# UCK358E – INTR. TO ARTIFICIAL INTELLIGENCE
## SPRING '23

### LECTURE 3
#### LINEAR AND POLYNOMIAL REGRESSION

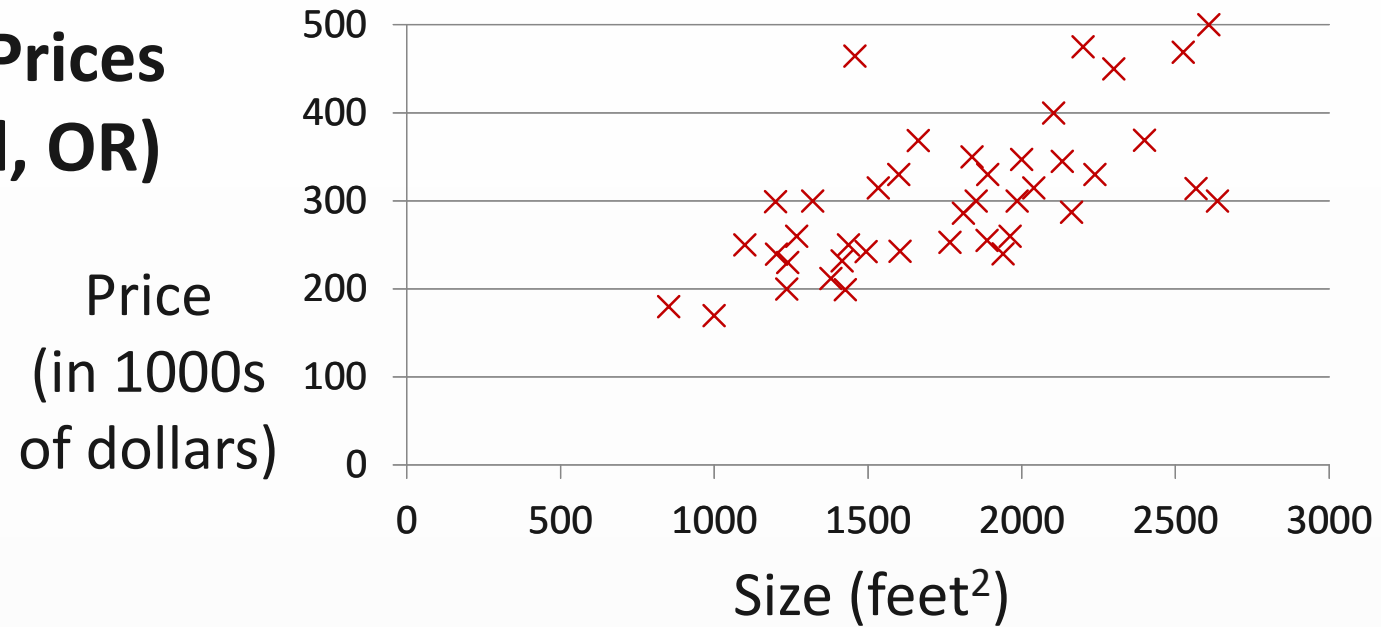Instructor: Asst. Prof. Barış Başpınar

**Housing Prices
(Portland, OR)**

Price
(in 1000s
of dollars)



Size (feet$^2$)

Supervised Learning

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

Regression Problem

Predict real-valued output

| | Size in feet² (x) | Price ($) in 1000's (y) |
|---|---|---|
| **Training set of housing prices (Portland, OR)** | 2104 | 460 |
| | 1416 | 232 |
| | 1534 | 315 |
| | 852 | 178 |
| | ... | ... |

$m = 50$

Notation:

$m$ = Number of training examples

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

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

$$\left(x^{(1)}, y^{(1)}\right) = (2104, 460)$$

$$\left(x^{(2)}, y^{(2)}\right) = (1416, 232)$$
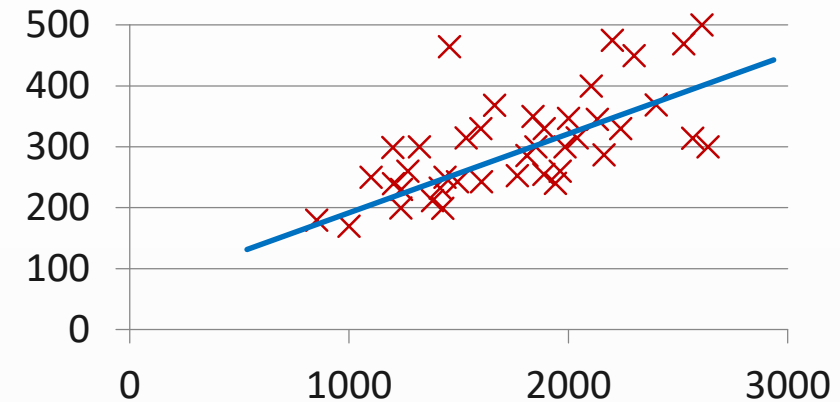
$$\left(x^{(i)}, y^{(i)}\right) \rightarrow i^{th} \text{ training example}$$

# Model Representation: linear regression



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

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

Linear regression with one variable.
Univariate linear regression.

Training Set

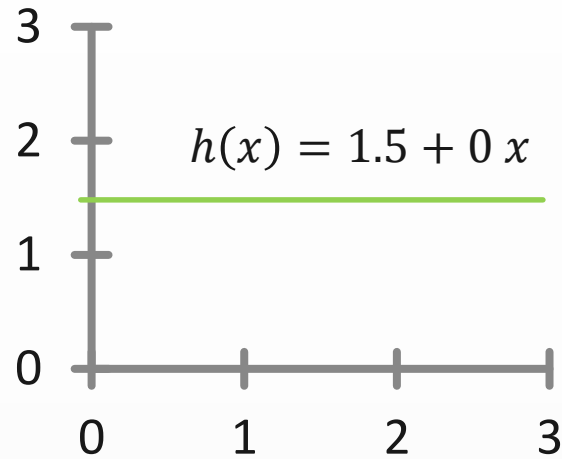| Size in feet² (x) | Price (\$) in 1000's (y) |
|---|---|
| 2104 | 460 |
| 1416 | 232 |
| 1534 | 315 |
| 852 | 178 |
| ... | ... |

**m** = 50

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

$\theta_i$'s:    Parameters
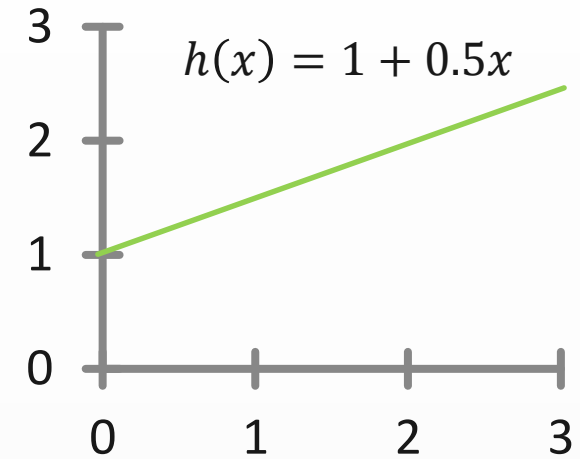
How to choose $\theta_i$'s ?

# Cost Function

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



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

$h(x) = 0.5x$

$h(x) = 1 + 0.5x$

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

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

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

Number of training examples    Squared error function

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

Cost function

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

$\theta_1, \theta_2$

y

x

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

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

# Cost Function Intuition

Hypothesis:

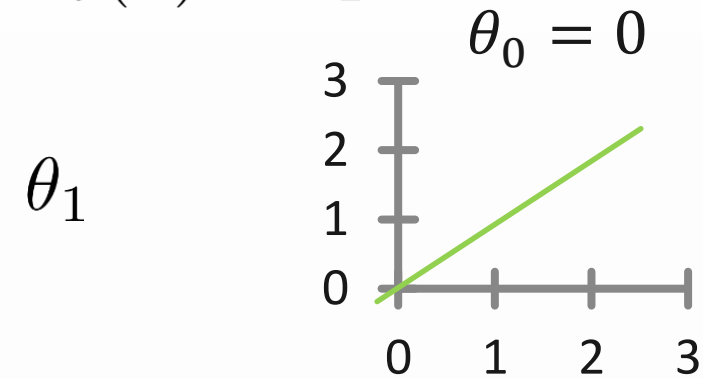$$h_\theta(x) = \theta_0 + \theta_1 x$$

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

$$\theta_0 = 0$$

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

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

# Cost Function Intuition

$$h_\theta(x)$$

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

$$J(\theta_1)$$

(function of the parameter $\theta_1$)



$\theta_1 = 1$

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

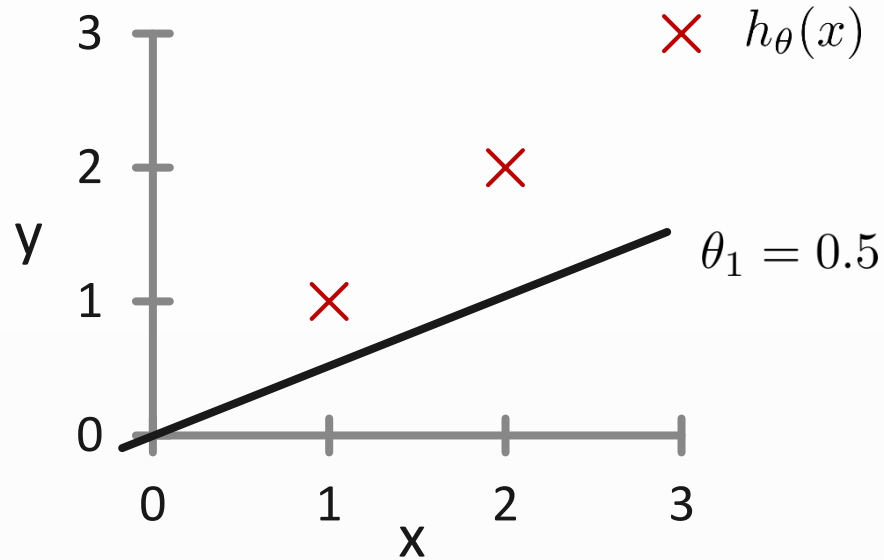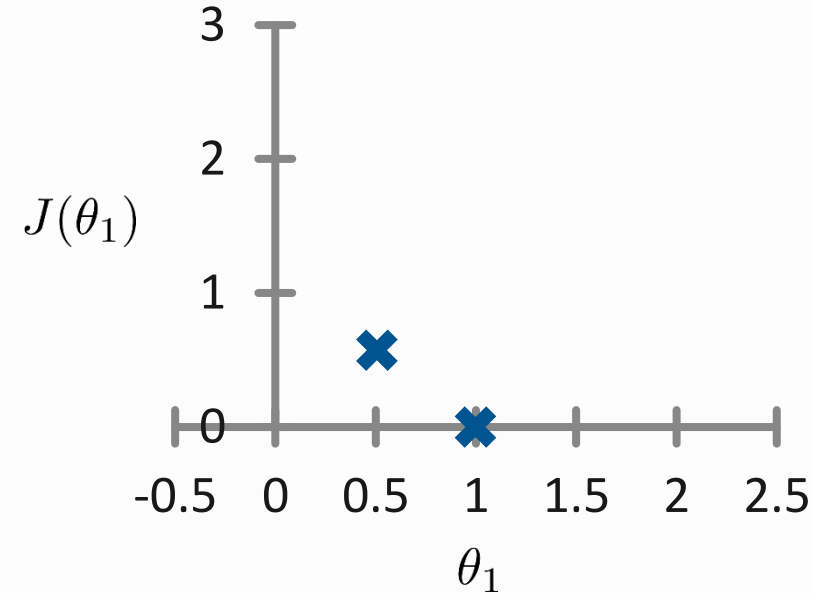$$= \frac{1}{2m}\sum_{i=1}^{m}\left(\theta_1 x^{(i)} - y^{(i)}\right)^2$$

$\longrightarrow$

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

# Cost Function Intuition

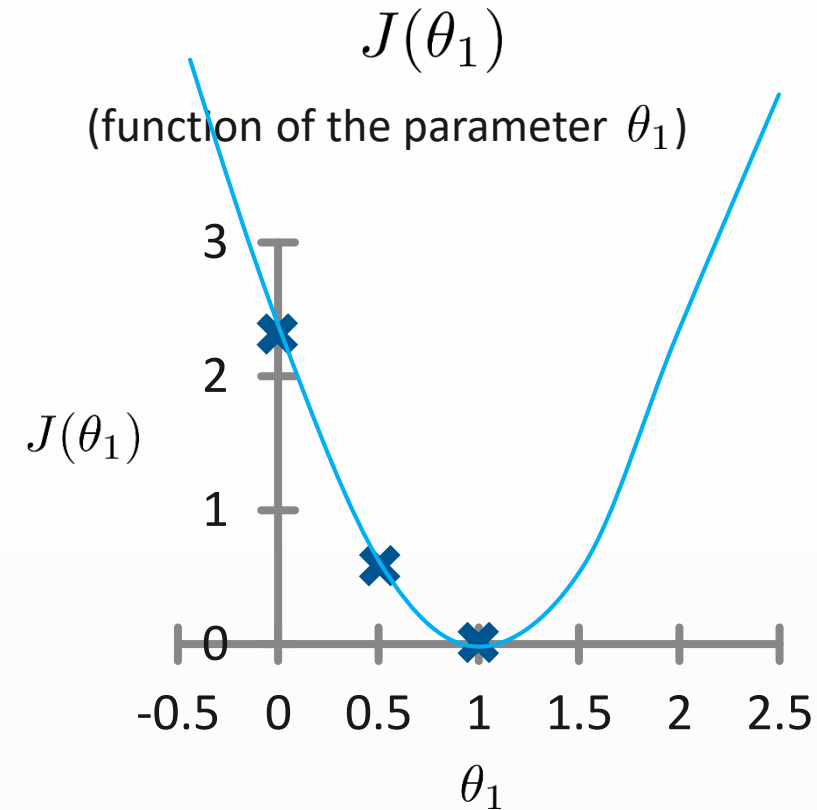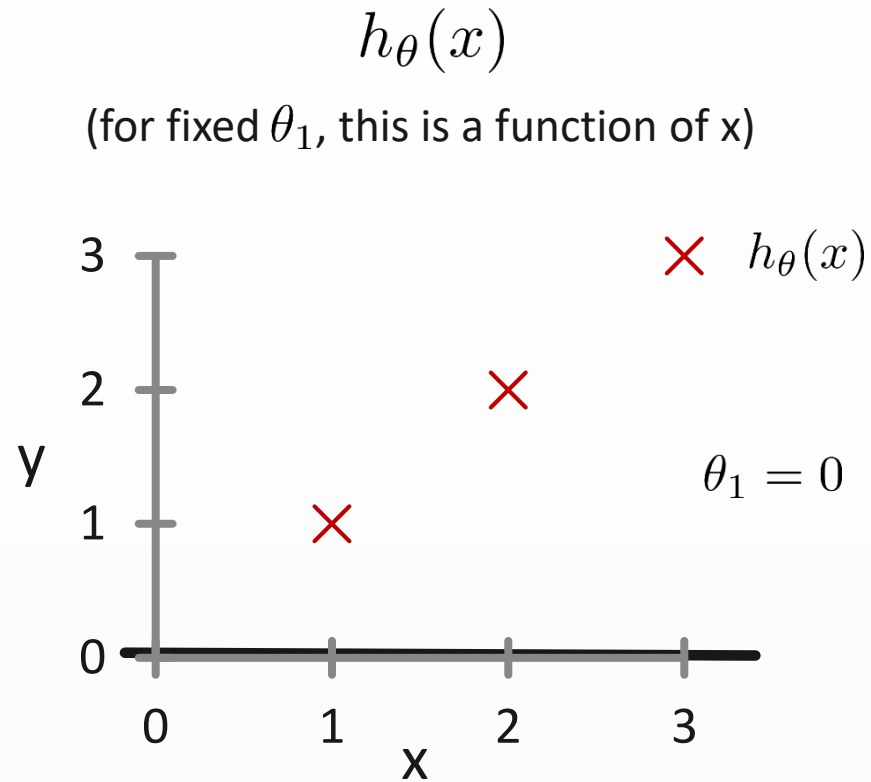$$h_\theta(x)$$

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

$$J(\theta_1)$$

(function of the parameter $\theta_1$)



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

$\theta_1 = 0.5$

$J(\theta_1)$

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

$$h_\theta(x)$$

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

$$J(\theta_1)$$

(function of the parameter $\theta_1$)



$\theta_1 = 0$

$J(\theta_1)$

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

$$J(0) = \frac{1}{2 \times 3}[(0-1)^2 + (0-2)^2 + (0-3)^2] = \frac{14}{6} = 2.33$$

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

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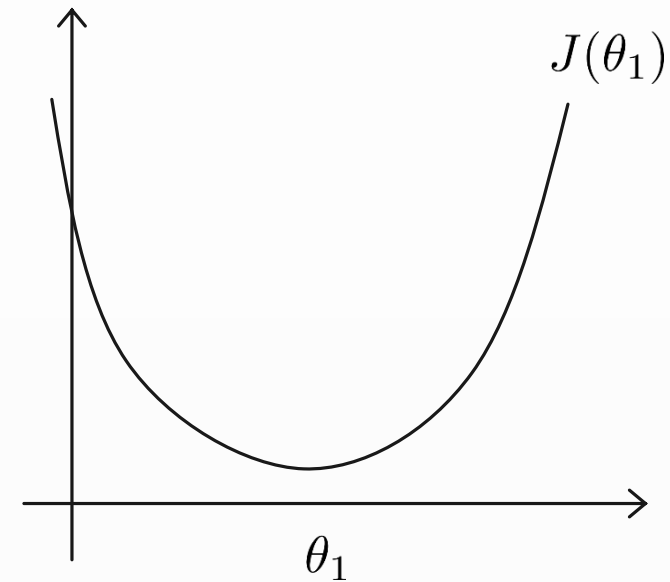
(for fixed $\theta_0, \theta_1$, this is a function of x)

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

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



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

$$\boxed{\theta_0, \theta_1}$$

# Cost Function Intuition

Contour plots



$J(\theta_0, \theta_1)$
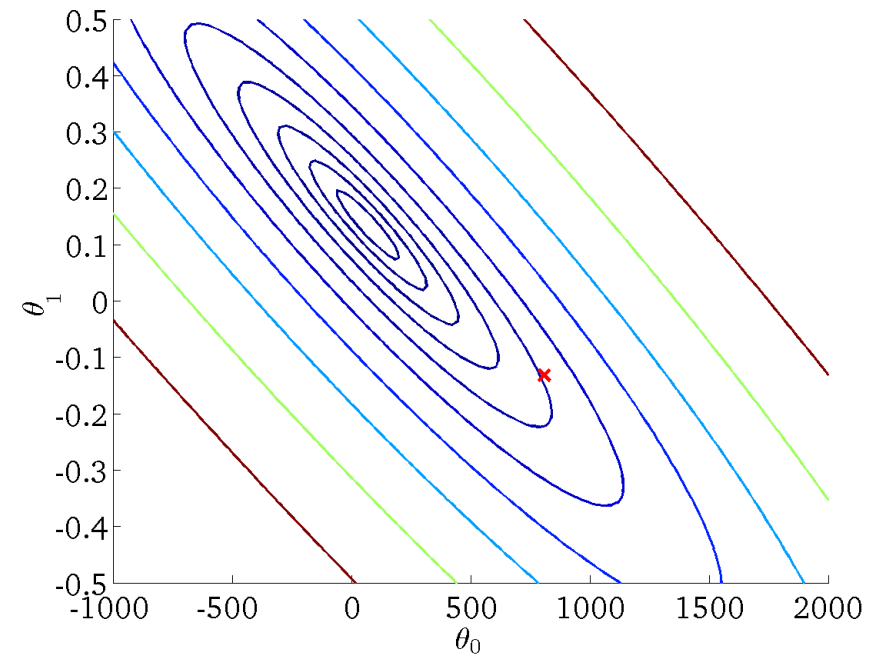
$h_\theta(x)$

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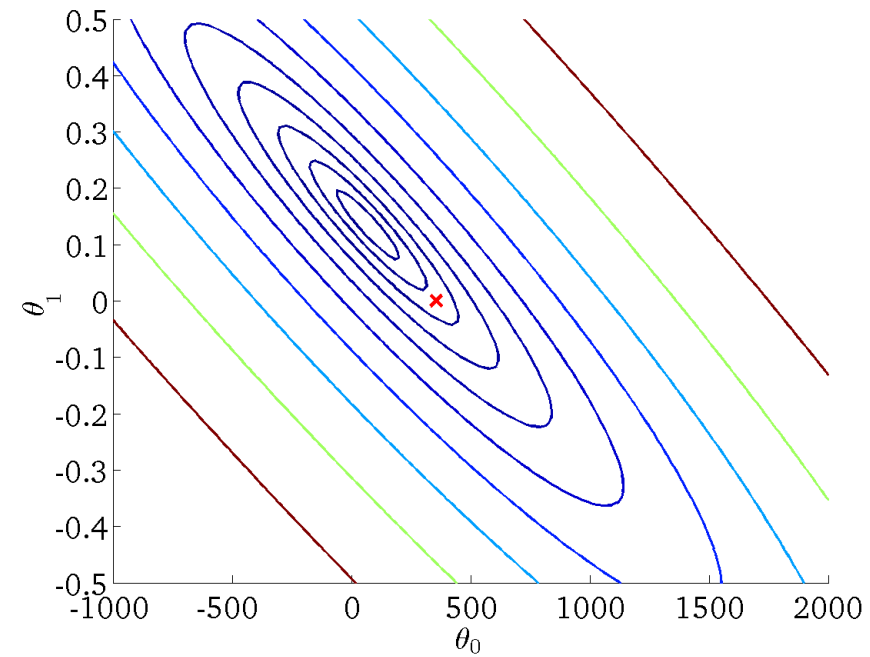
$J(\theta_0, \theta_1)$

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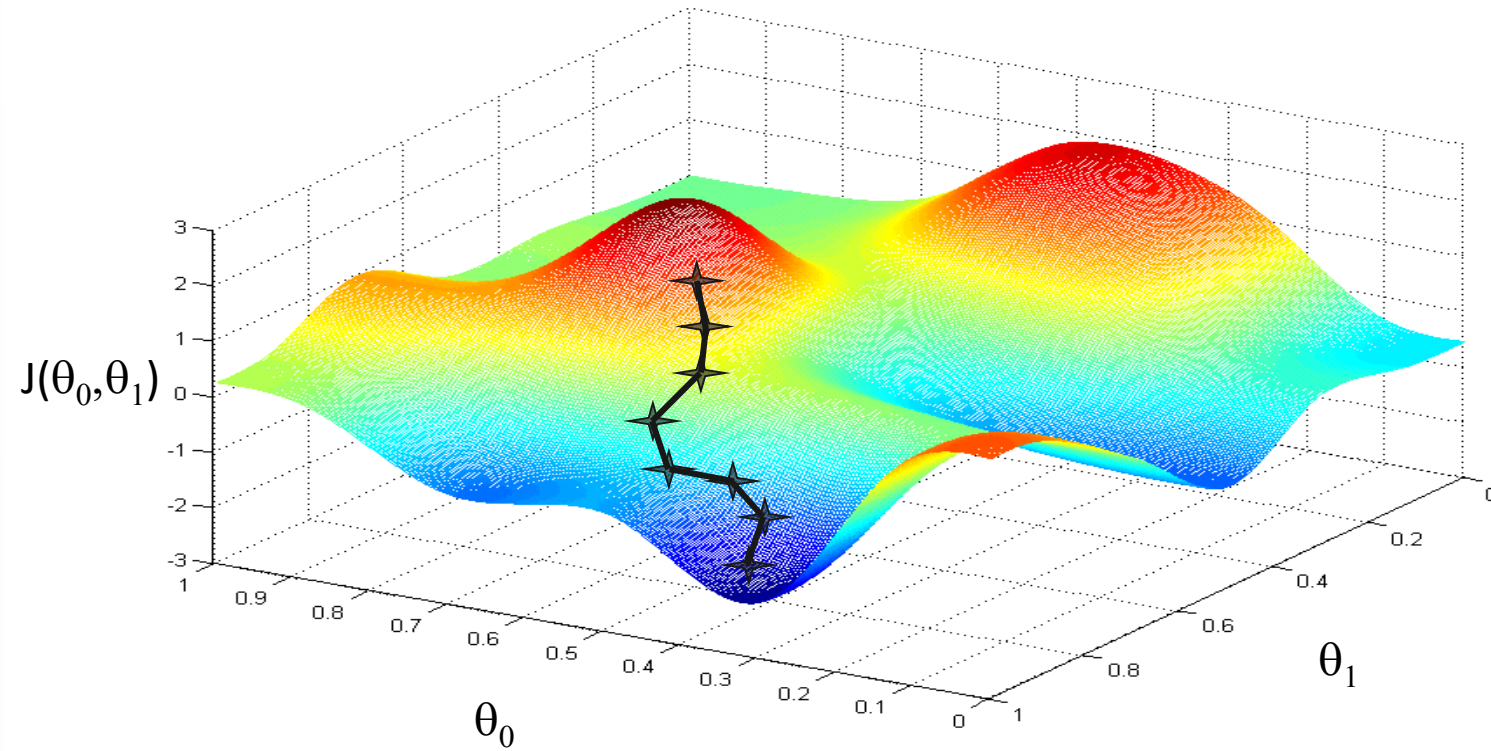
Have some function $J(\theta_0, \theta_1)$

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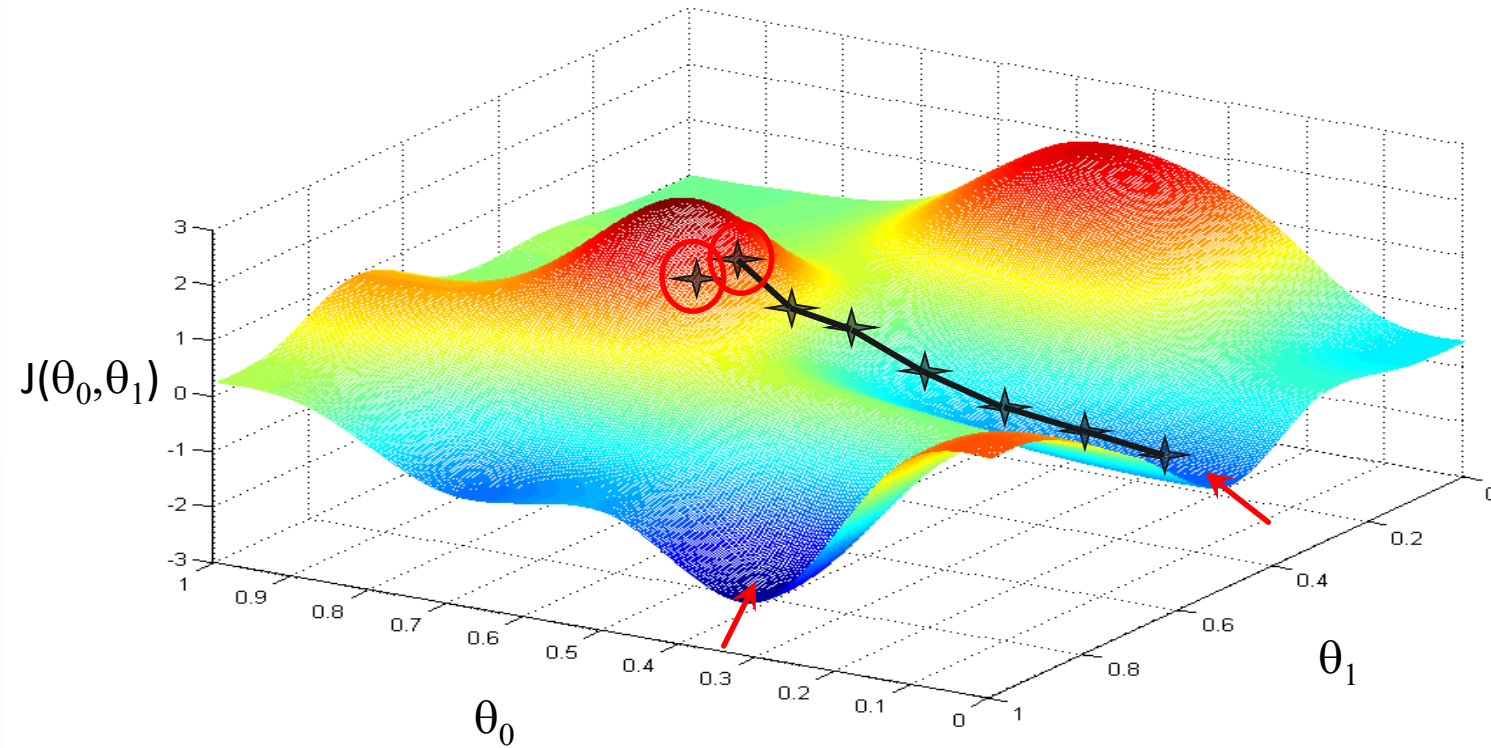
**Outline:**

- 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

repeat until convergence {

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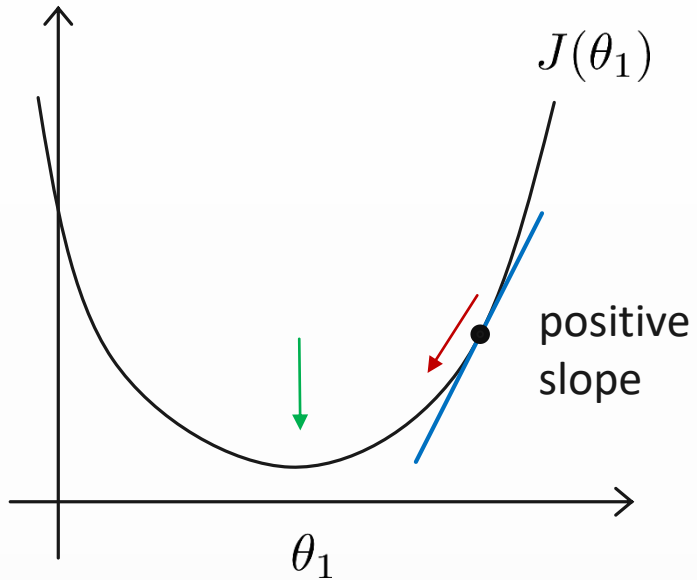
}

learning rate

---

Correct: Simultaneous update

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

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

$\theta_0 := \text{temp0}$
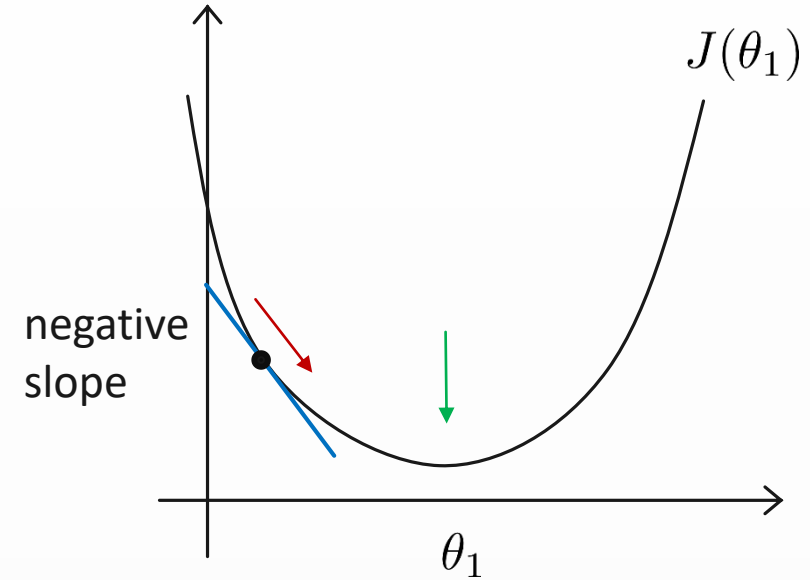
$\theta_1 := \text{temp1}$

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 Intuition



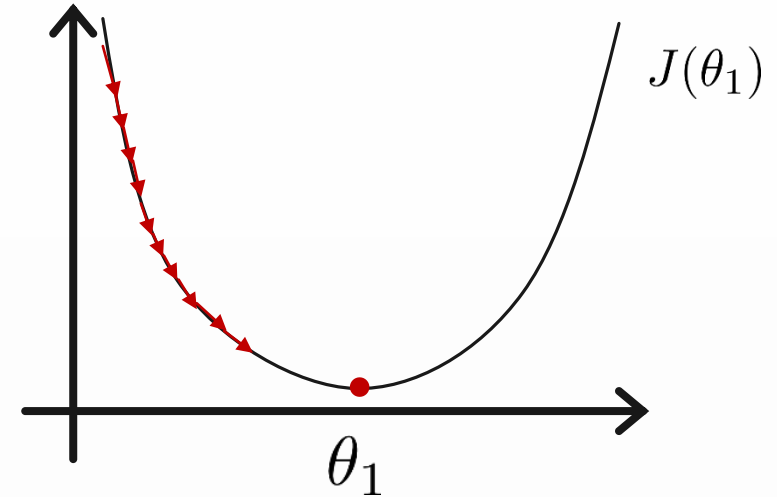$\theta_1 := \theta_1 - \alpha \boxed{\frac{\partial}{\partial \theta_1} J(\theta_1)}$

$\geq 0$

$\theta_1 := \theta_1 - \alpha \boxed{\frac{\partial}{\partial \theta_1} J(\theta_1)}$
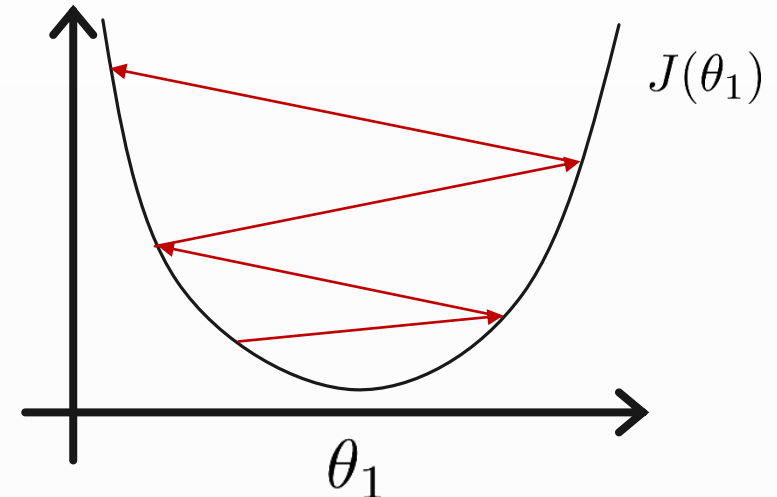
$\leq 0$

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

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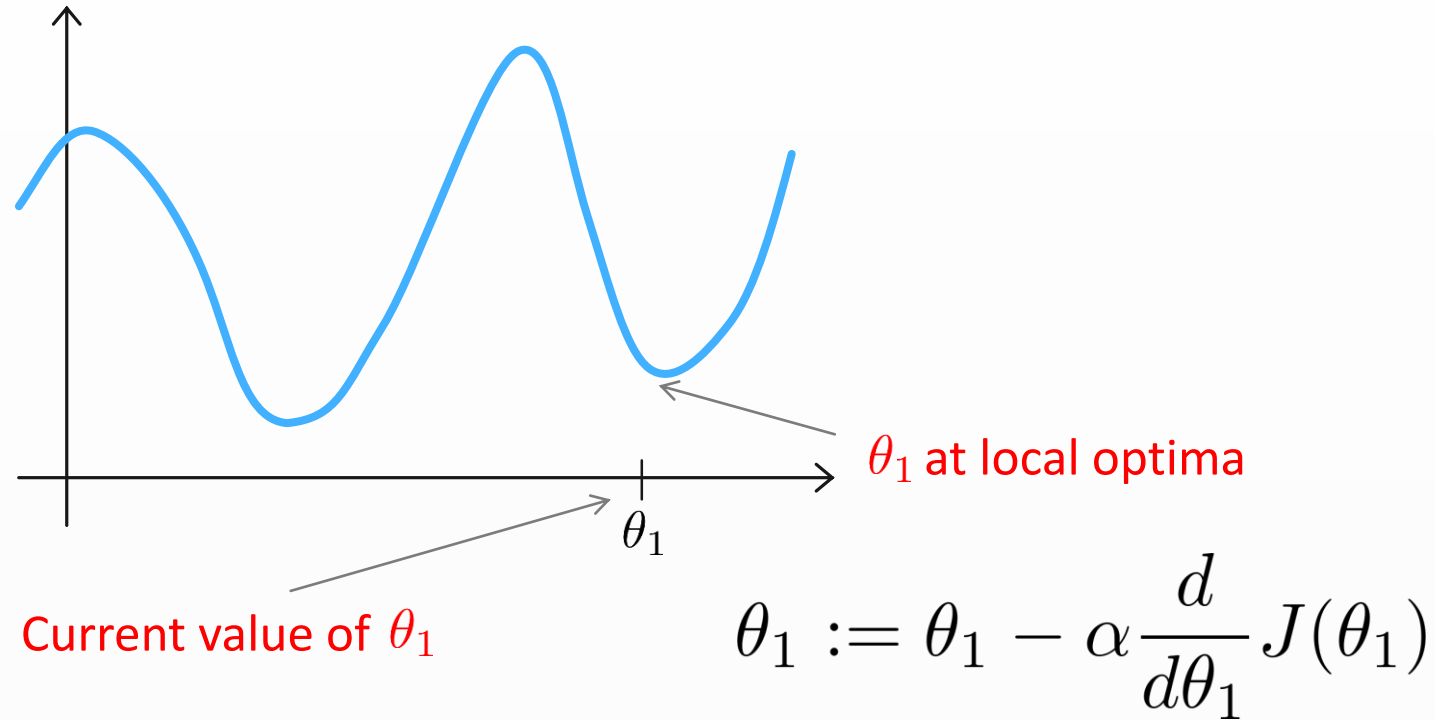
If α is too large, gradient descent can overshoot the minimum. It may fail to converge, or even diverge.

# Gradient Descent Intuition

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



$\theta_1$ at local optima

Current value of $\theta_1$

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

- As we approach a local minimum, gradient descent will automatically take smaller steps

### 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

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

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

# Gradient Descent for Linear Regression

In order to implement this algorithm, we need to calculate the partial derivatives:

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

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

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

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

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

repeat until convergence {

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

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

}

update
$\theta_0$ and $\theta_1$
simultaneously

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

$h_\theta(x)$

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

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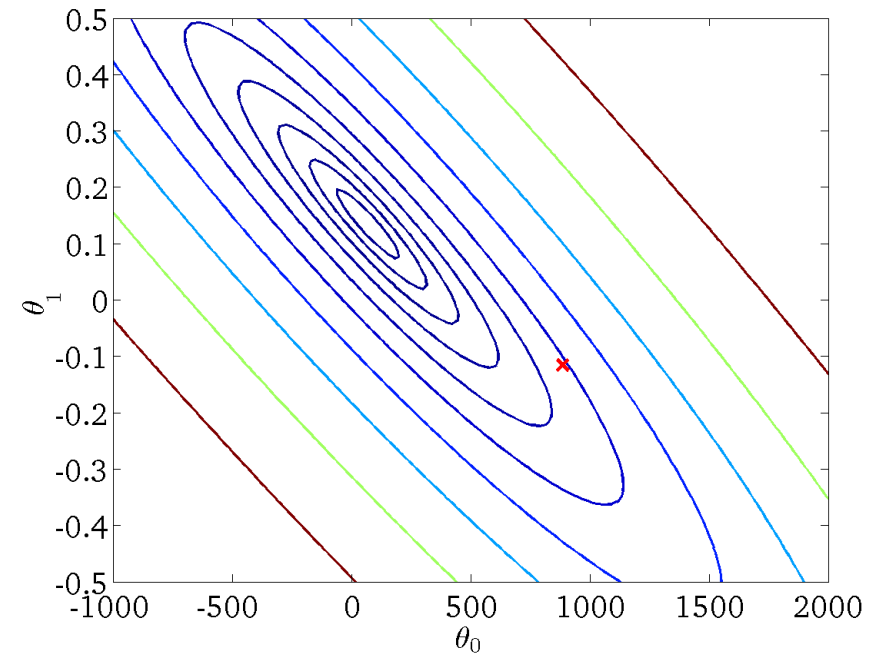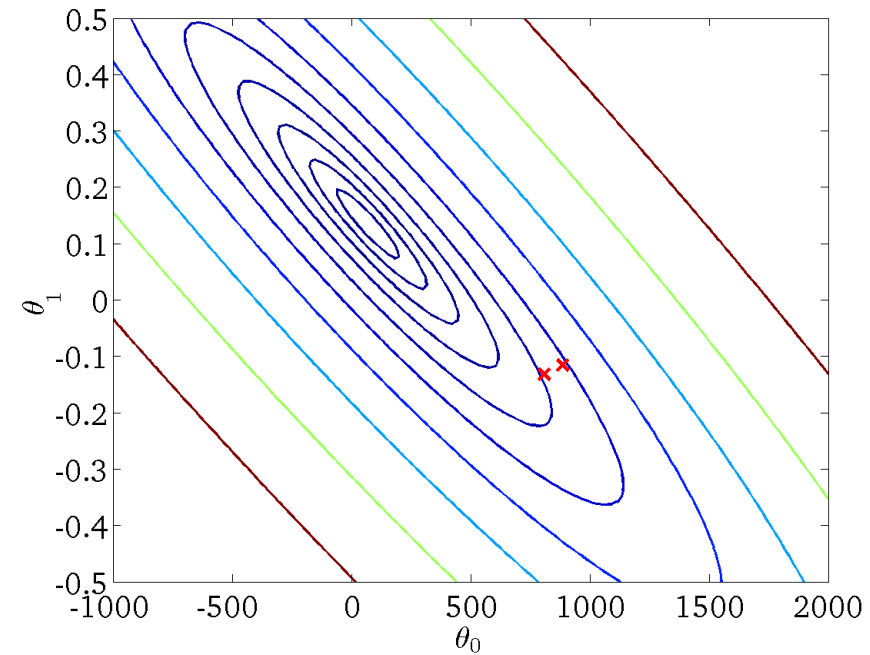
(function of the parameters $\theta_0, \theta_1$)

$$h_\theta(x)$$

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

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

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

$$h_\theta(x)$$

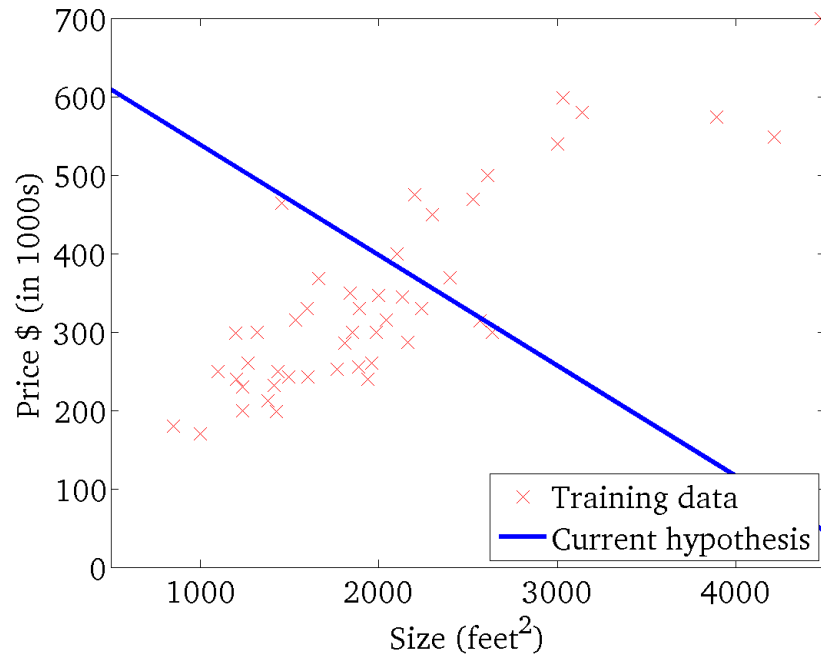(for fixed $\theta_0, \theta_1$, this is a function of x)

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

(function of the parameters $\theta_0, \theta_1$)
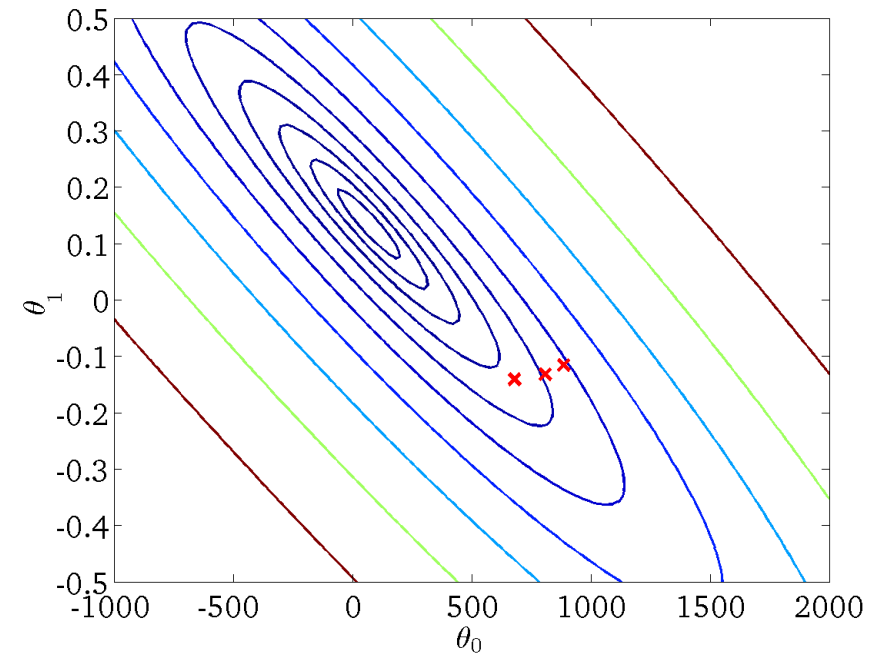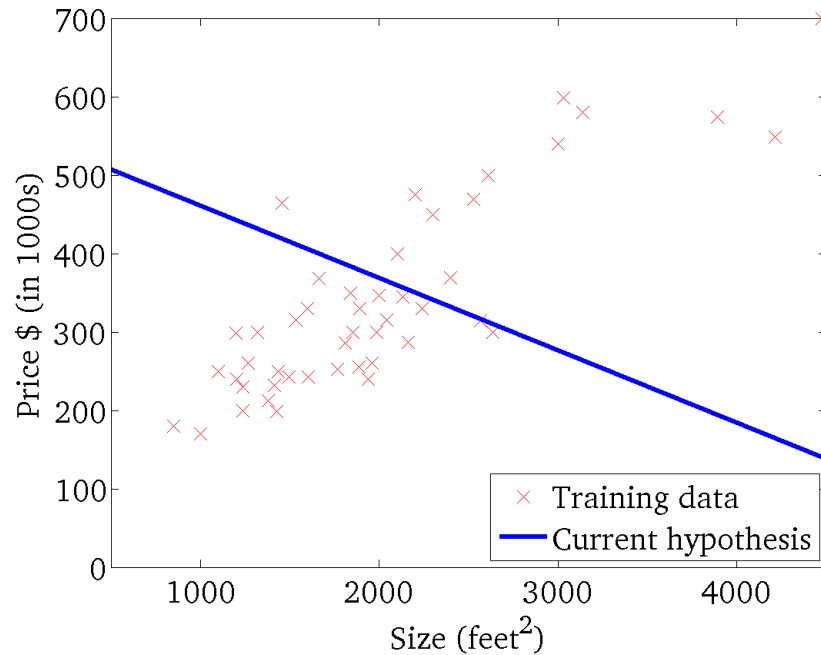
# Gradient Descent for Linear Regression

$$h_\theta(x)$$

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

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

(function of the parameters $\theta_0, \theta_1$)
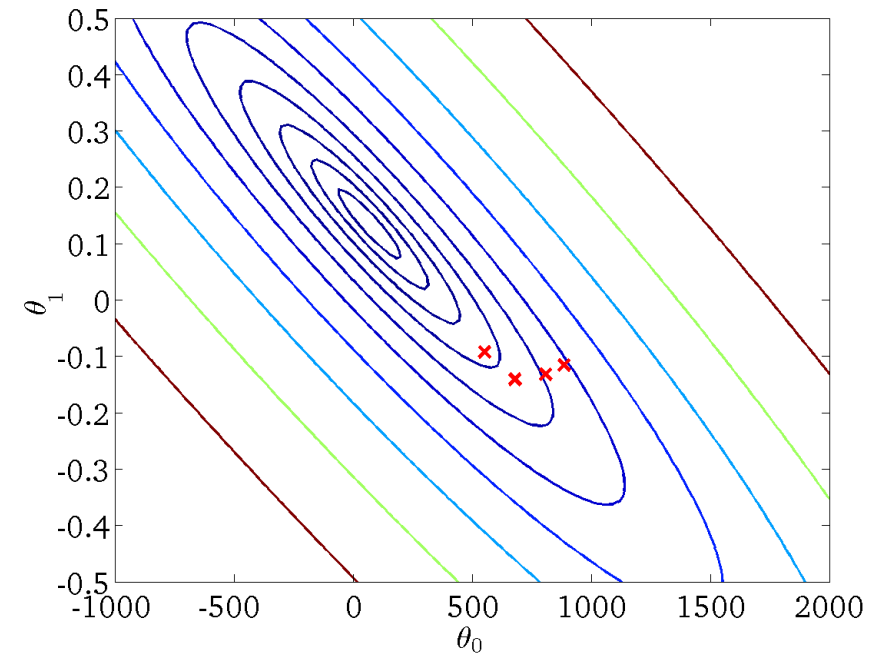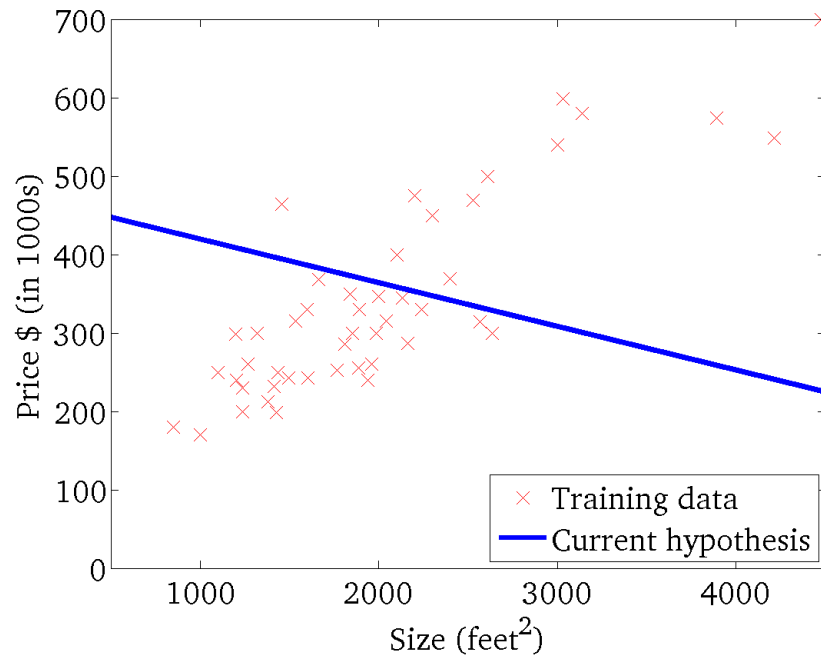
# Gradient Descent for Linear Regression

$h_\theta(x)$

(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 $\theta_0, \theta_1$)

# Gradient Descent for Linear Regression

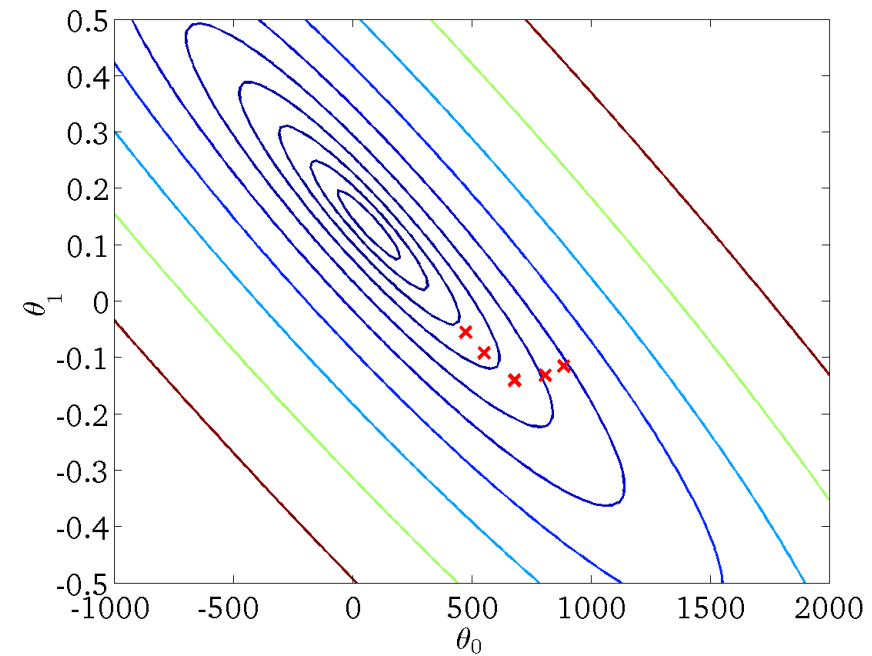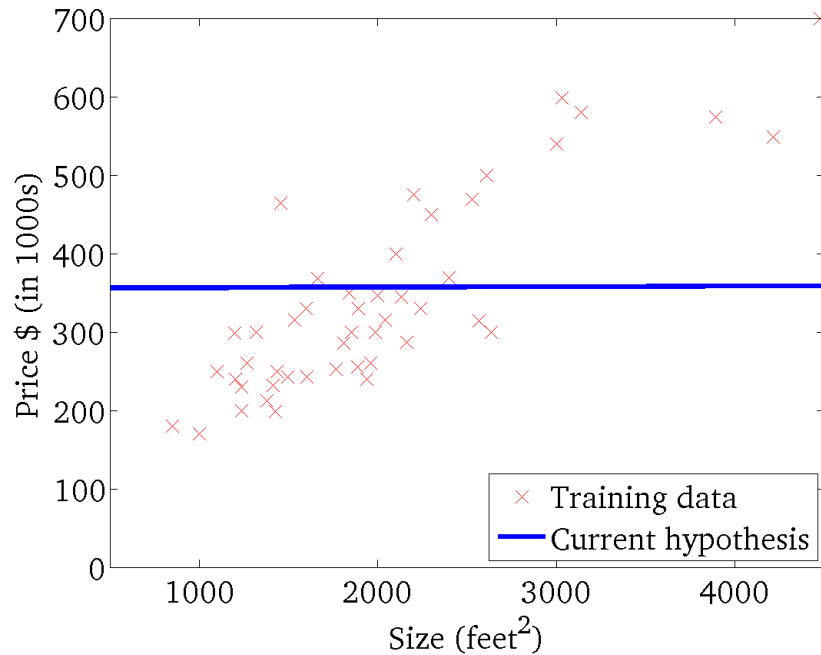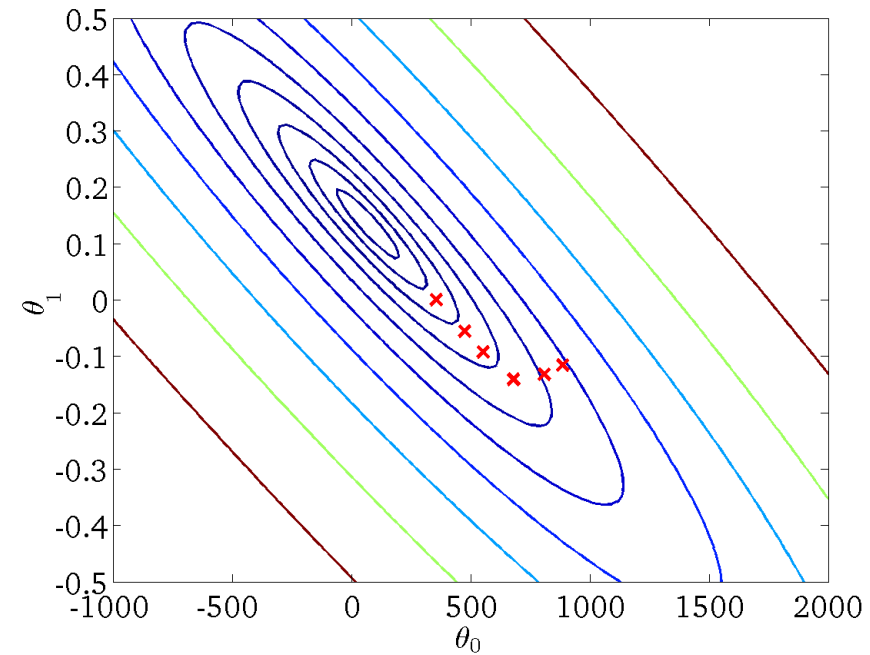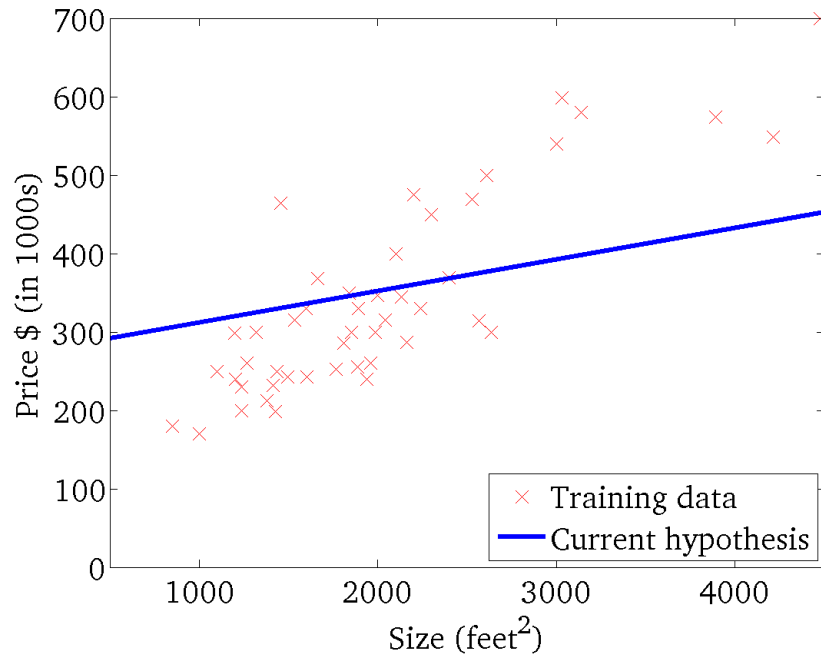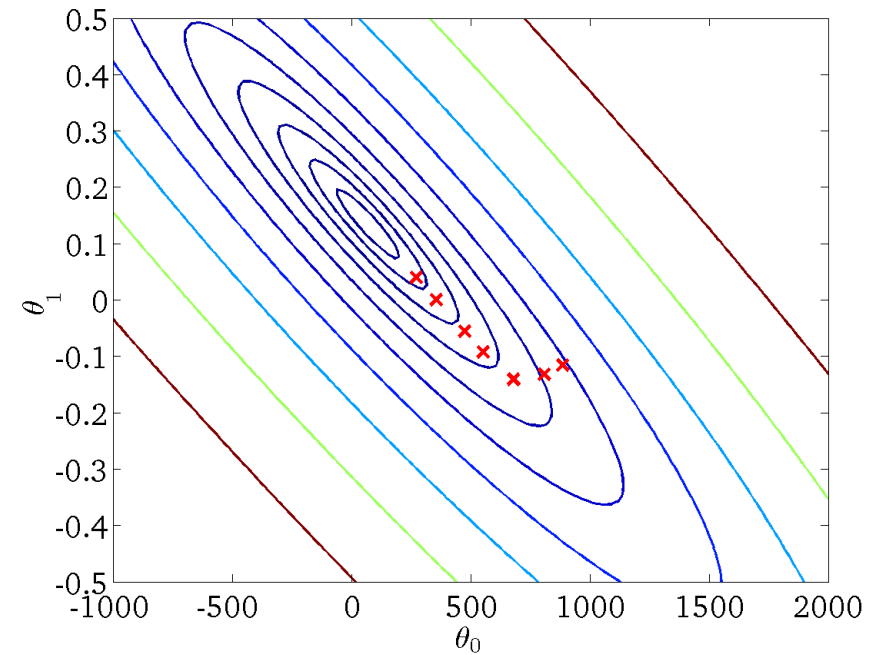# Gradient Descent for Linear Regression

# 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 |
| ... | ... | ... | ... | ... |

Notation:

$n$ = number of features

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

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

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

$$x_3^{(2)} = 2$$

# Multiple Features (Variables)

Hypothesis:

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

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

For convenience of notation, define  $x_0 = 1$

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

Multivariate linear regression

**Gradient Descent**

Previously (n=1):

Repeat $\{$

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

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

(simultaneously update $\theta_0, \theta_1$)

$\}$

New algorithm $(n \geq 1)$:

Repeat $\{$

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

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

$\}$

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

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

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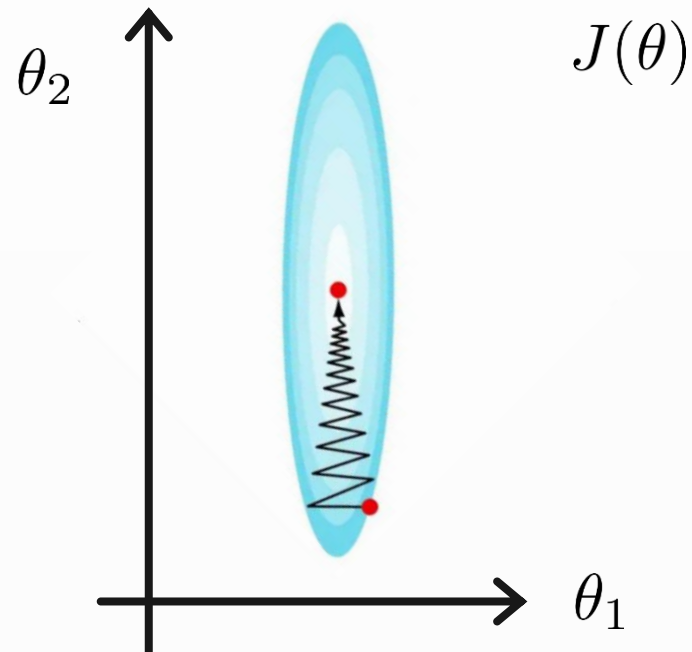
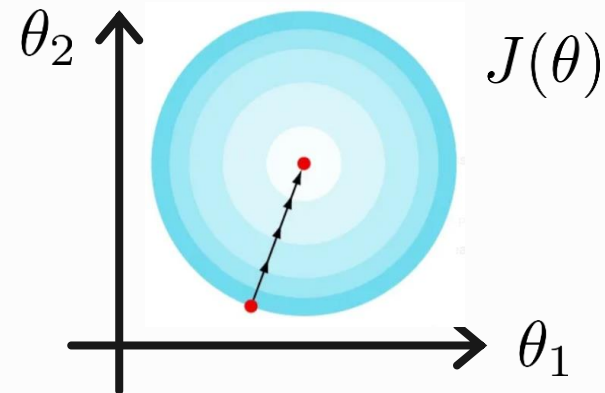$\ldots$

Idea: Make sure features are on a similar scale

E.g. $x_1$ = size (0-2000 feet²)

$x_2$ = number of bedrooms (1-5)

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

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

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

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

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

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

$$x_1 \leftarrow \frac{x_1 - \mu_1}{s_1} \qquad\qquad x_2 \leftarrow \frac{x_2 - \mu_2}{s_2}$$
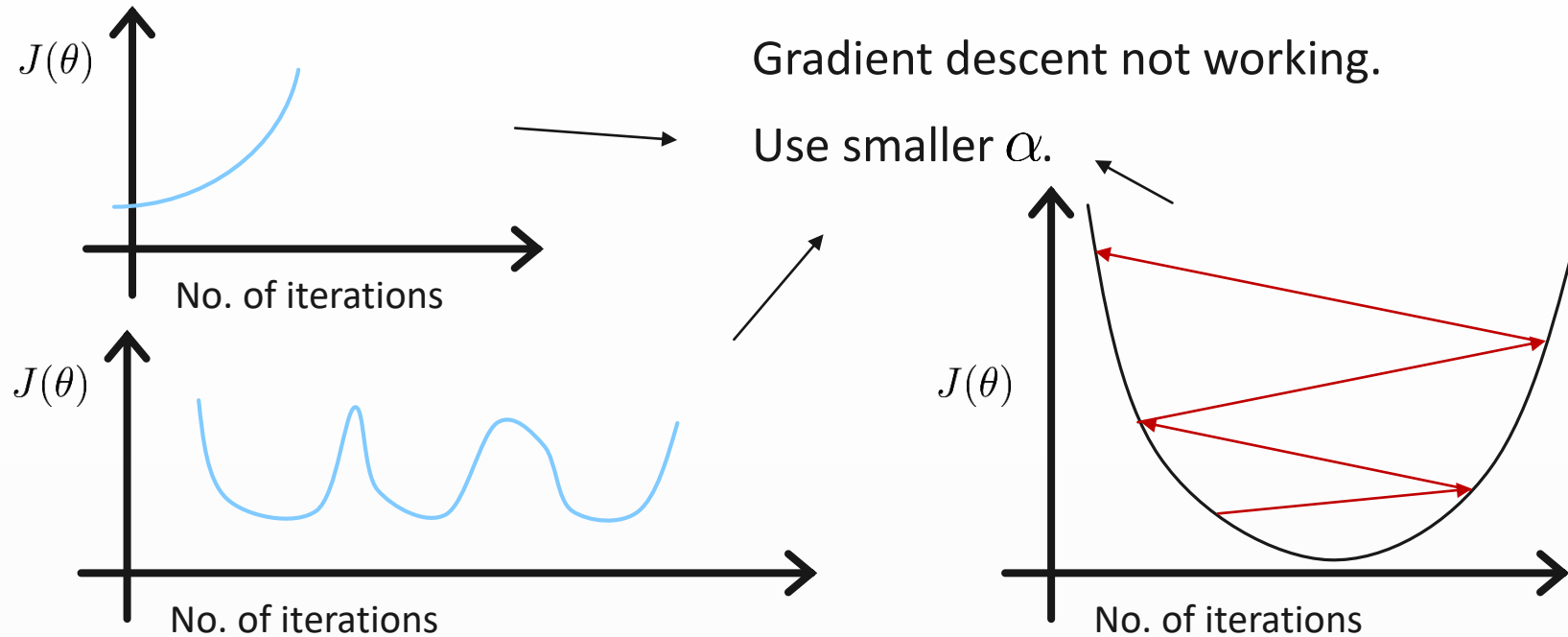
$\mu_i$: average value of $x_i$ in training set

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

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



Gradient descent not working.

Use smaller $\alpha$.

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

- 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

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

$\xrightarrow{\phantom{xxx}}$
3x
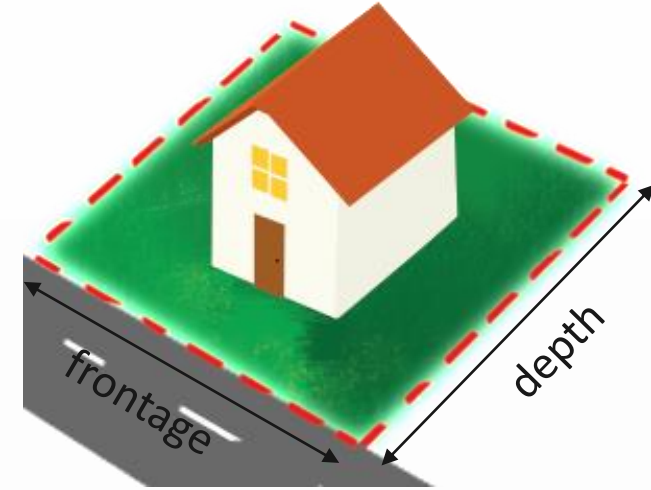
## Housing prices prediction

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



Area:

$x$ = frontage * depth

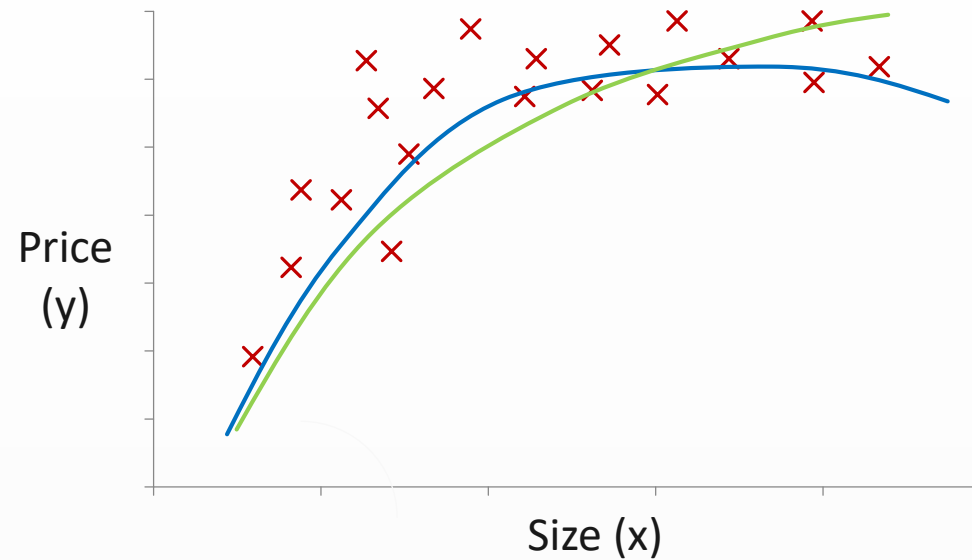$$h_\theta(x) = \theta_0 + \theta_1 x$$

## Choice of features



$$h_\theta(x) = \theta_0 + \theta_1(size) + \theta_2(size)^2$$

$$h_\theta(x) = \theta_0 + \theta_1(size) + \theta_2\sqrt{(size)}$$

Normal equation: Method to solve for analytically.

Intuition: If 1D $(\theta \in \mathbb{R})$

$$J(\theta) = a\theta^2 + b\theta + c$$

$$\frac{\partial}{\partial \theta} J(\theta) = 0$$

Solve for $\theta$



$\theta \in \mathbb{R}^{n+1}$ $\qquad J(\theta_0, \theta_1, \ldots, \theta_m) = \frac{1}{2m} \sum_{i=1}^{m} (h_\theta(x^{(i)}) - y^{(i)})^2$

$$\frac{\partial}{\partial \theta_j} J(\theta) = \cdots = 0 \qquad \text{(for every } j\text{)}$$

Solve for $\theta_0, \theta_1, \ldots, \theta_n$

**Examples:** $m = 5$.

| $x_0$ | Size (feet²) $x_1$ | Number of bedrooms $x_2$ | Number of floors $x_3$ | Age of home (years) $x_4$ | Price ($1000) $y$ |
|---|---|---|---|---|---|
| 1 | 2104 | 5 | 1 | 45 | 460 |
| 1 | 1416 | 3 | 2 | 40 | 232 |
| 1 | 1534 | 3 | 2 | 30 | 315 |
| 1 | 852 | 2 | 1 | 36 | 178 |
| 1 | | | | | |

$$X = \begin{bmatrix} 1 & 2104 & 5 & 1 & 45 \\ 1 & 1416 & 3 & 2 & 40 \\ 1 & 1534 & 3 & 2 & 30 \\ 1 & 852 & 2 & 1 & 36 \\ 1 & 3000 & 4 & 1 & 38 \end{bmatrix} \qquad y = \begin{bmatrix} 460 \\ 232 \\ 315 \\ 178 \\ 540 \end{bmatrix}$$

$$\Theta = (X^T X)^{-1} X^T y$$

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

### Gradient Descent

- Need to choose $\alpha$.
- Needs many iterations.
- Works well even when $n$ is large.

### Normal Equation

- No need to choose $\alpha$.
- Don't need to iterate.
- Need to compute $(X^T X)^{-1}$
- Slow if $n$ is very large.

# References

- A. Ng. Machine Learning, Lecture Notes.

- I. Goodfellow, Y. Bengio and A. Courville, "Deep Learning", 2016.