**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

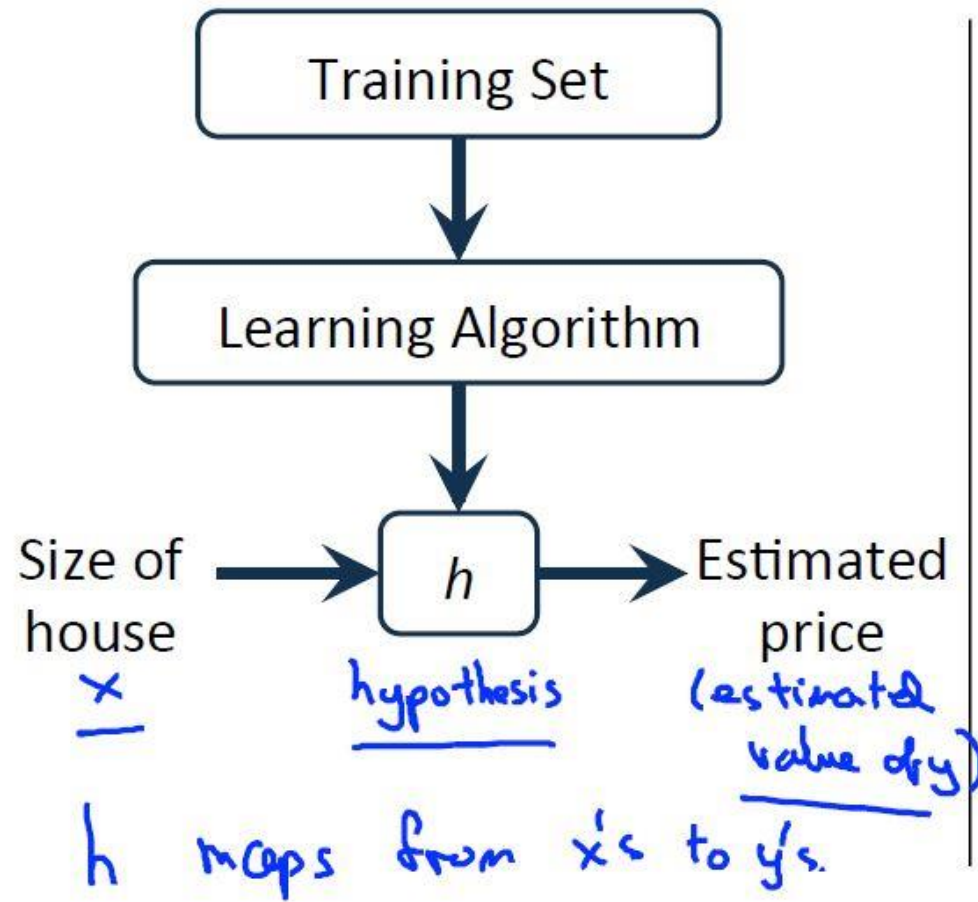$(x, y)$ - one training example

$(x^{(i)}, y^{(i)})$ - $i^{th}$ training example
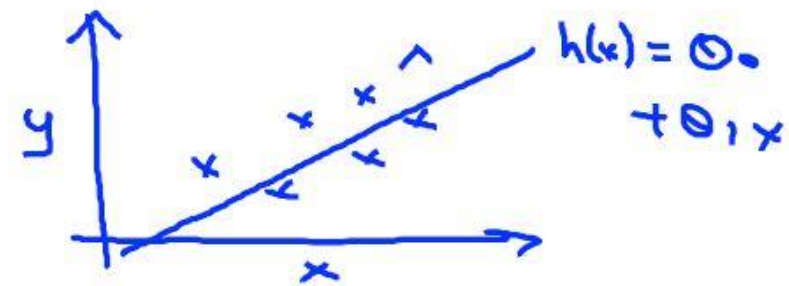
$x^{(1)} = 2104$

$x^{(2)} = 1416$

$y^{(1)} = 460$

Training Set

Learning Algorithm

Size of house $\xrightarrow{\text{x}}$ $h$ $\xrightarrow{}$ Estimated price (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. $(x)$
Univariate linear regression.

$\llcorner$ one variable

Training Set

| 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 ?

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



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

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

$h(x) = 0.5x$

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

$h_\theta(x)$

$h(x)$

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

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

\# training examples

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

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

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

$x, y$

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

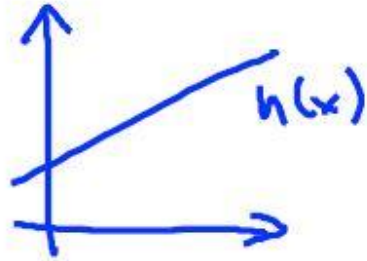$$\underset{\theta_0, \theta_1}{\text{minimize}} \quad J(\theta_0, \theta_1)$$

Cost function

Squared error function

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

## Simplified

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

$$\theta_0 = 0$$

$$\theta_1$$

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

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

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

$\rightarrow h_\theta(x)$

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

$h_\theta(x)$

$\theta_1 = 1$

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

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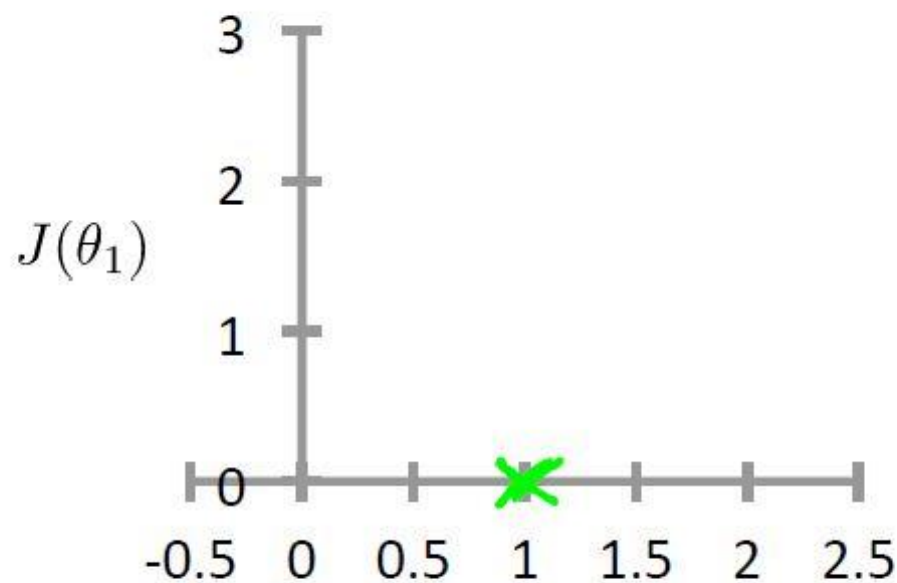
$$J(\theta_1) = \frac{1}{2m} \sum_{i=1}^{m} (h_\theta(x^{(i)}) - y^{(i)})^2$$

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

$\rightarrow J(\theta_1)$

(function of the parameter $\theta_1$)

$J(\theta_1)$

$\theta_1$

$\theta_1 = 0.5?$

$J(1) = 0$

# $h_\theta(x)$

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



# $J(\theta_1)$

(function of the parameter $\theta_1$)



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

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

$$\theta_1 = 0?$$

$$J(0) = ?$$

# $h_\theta(x)$

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



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

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

# $J(\theta_1)$

(function of the parameter $\theta_1$)

$5.25$

$\theta_1 = 1$

$h(x) = -0.5x$

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

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)

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

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



$\theta_0 = 50$

$\theta_1 = 0.06$

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

$\theta_1$

$\theta_0, \theta_1$

Contour plots
Contour figures.

$J(\theta_0, \theta_1)$

$J(\theta_0, \theta_1)$

$\theta_1$

$\theta_0$

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

Want $\min_{\theta_0, \theta_1} J(\theta_0, \theta_1)$    $\min_{\theta_0 \cdots \theta_n} J(\theta_0, \ldots, \theta_n)$

**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

# Gradient descent algorithm

$a := b$

$a := a+1$

$a = b$

$a = a+1$

$\theta_0, \theta_1$

repeat until convergence {

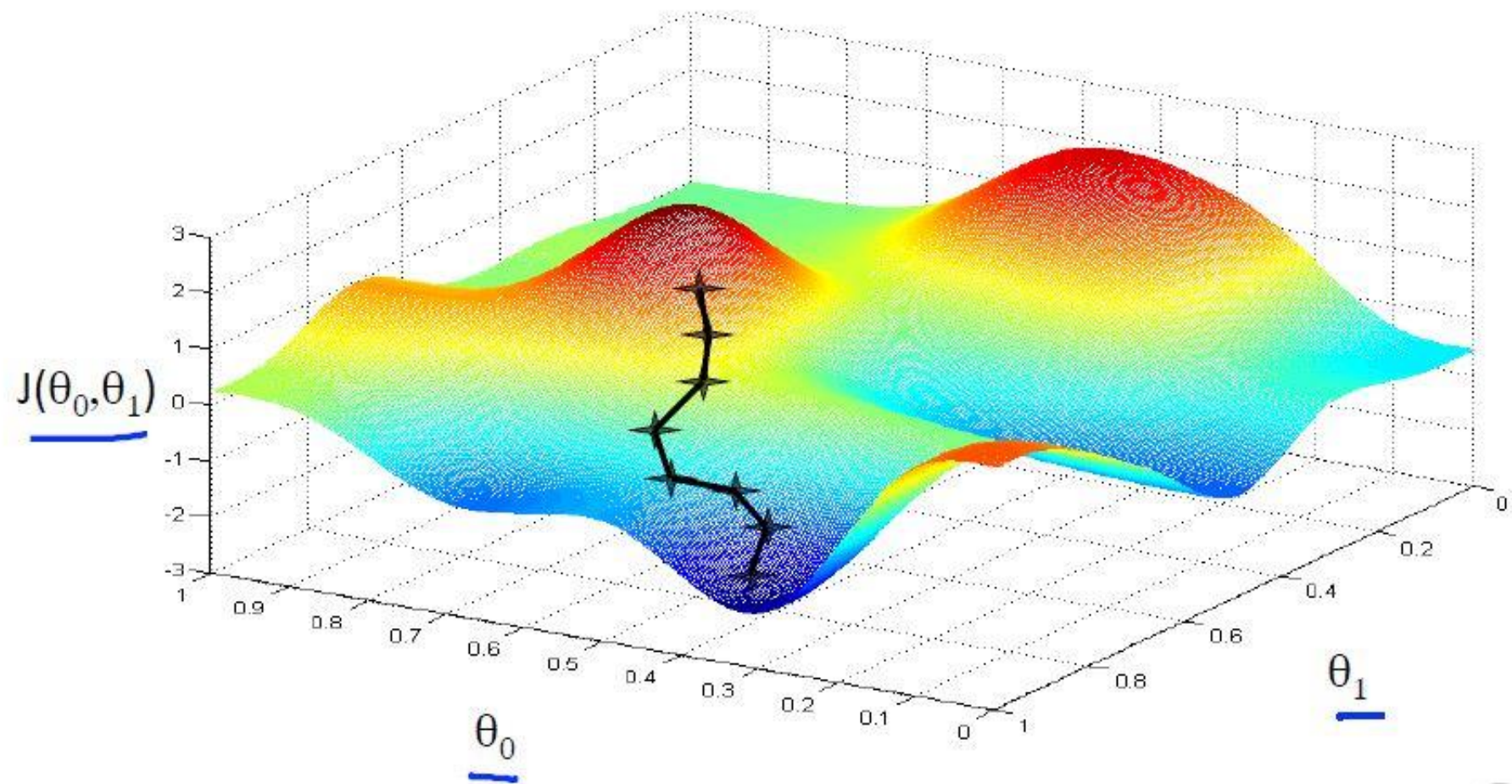$$\theta_j := \theta_j - \alpha \frac{\partial}{\partial \theta_j} J(\theta_0, \theta_1) \qquad \text{(for } j = 0 \text{ and } j = 1)$$

}

learning rate

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:

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

$\theta_0 := \text{temp0}$

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

$\theta_1 := \text{temp1}$

# 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

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

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

$$\theta_1 := \theta_1 - \boxed{\alpha} \; \boxed{\dfrac{d}{d\theta_1} J(\theta_1)} \underset{\geq 0}{}$$

$\dfrac{d}{d\theta_1}$

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

negative slope

$J(\theta_1)$

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

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

$\theta_1 \rightarrow$

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

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.

Current value of $\theta_1$

$\theta_1$ at local optima

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

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.

$J(\theta_1)$

$J(\theta_1)$

$\theta_1$

## 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

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

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

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

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

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

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

# Gradient descent algorithm

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

repeat until convergence {

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

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

}

update
$\theta_0$ and $\theta_1$
simultaneously

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

"Convex function"

Bowl-shaped

# "Batch" Gradient Descent

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

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

# Multiple features (variables).

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

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

# Multiple features (variables).

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

$m = 47$

Notation:

$n$ = number of features $\qquad n = 4$

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

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

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

$x_3^{(2)} = 2$

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

age

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

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

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

$$\underbrace{\begin{bmatrix} \theta_0 & \theta_1 & \cdots & \theta_n \end{bmatrix}}_{\theta^T} \begin{bmatrix} x_0 \\ x_1 \\ \vdots \\ x_n \end{bmatrix}$$

$(n+1) \times 1$ matrix

$\theta^T x$

$x$

$$h_\theta(x) = \theta_0 x_0 + \theta_1 x_1 + \cdots + \theta_n x_n \Longleftarrow$$

$\searrow^{=1}$

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

Multivariate linear regression. $\longleftarrow$

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: $\theta_0, \theta_1, \ldots, \theta_n$    $\theta$    $n+1$ - 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 {

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

}

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

# Gradient Descent

Previously (n=1):

Repeat {

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

$$\underbrace{\phantom{\frac{1}{m} \sum_{i=1}^{m} (h_\theta(x^{(i)}) - y^{(i)})}}_{\frac{\partial}{\partial \theta_0} J(\theta)}$$

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

(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 $j = 0, \ldots, n$)

}

---

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

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

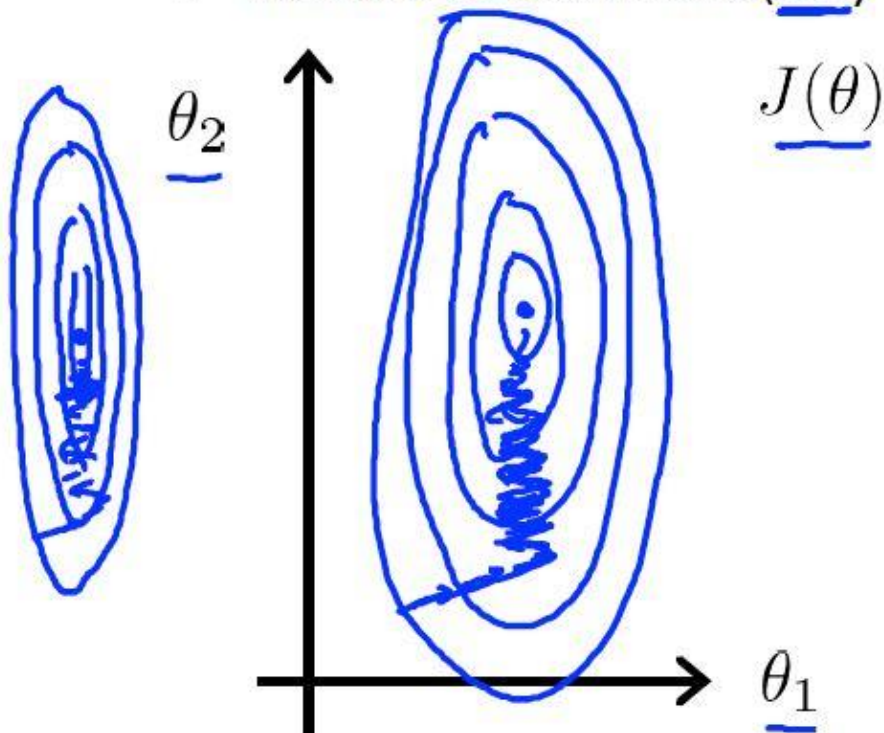$$\theta_1 := \theta_1 - \alpha \frac{1}{m} \sum_{i=1}^{m} (h_\theta(x^{(i)}) - y^{(i)}) x_1^{(i)}$$

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

...

# Feature Scaling

Idea: Make sure features are on a similar scale.
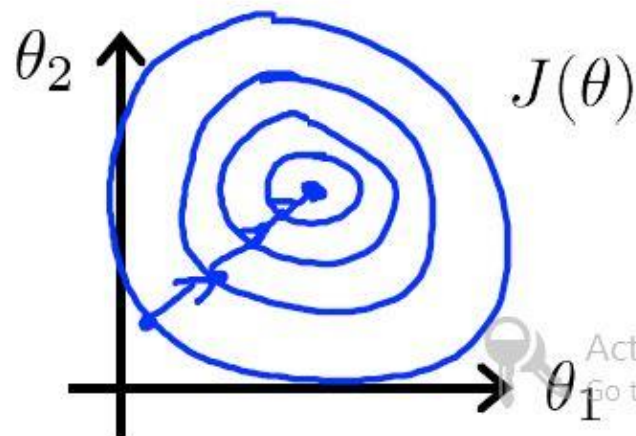
E.g. $x_1$ = size (0-2000 feet$^2$) ←

$x_2$ = number of bedrooms (1-5) ←

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

$$x_2 = \frac{\text{number of bedrooms}}{5}$$

$$0 \le x_1 \le 1 \qquad 0 \le x_2 \le 1$$



$J(\theta)$

$\theta_2$

$\theta_1$

## Feature Scaling

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

$$x_0 = 1$$

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

$x_2 = \dfrac{\#bedrooms - 2}{5 \quad 4}$

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

Averya $size = 100$

$1 - 5$ bedrooms

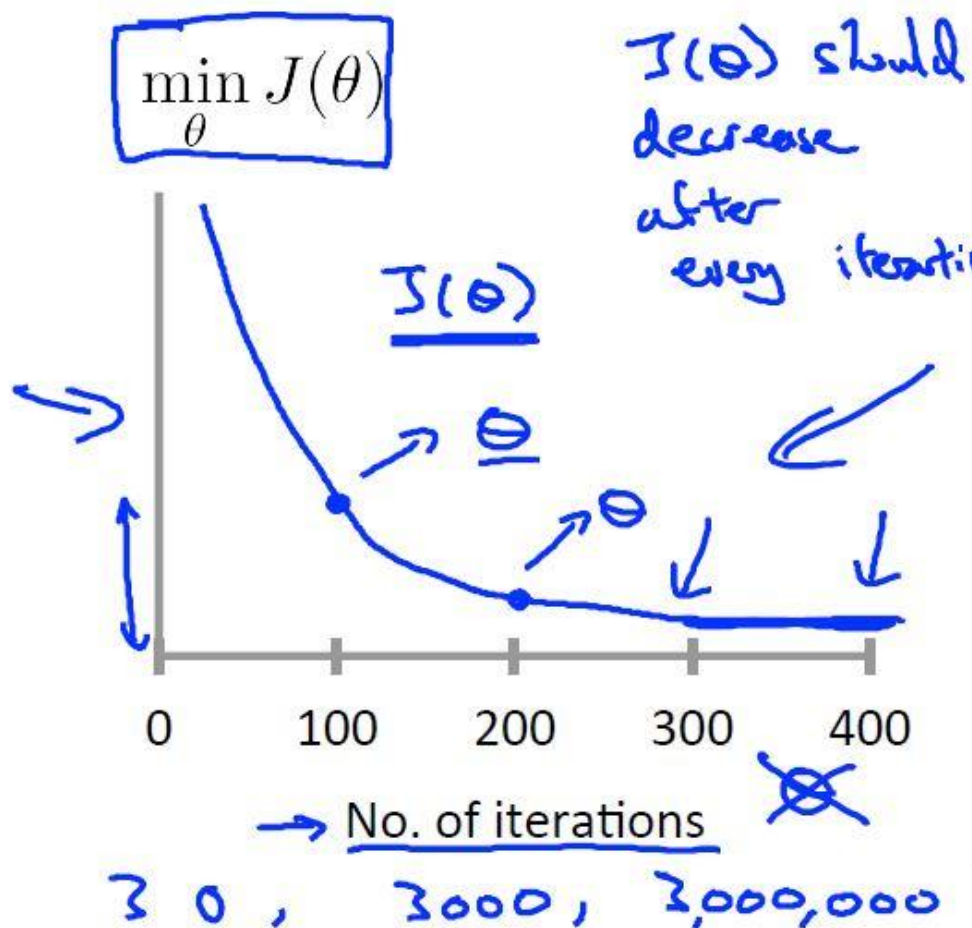$x_1 \leftarrow \dfrac{x_1 - \mu_1}{S_1}$    avg value of $x_1$ in traning set

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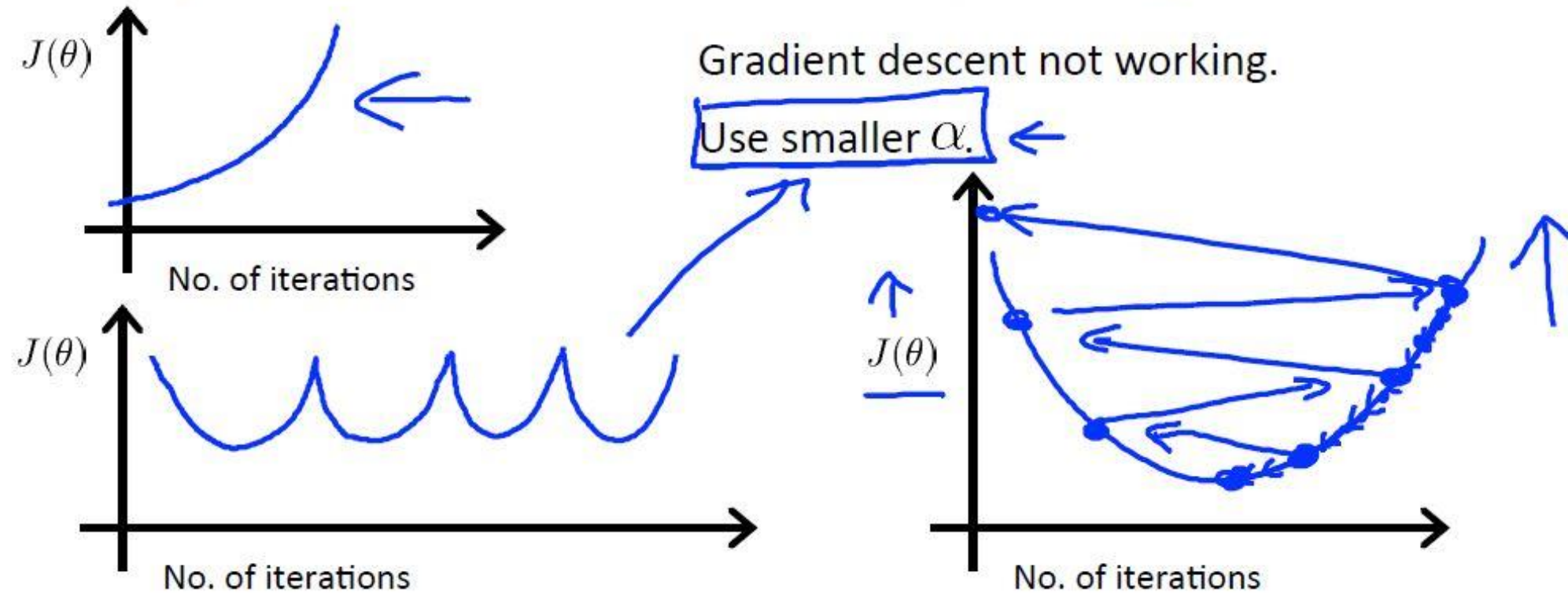(or standard deviation)

$x_2 \leftarrow \dfrac{x_2 - \mu_1}{S_2}$

Activate Win
Go to PC settings

# Making sure gradient descent is working correctly.

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

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

$J(\theta)$

→ Example automatic convergence test:

→ Declare convergence if $J(\theta)$ decreases by less than $10^{-3}$ in one iteration.

$\varepsilon$

0    100    200    300    400

→ No. of iterations

3 0 ,    3000,    3,000,000

# Making sure gradient descent is working correctly.



- For sufficiently small $\alpha$, $J(\theta)$ should decrease on every iteration.
- But if $\alpha$ is too small, gradient descent can be slow to converge.

## Summary:

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

To choose $\alpha$, try

$$\ldots, 0.001, 0.003, 0.01, 0.03, 0.1, 0.3, 1, \ldots$$