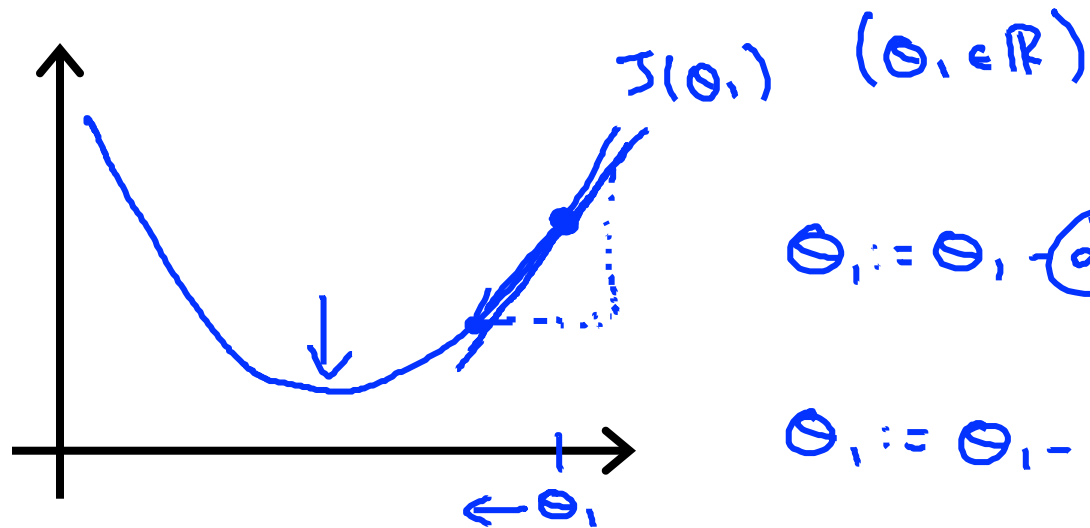$$\theta_j := \theta_j - \alpha \frac{\partial}{\partial \theta_j} J(\theta_0, \theta_1)$$

(simultaneously update $j = 0$ and $j = 1$)

}

learning rate

derivative

$$\min_{\theta_1} J(\theta_1) \qquad \theta_1 \in \mathbb{R}.$$

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

$\leftarrow \theta_1$

$\theta_1 := \theta_1 - \textcircled{$\alpha$} \boxed{\dfrac{\partial}{\partial \theta_1} J(\theta_1)} \quad \dfrac{\partial}{\partial \theta_1} \leftarrow$

$\geq 0$

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

negative slope

$J(\theta_1)$

$\theta_1 \longrightarrow$

$\dfrac{\dfrac{\partial}{\partial \theta_1} J(\theta_1)}{\leq 0}$

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

Andrew Ng

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

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

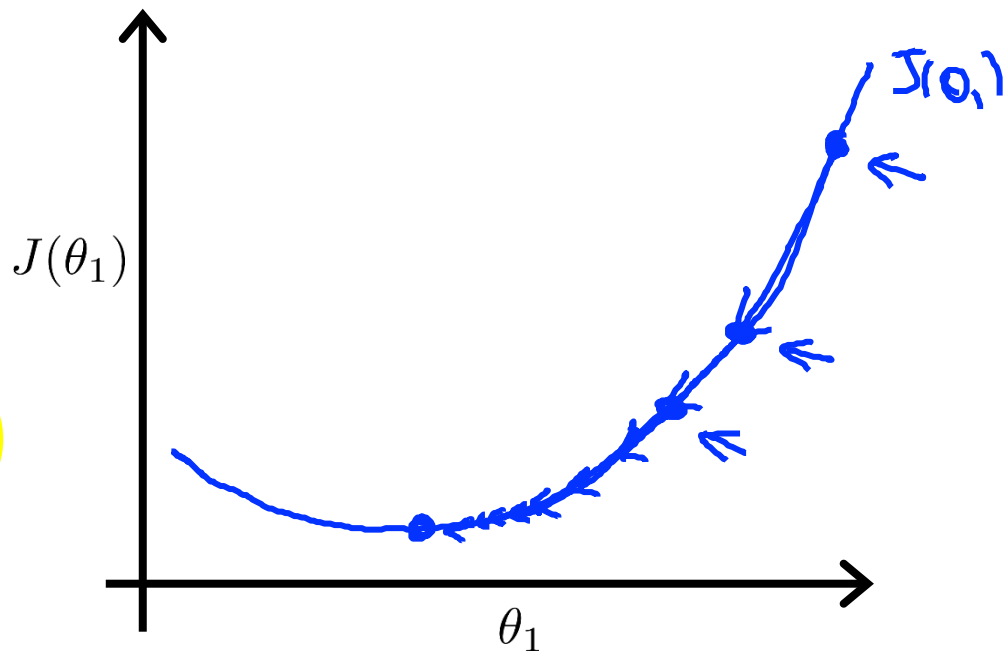If α is too large, gradient descent can overshoot the minimum. It may fail to converge, or even diverge.



$J(\theta_1)$



$J(\theta_1)$

$\theta_1$ at local optima

Current value of $\theta_1$

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

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.

$J(\theta_1)$

$J(\theta_1)$

$\theta_1$

Linear regression
with one variable

Gradient descent for
linear regression

Machine Learning

## Gradient descent algorithm

$\text{repeat until convergence } \{$

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

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

$\}$

## Linear Regression Model

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

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

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

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

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

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

# **Gradient descent algorithm**

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

repeat until convergence {

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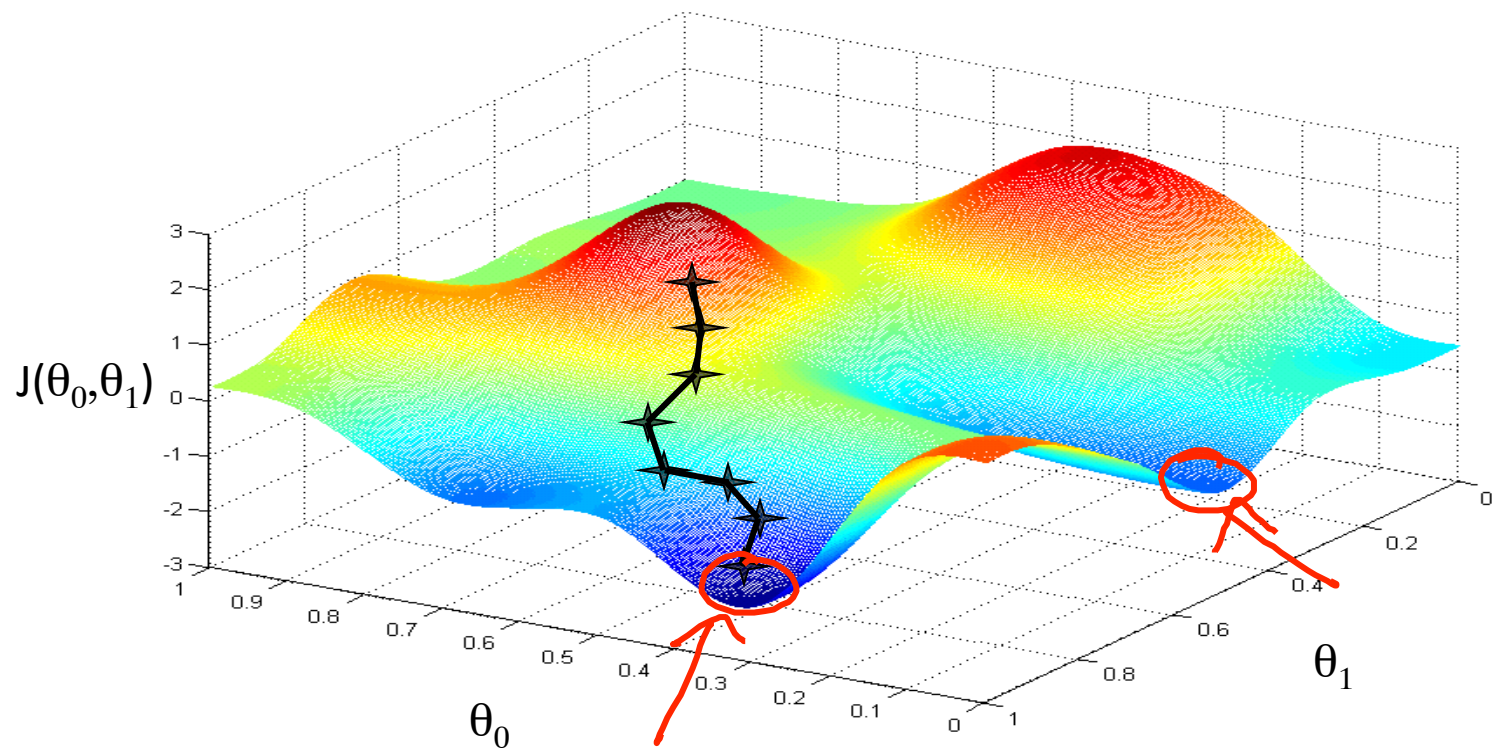
$$\theta_1 := \theta_1 - \alpha \frac{1}{m} \sum_{i=1}^{m} \left( h_\theta(x^{(i)}) - y^{(i)} \right) \cdot x^{(i)}$$

}

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

update $\theta_0$ and $\theta_1$ simultaneously

$J(\theta_0,\theta_1)$

$\theta_0$

$\theta_1$

Andrew Ng

$J(\theta_0, \theta_1)$

$\theta_0$

$\theta_1$

Andrew Ng

the cost function for linear regression
is always going to be a bow shaped function !!
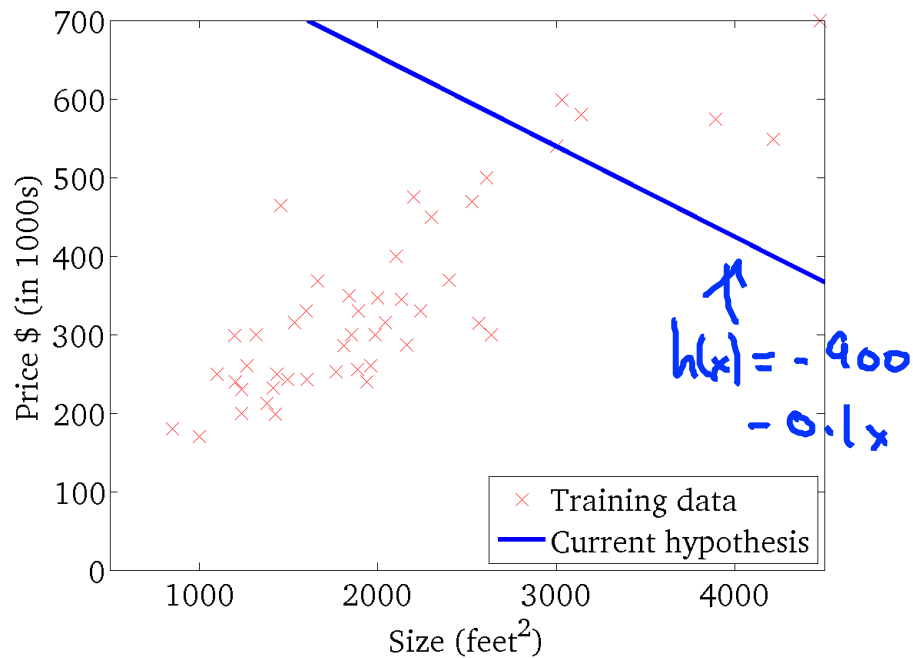there are no local optima (other than the
global optimum)

"Convex function"

Bowl-shaped

$J(\theta_0, \theta_1)$

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 $\theta_0, \theta_1$)

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

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 $\theta_0, \theta_1$)



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 $\theta_0, \theta_1$)

Price $ (in 1000s)

1250

Size (feet$^2$)

Training data
Current hypothesis

$\theta_1$

$\theta_0$

Andrew Ng

# "Batch" Gradient Descent

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

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