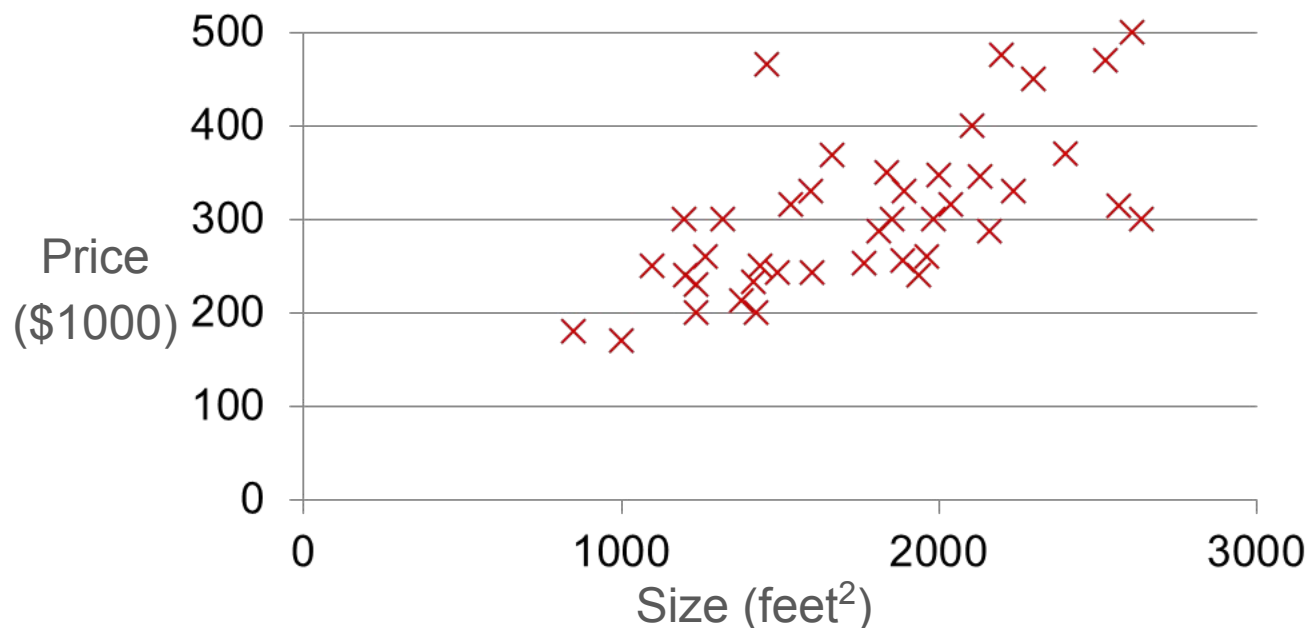# CENG 463
# Machine Learning

## Lecture 04 - Linear Regression

# Model with One Variable

**Example:** House prices according to area

- Supervised Learning: The "right answer" for each example in the data is given.
- Regression Problem:  Predict real-valued output

# Model with One Variable

Training set for house pricing:

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

**Notation:**

**m** = number of training examples

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

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

**(x,y)** = one training example

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

# Model with One Variable



**h** is a function from **x** to **y**.
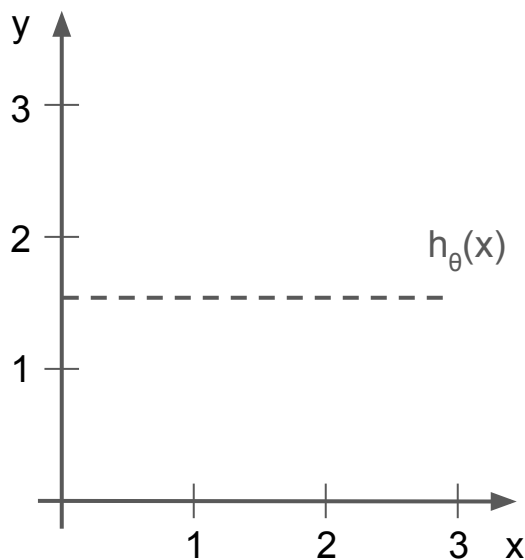How can we represent **h**?
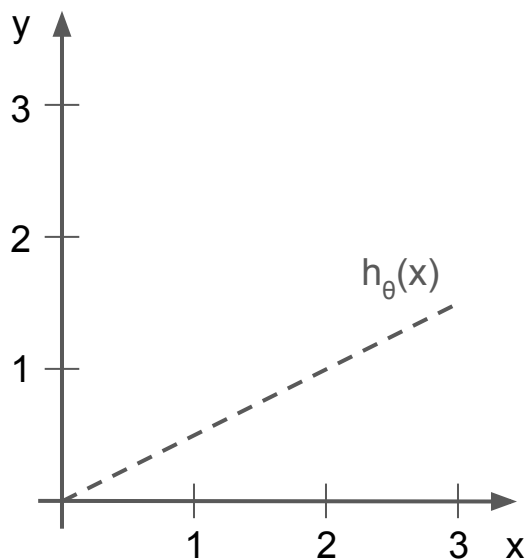
$$h_\theta(x) = \theta_0 + \theta_1 x$$

Linear regression with one variable.
i.e. Univariate linear regression.
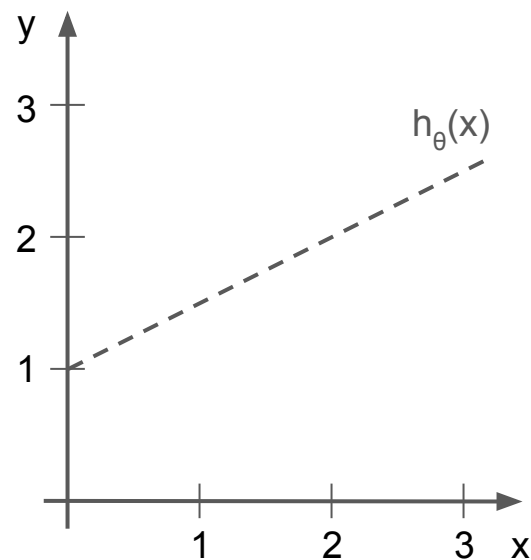
# Model with One Variable

- Hypothesis: $h_\theta(x) = \theta_0 + \theta_1 x$
- $\theta_i$'s are **parameters**.



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

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

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

# Cost Function

- Hypothesis: $h_\theta(x) = \theta_0 + \theta_1 x$
- $\theta_i$'s are **parameters**. How to choose them?
- Choose $\theta_i$'s so that **$h_\theta(x)$** is close to **y** <u>for our training examples</u>.

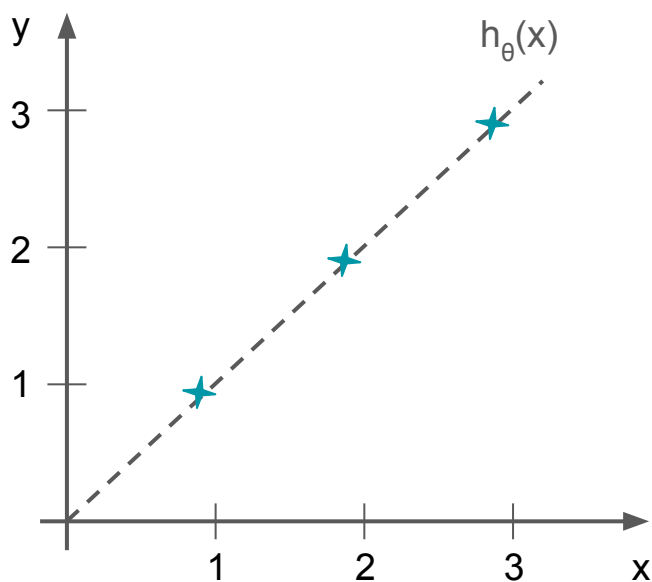$$\min_{\theta_0 \theta_1} \frac{1}{2m} \sum_{i=1}^{m} \left( h_\theta(x^{(i)}) - y^{(i)} \right)^2$$

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

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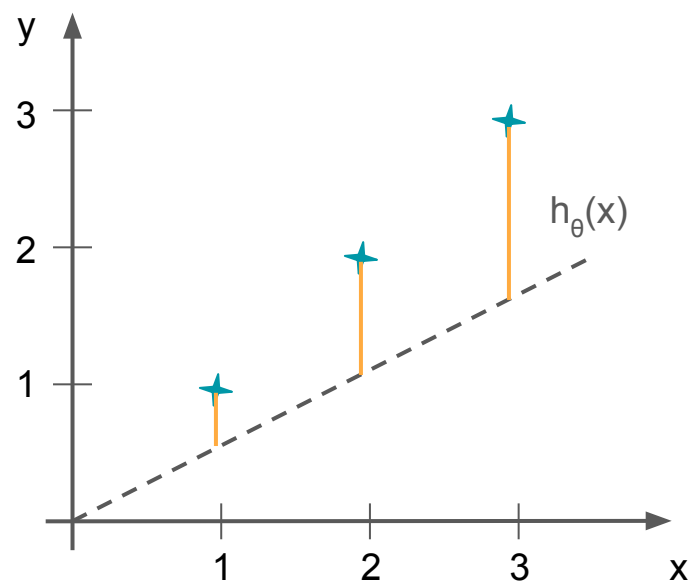
# Cost Function

$$J(\theta_0, \theta_1) = \frac{1}{2m} \sum_{i=1}^{m} \left( h_\theta(x^{(i)}) - y^{(i)} \right)^2$$
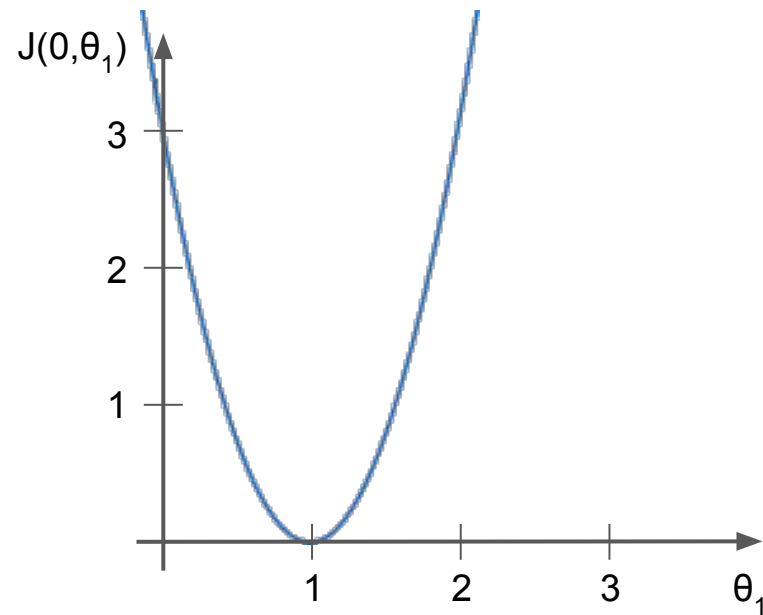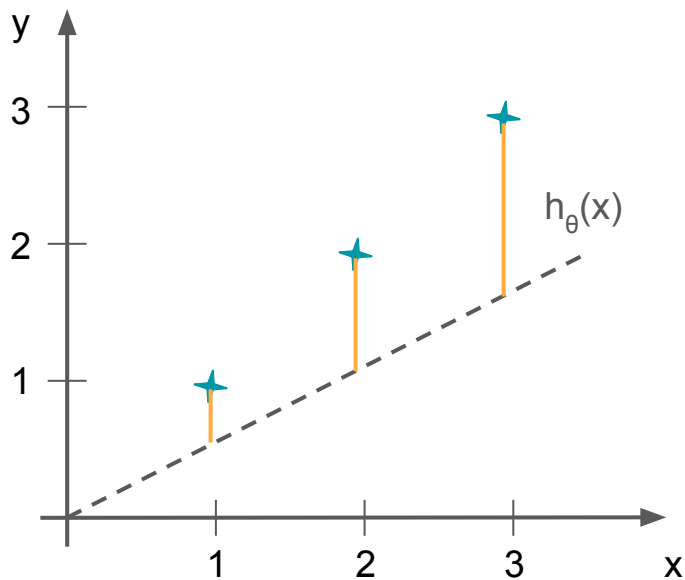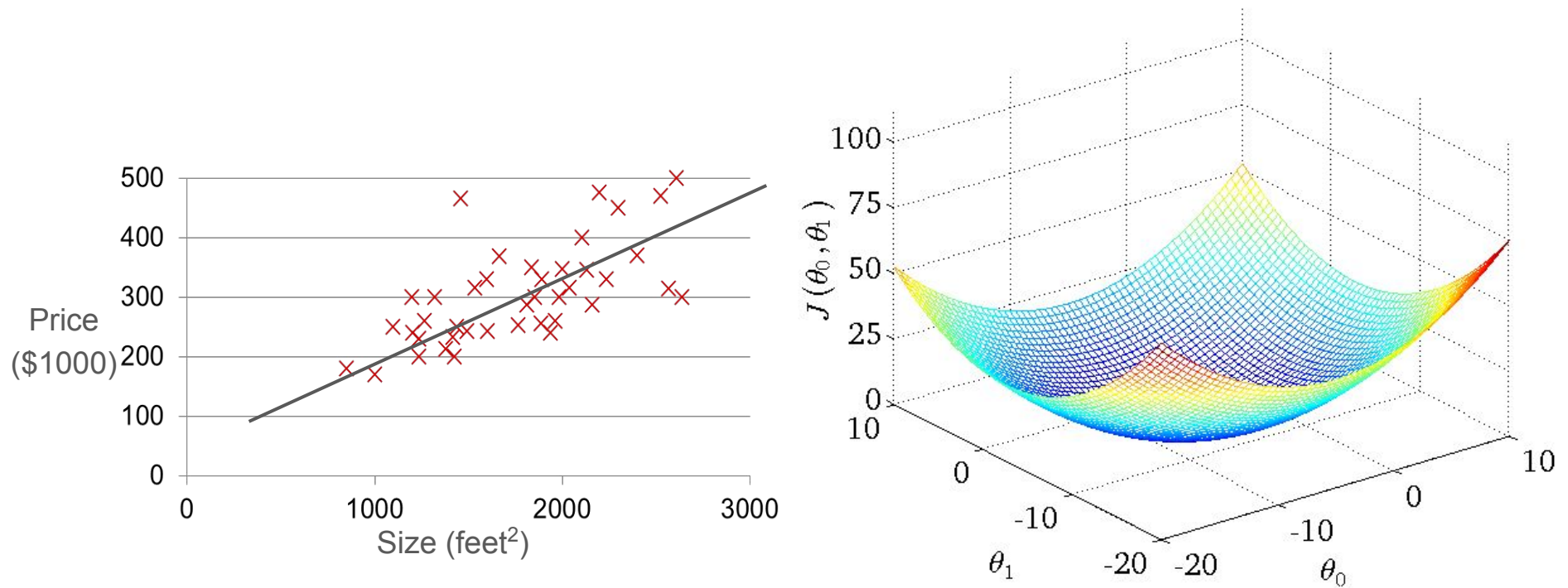
What is J(0,1)?

What is J(0,0.5)?

# Cost Function

Change of cost function according to $\theta_1$:

# Cost Function

Change of cost function according to $\theta_0$ and $\theta_1$:

# Cost Function

Cost with changing $\theta_0$ and $\theta_1$ (shown with a contour plot)
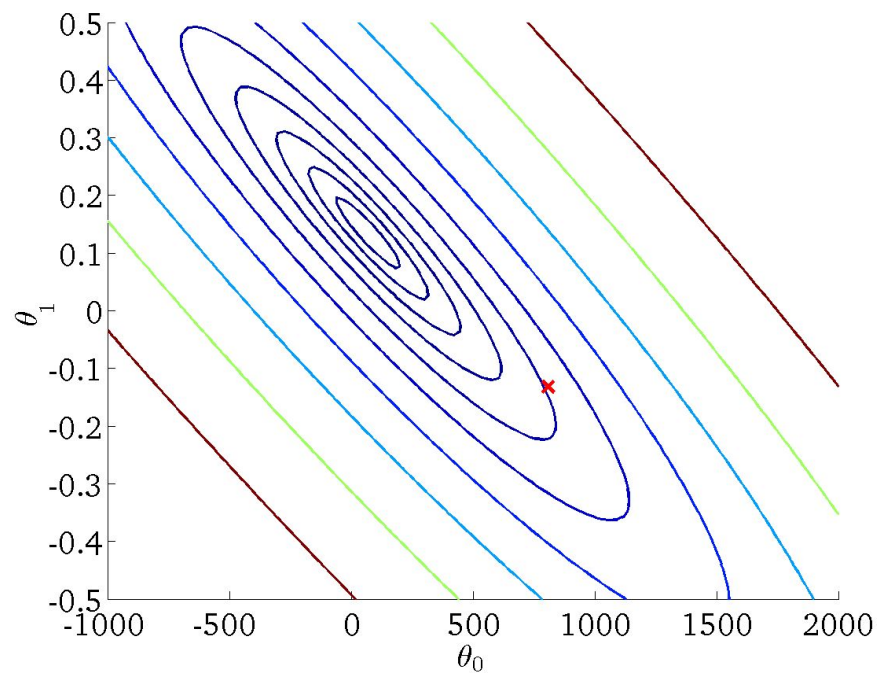
# Cost Function

Cost with changing $\theta_0$ and $\theta_1$ (shown with a contour plot)

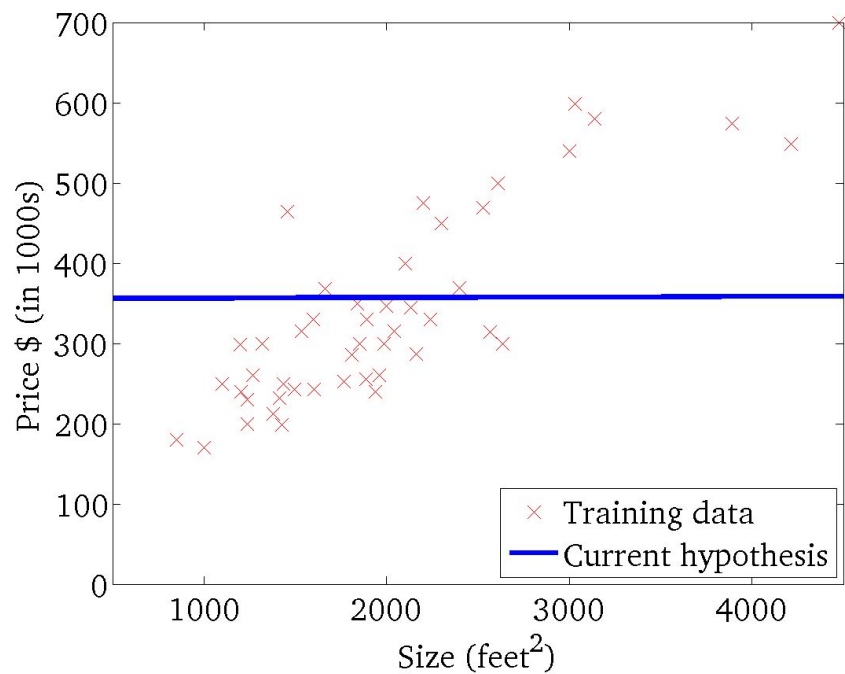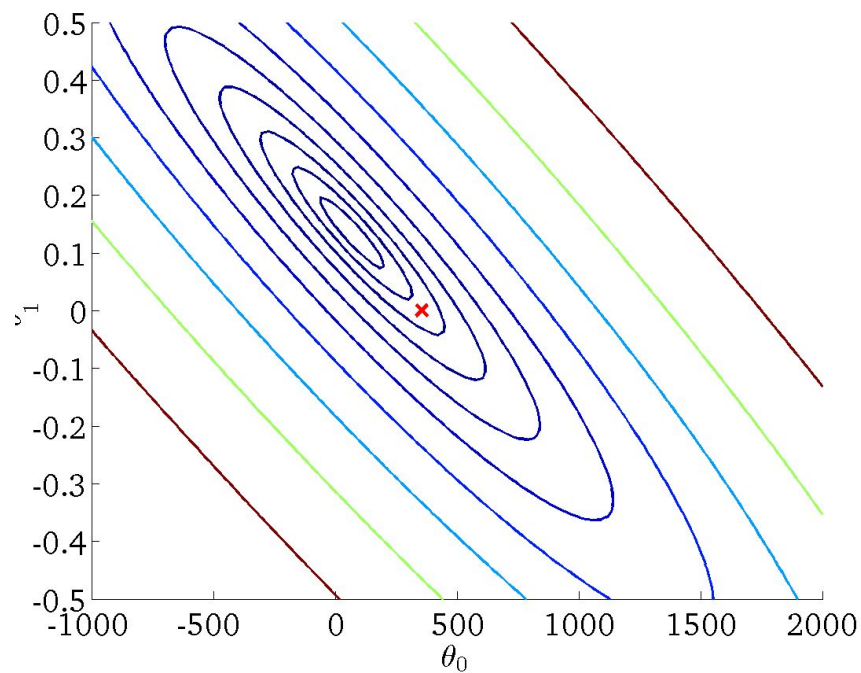$h_\theta(x)$                                    $J(\theta_0, \theta_1)$

# Cost Function

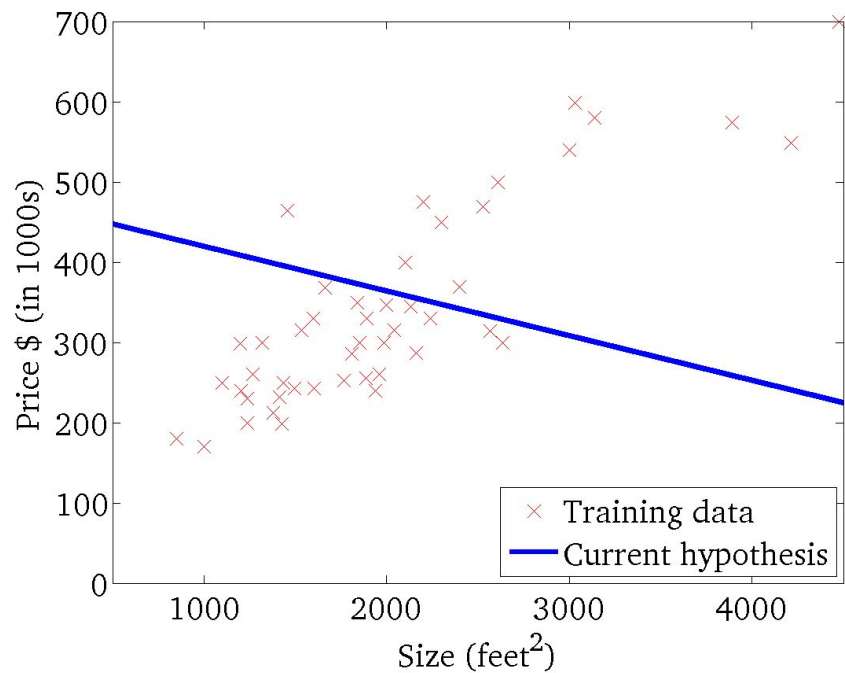Cost with changing $\theta_0$ and $\theta_1$ (shown with a contour plot)
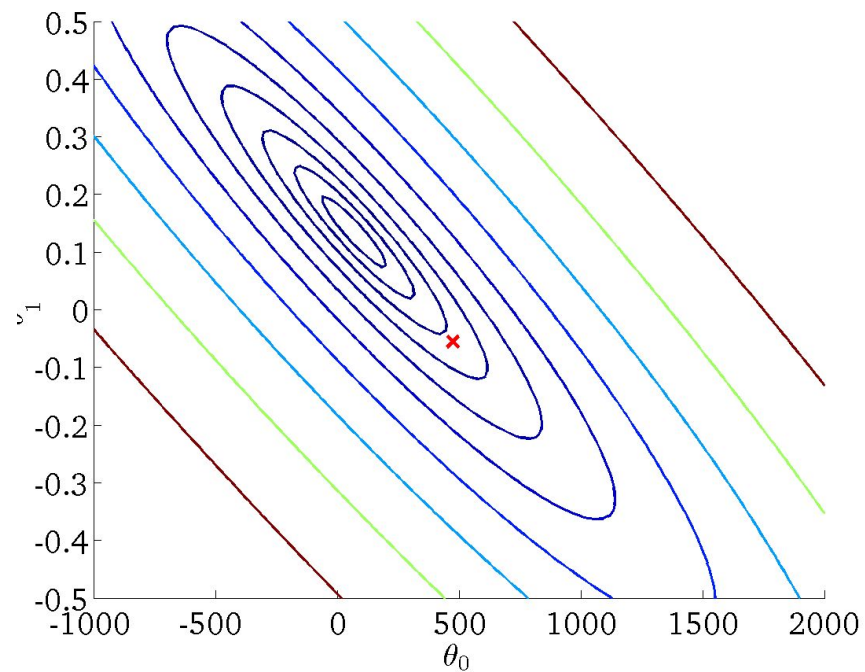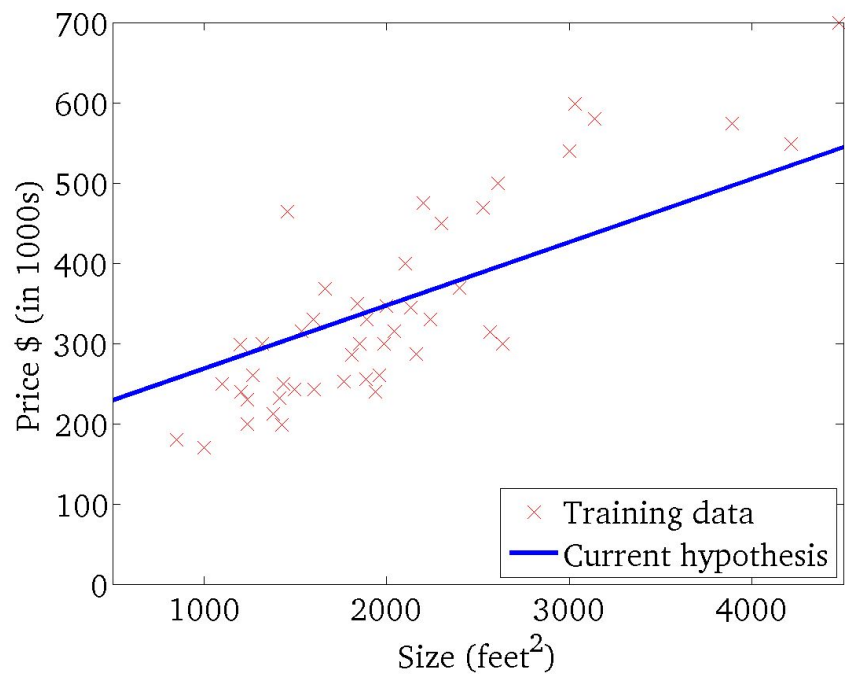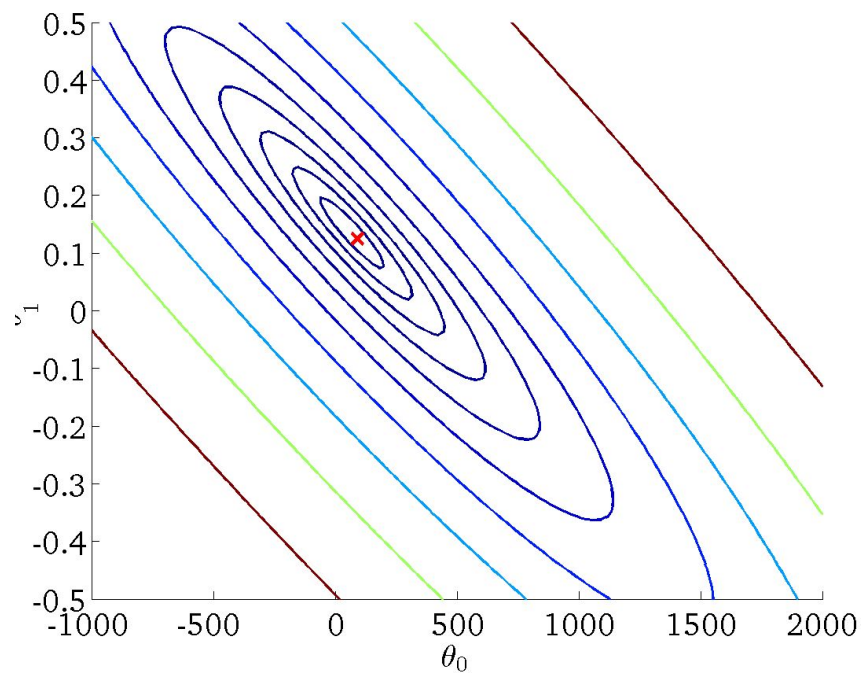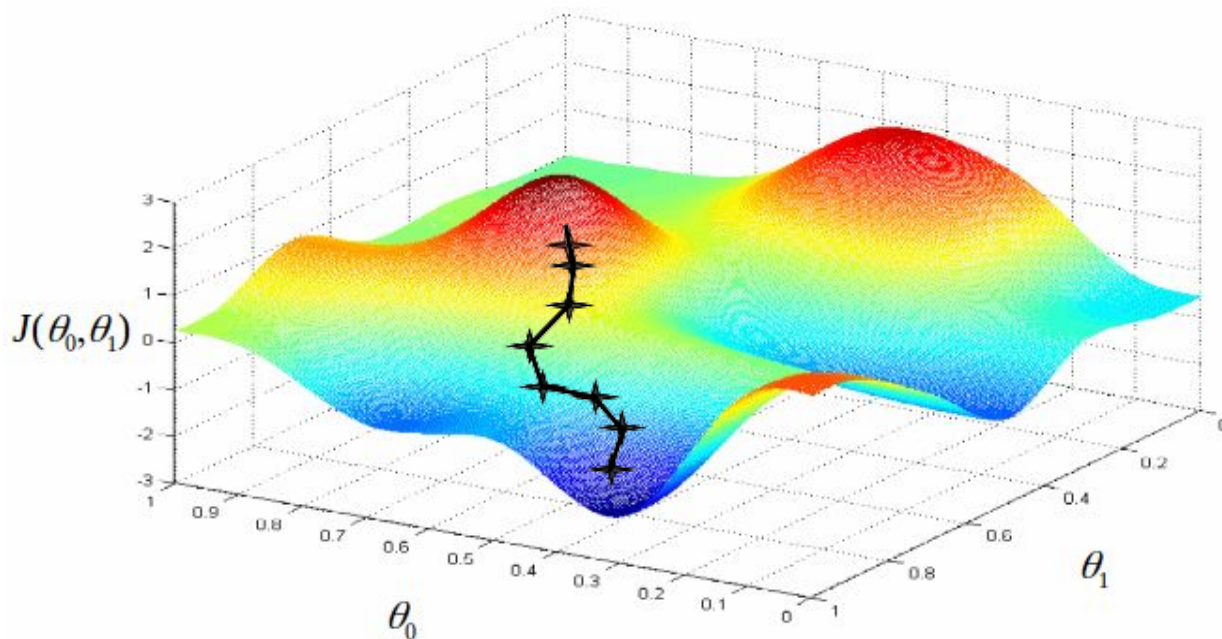
$h_\theta(x)$ $J(\theta_0,\theta_1)$

# Cost Function

Cost with changing $\theta_0$ and $\theta_1$ (shown with a contour plot)
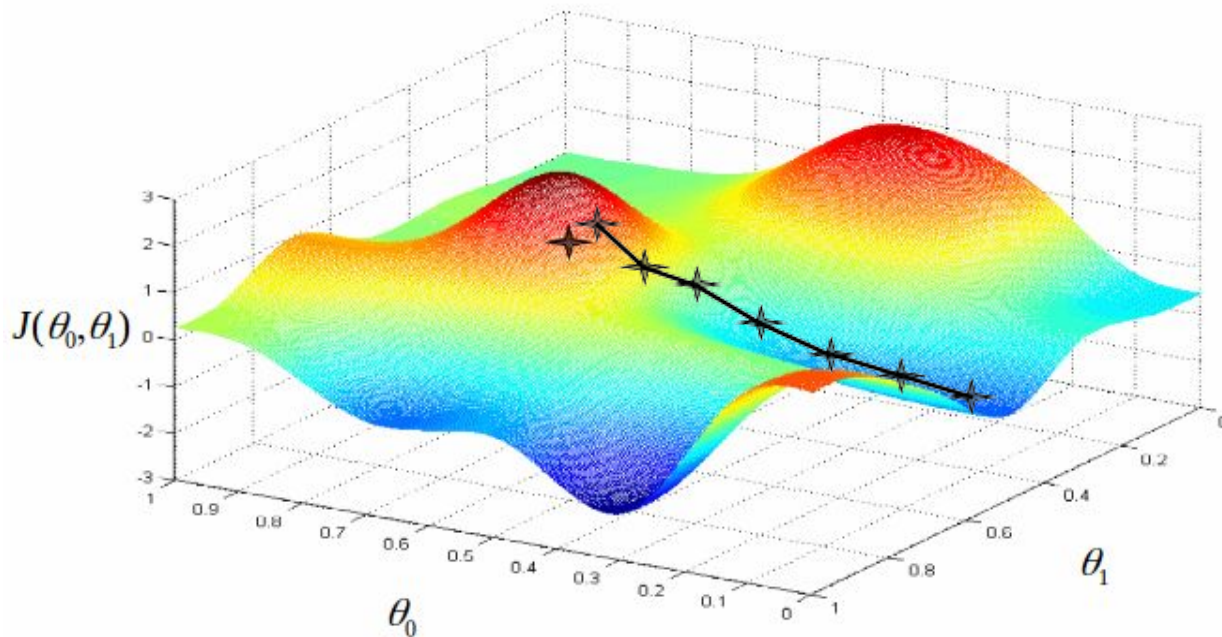
# Gradient Descent

- We have some function $J(\theta_0, \theta_1)$ that we want to minimize.
  - Find $\theta_0$, $\theta_1$ parameters that minimize J:
    - Start with some $\theta_0$, $\theta_1$.
    - Keep changing $\theta_0$, $\theta_1$ to reduce $J(\theta_0, \theta_1)$ until hopefully we end up at a minimum.

# Gradient Descent

- Gradient descent does **NOT** guarantee to reach the **global minimum**!

# Gradient Descent

- Gradient descent algorithm for two-parameter case:

*repeat until convergence {*

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

*}*

This is the derivative term, indicates the direction of step.

This number is 'learning rate', controls the size of the step we take.
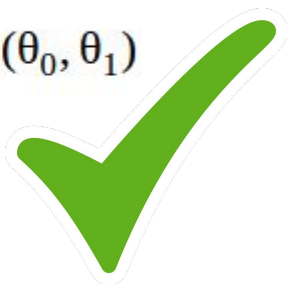
**Correct: Simultaneous update**

$$temp_0 := \theta_0 - \alpha \frac{\partial}{\partial \theta_0} J(\theta_0, \theta_1)$$

$$temp_1 := \theta_1 - \alpha \frac{\partial}{\partial \theta_1} J(\theta_0, \theta_1)$$

$$\theta_0 := temp_0$$

$$\theta_1 := temp_1$$

**Incorrect: Successive update**

$$temp_0 := \theta_0 - \alpha \frac{\partial}{\partial \theta_0} J(\theta_0, \theta_1)$$

$$\theta_0 := temp_0$$

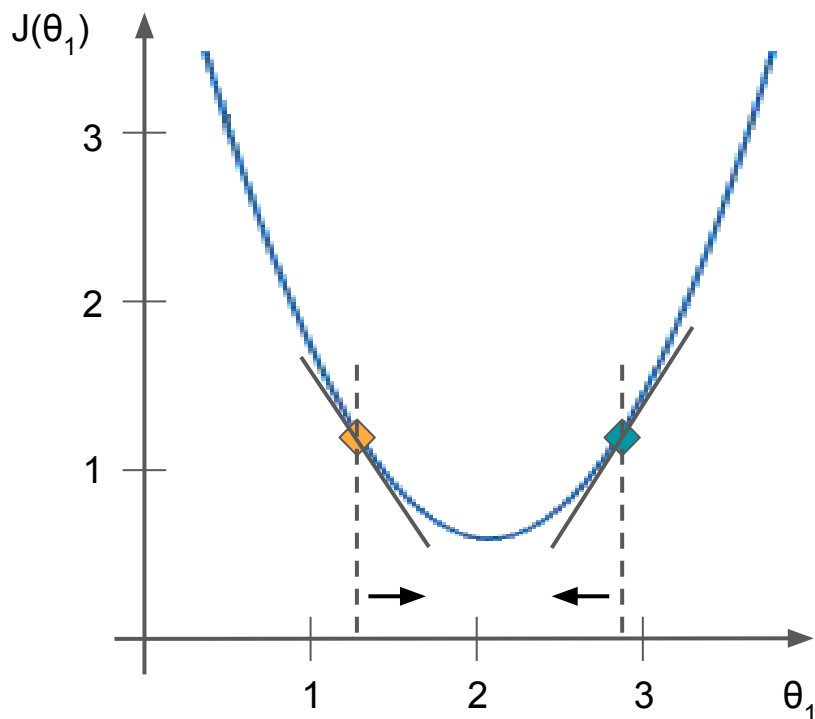$$temp_1 := \theta_1 - \alpha \frac{\partial}{\partial \theta_1} J(\theta_0, \theta_1)$$

$$\theta_1 := temp_1$$

# Gradient Descent

- Let's assume a cost function with one variable ($\theta_1$):

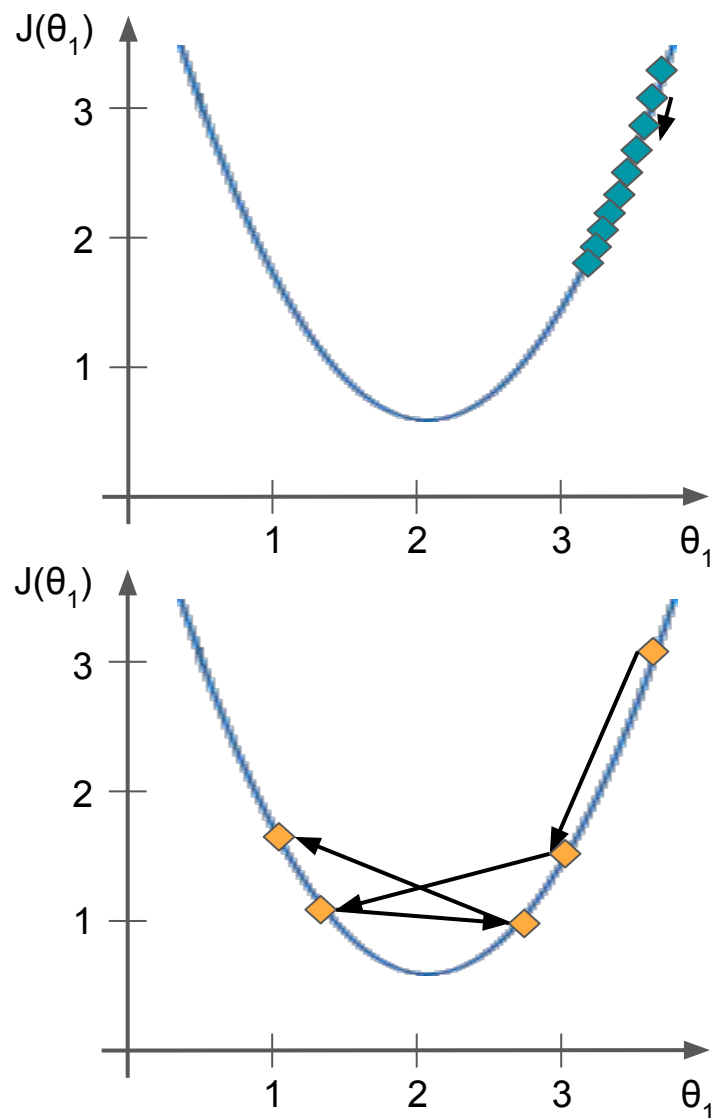$$\theta_1 := \theta_1 - \alpha \frac{\partial}{\partial \theta_1} J(\theta_1)$$

- The derivative gives the slope at that point.
  - Positive slope decreases the value of $\theta_1$.
  - Negative slope increases $\theta_1$.

# Gradient Descent

$$\theta_1 := \theta_1 - \alpha \frac{\partial}{\partial \theta_1} J(\theta_1)$$
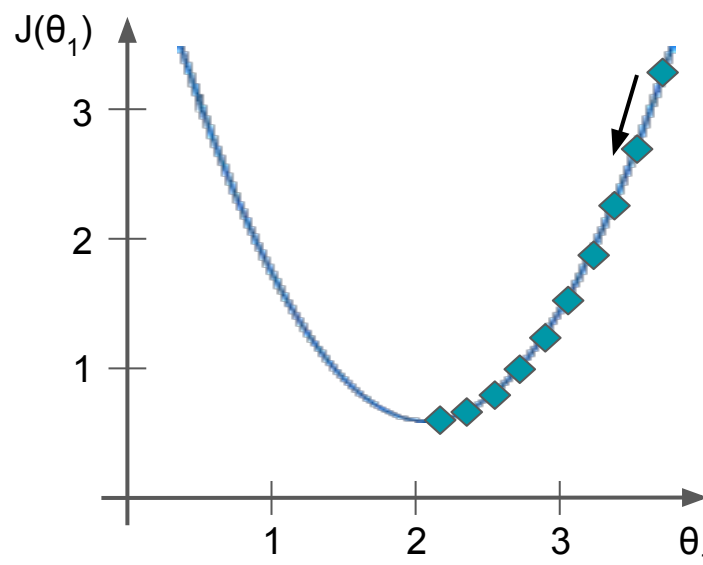
- If **α** is <u>too small</u>, gradient descent can be slow.
- If **α** is <u>too large</u>, gradient descent can overshoot the minimum. It may fail to converge, or even diverge.
- What does gradient descent do if we are already at the local minimum point?

# Gradient Descent

- Gradient descent can converge to a local minimum, even with the learning rate α <u>fixed</u>.
- As we approach a local minimum, gradient descent will automatically take smaller steps. So, no need to decrease α over time.

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

# Gradient Descent for Linear Regression

- We merge two things we have learned:

### Gradient Descent

*repeat until convergence {*

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

*}*

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

### Linear Regression

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

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

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

# Gradient Descent for Linear Regression

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

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

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

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

# Gradient Descent for Linear Regression

repeat until convergence {

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

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

}

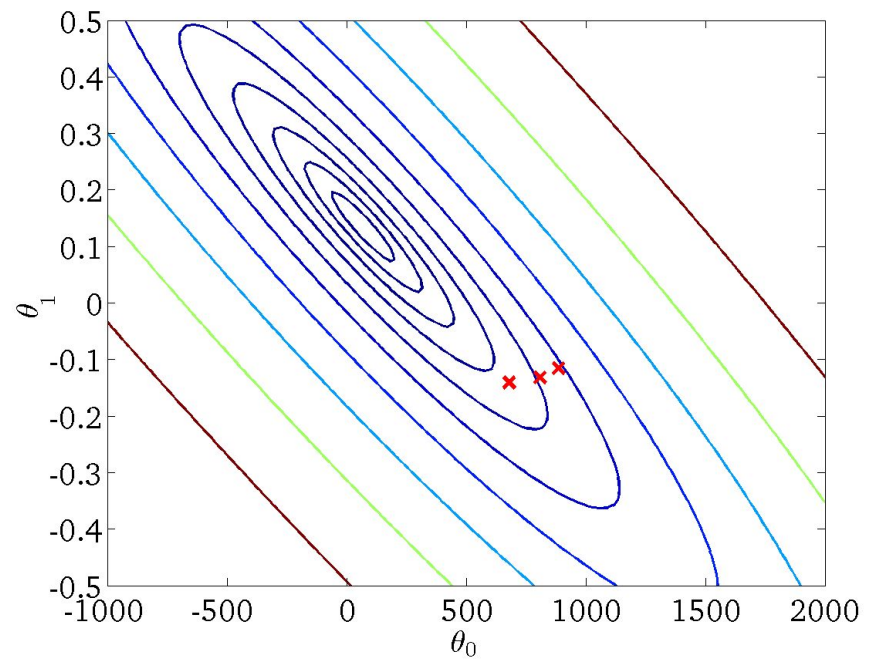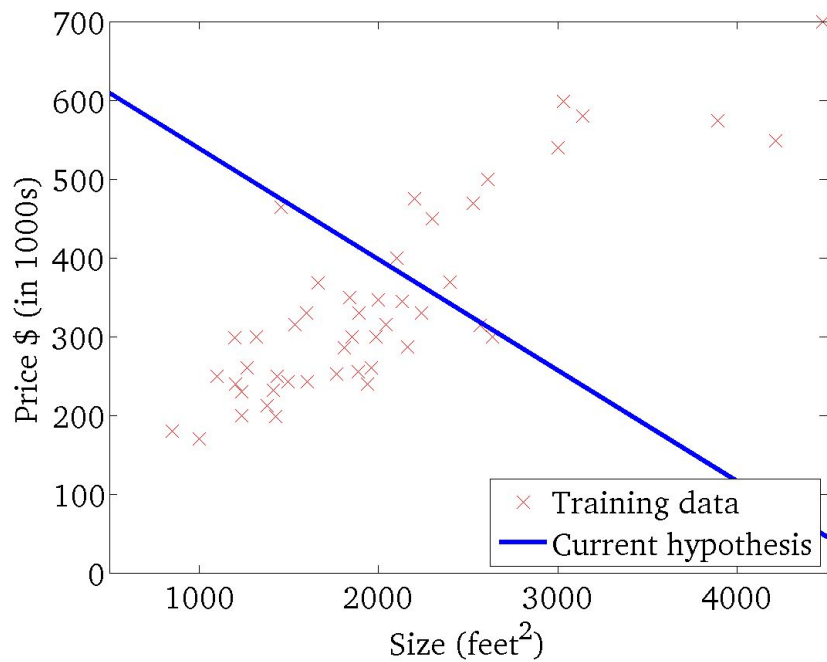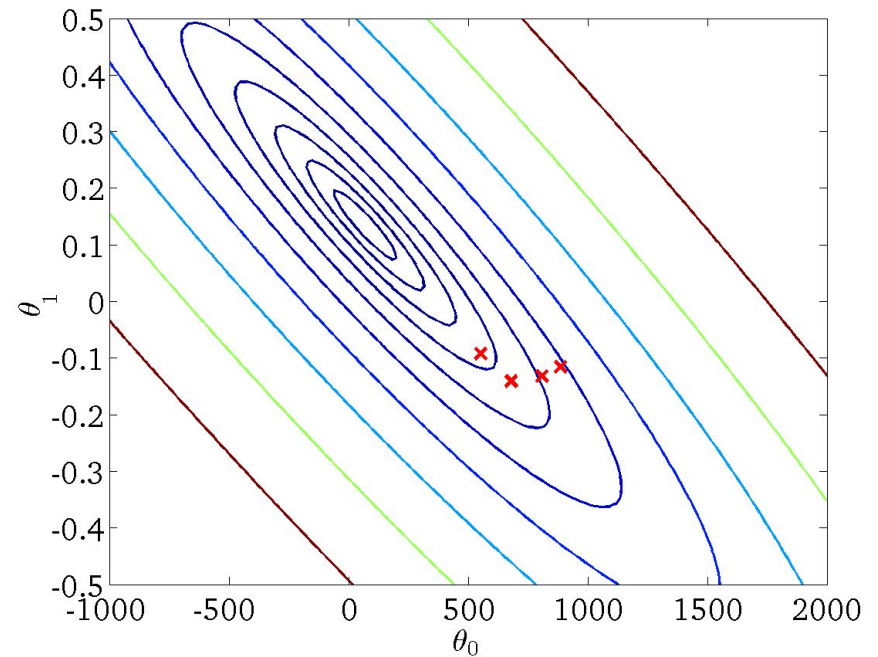Linear regression has a 'convex' (i.e. bowl-shaped) cost function, so we expect to reach global minimum.
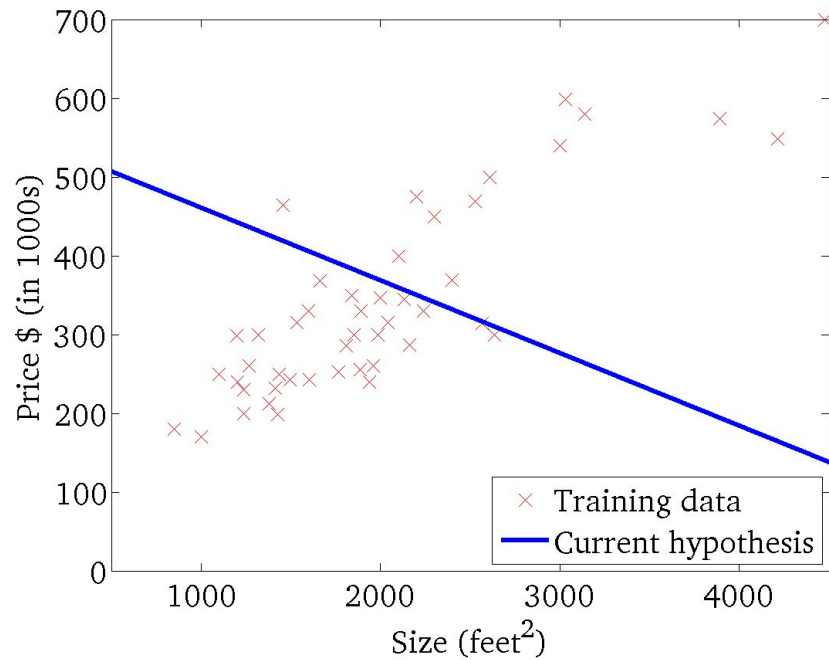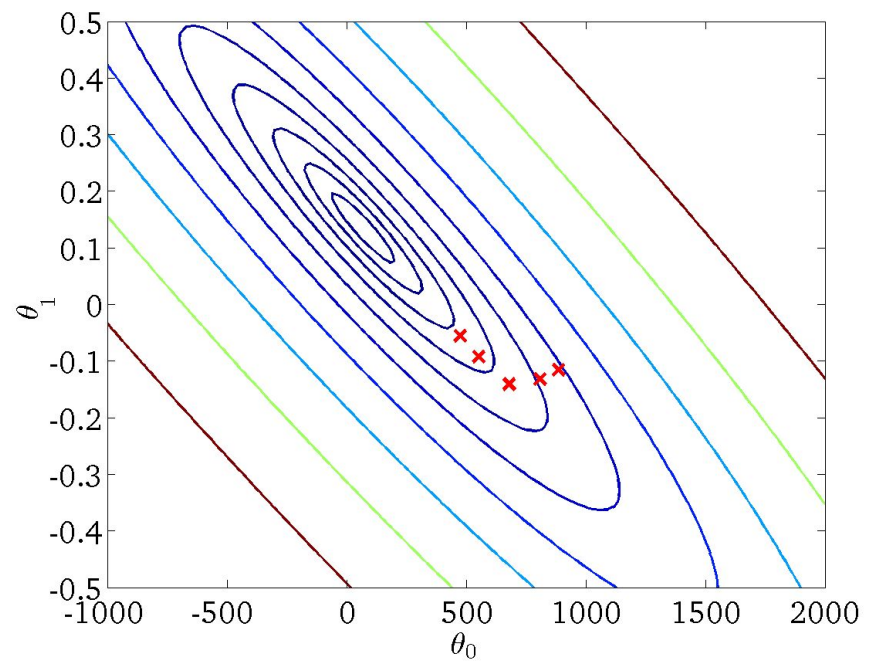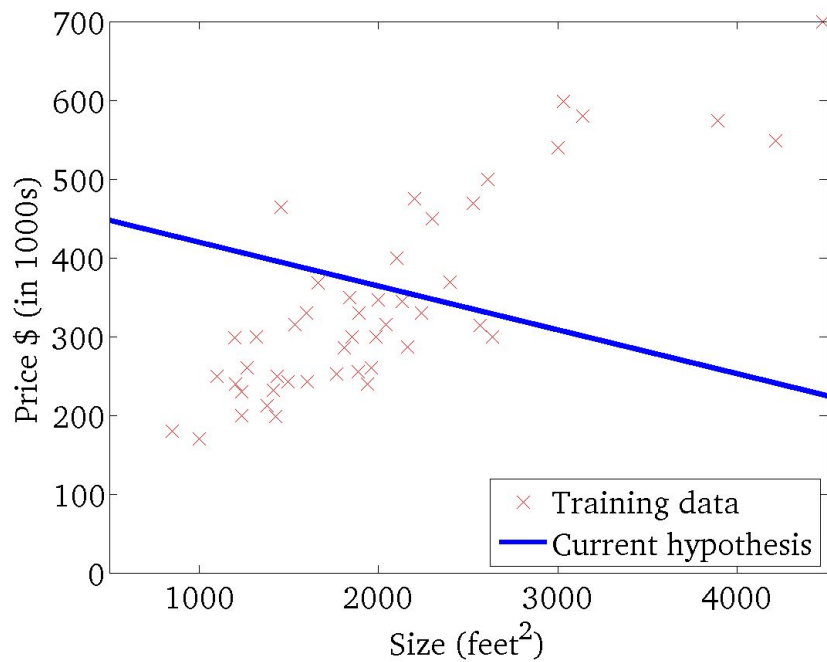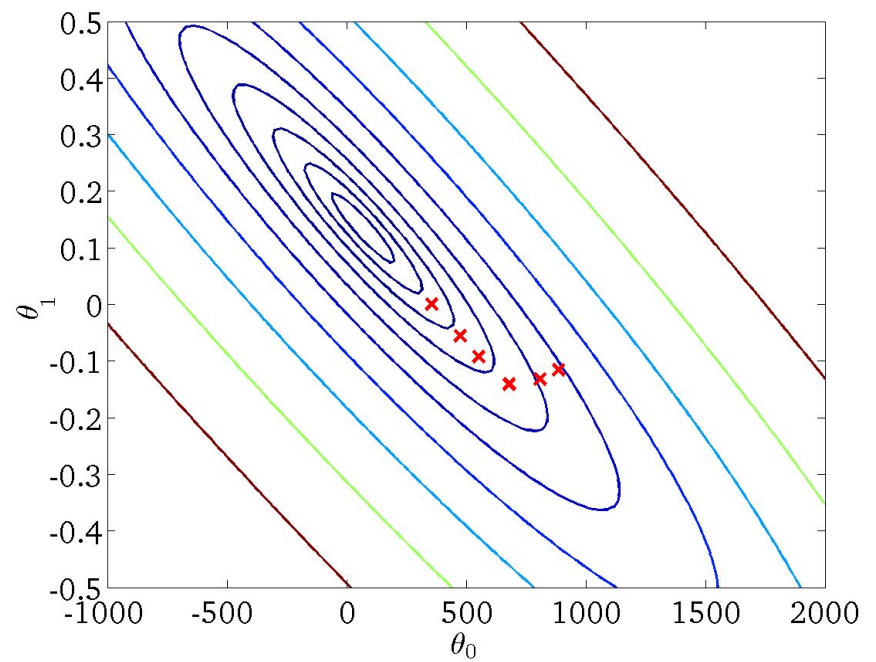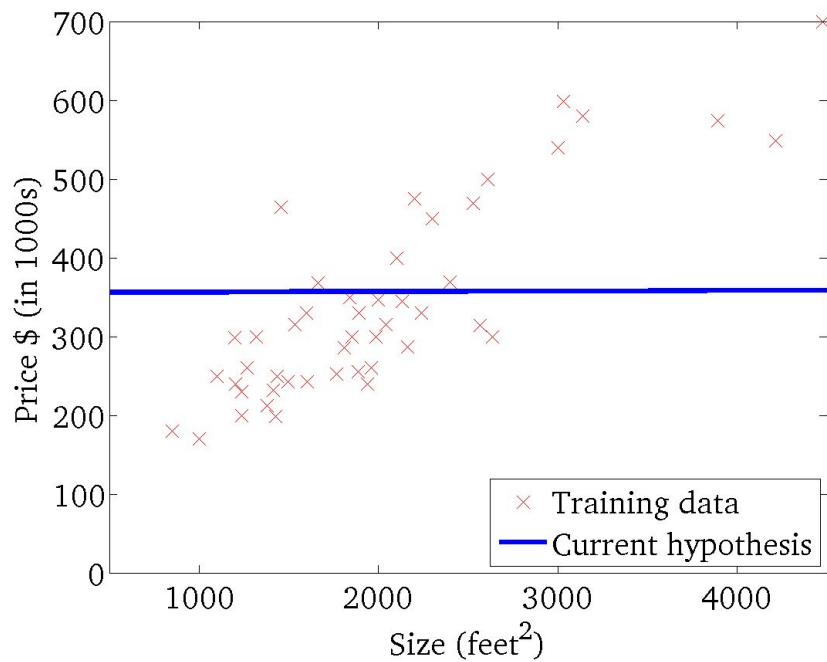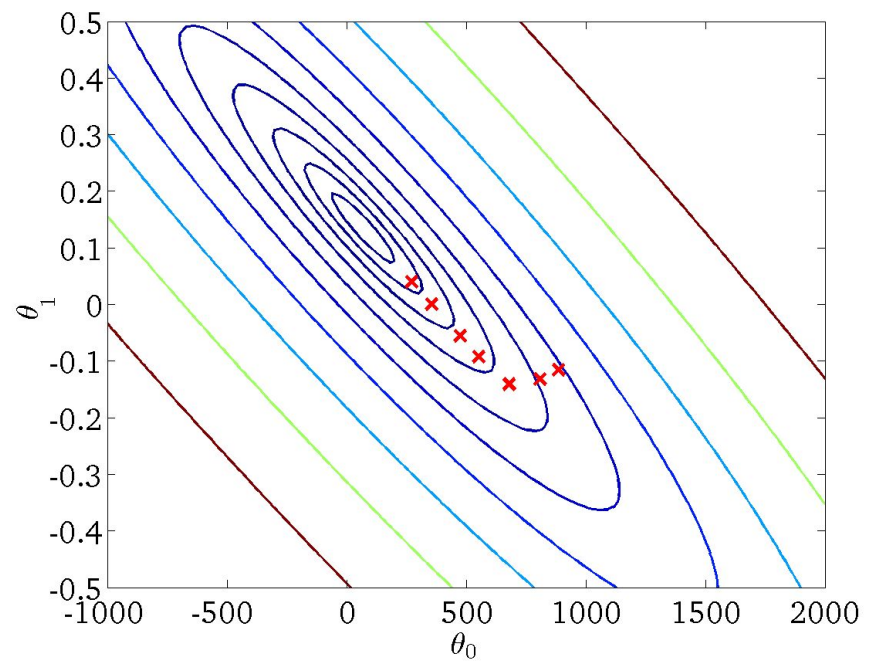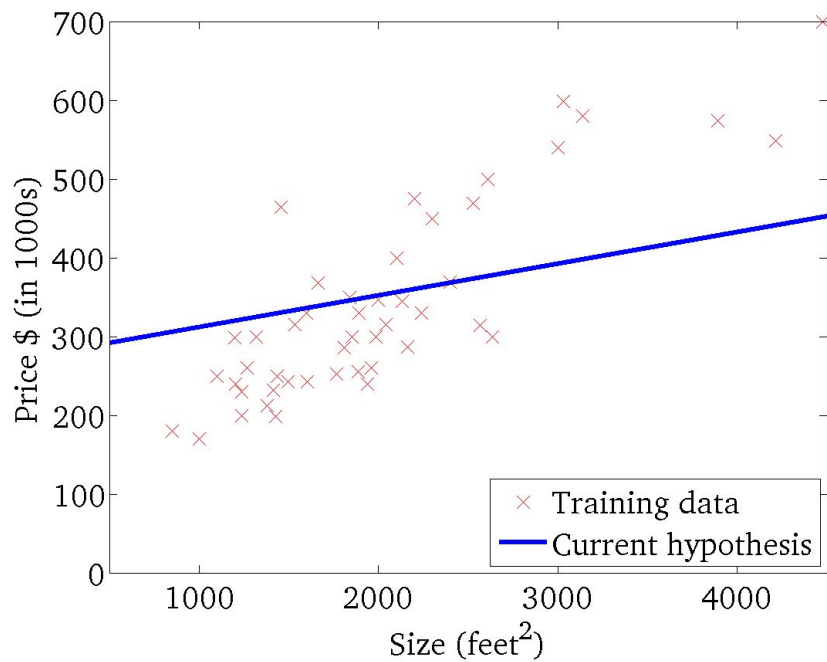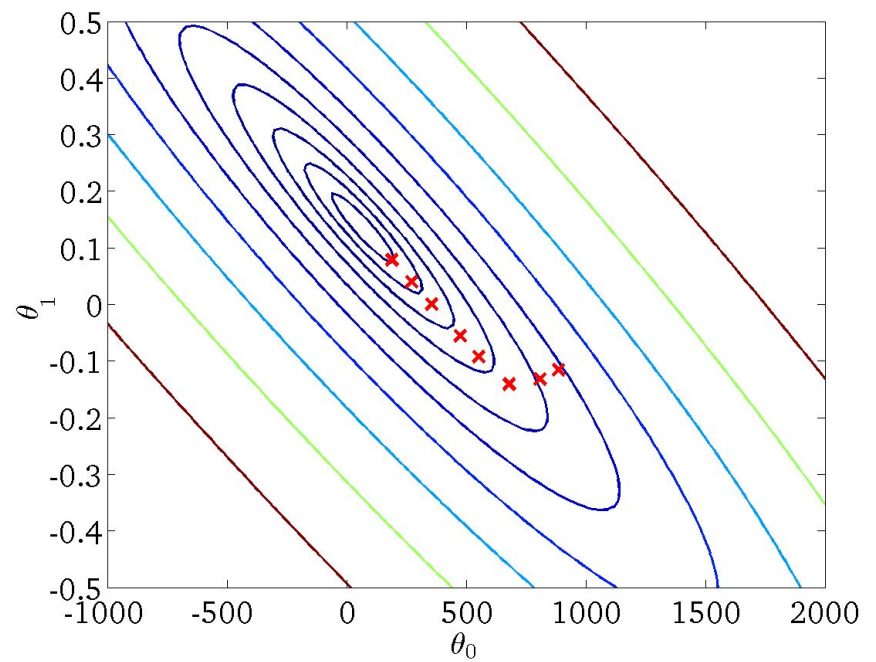
# Gradient Descent in Action

# Gradient Descent in Action

# Gradient Descent in Action

# Gradient Descent in Action
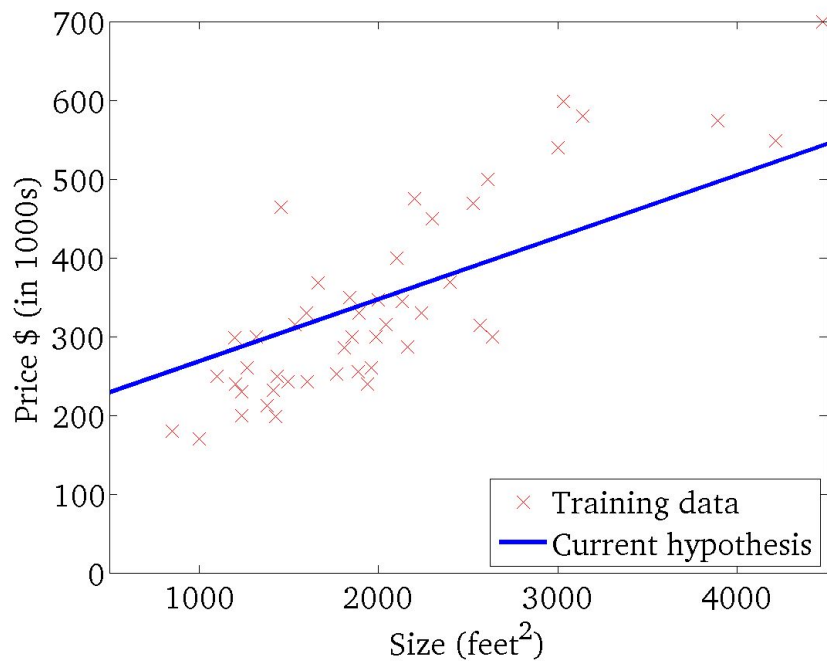
# Gradient Descent in Action
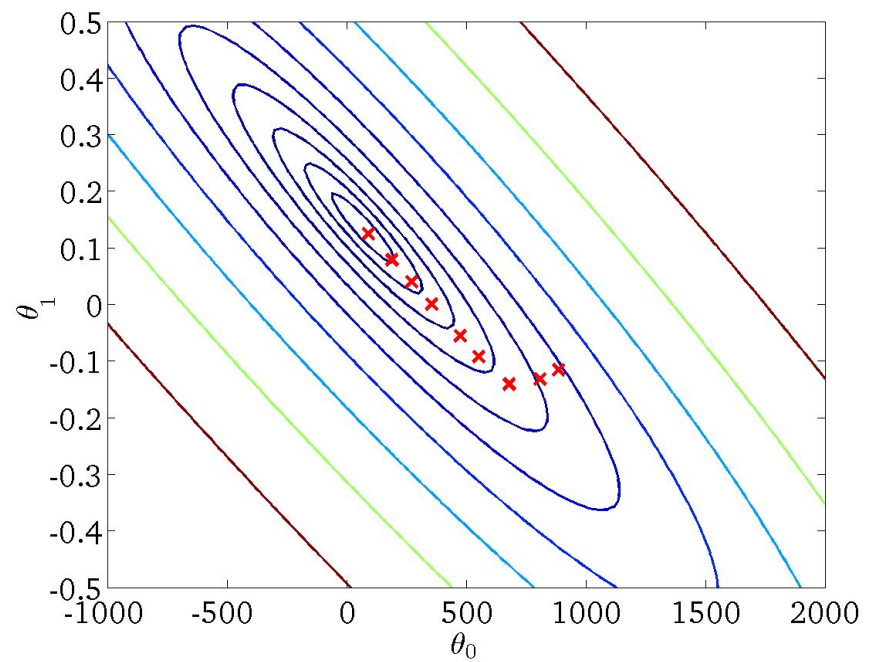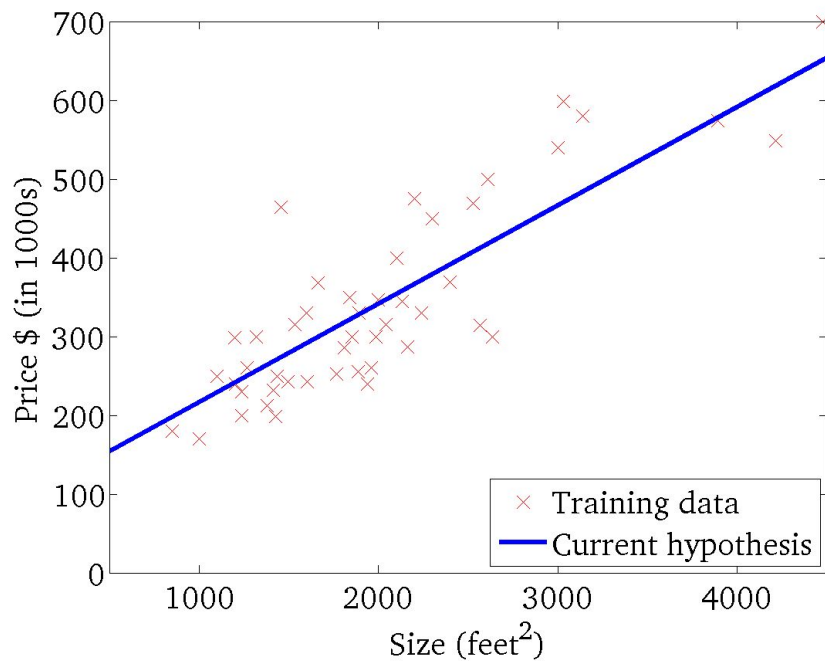
# Gradient Descent in Action

# Gradient Descent in Action

# Gradient Descent in Action

# Gradient Descent in Action

# Summary

- We have learned about:
    - Linear Regression Model with One Variable
    - Cost Function
    - Gradient Descent
    - Gradient Descent with Linear Regression