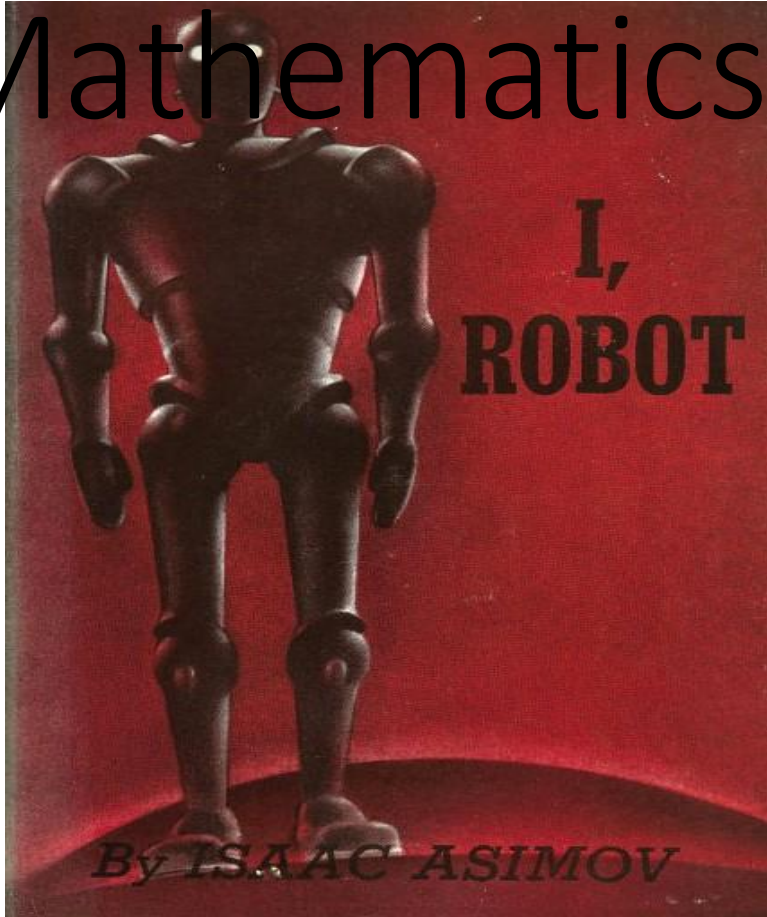# DL in Applied Mathematics

## Linear regression with one variable



Machine Learning

Lecture 3: Introducing Optimization.
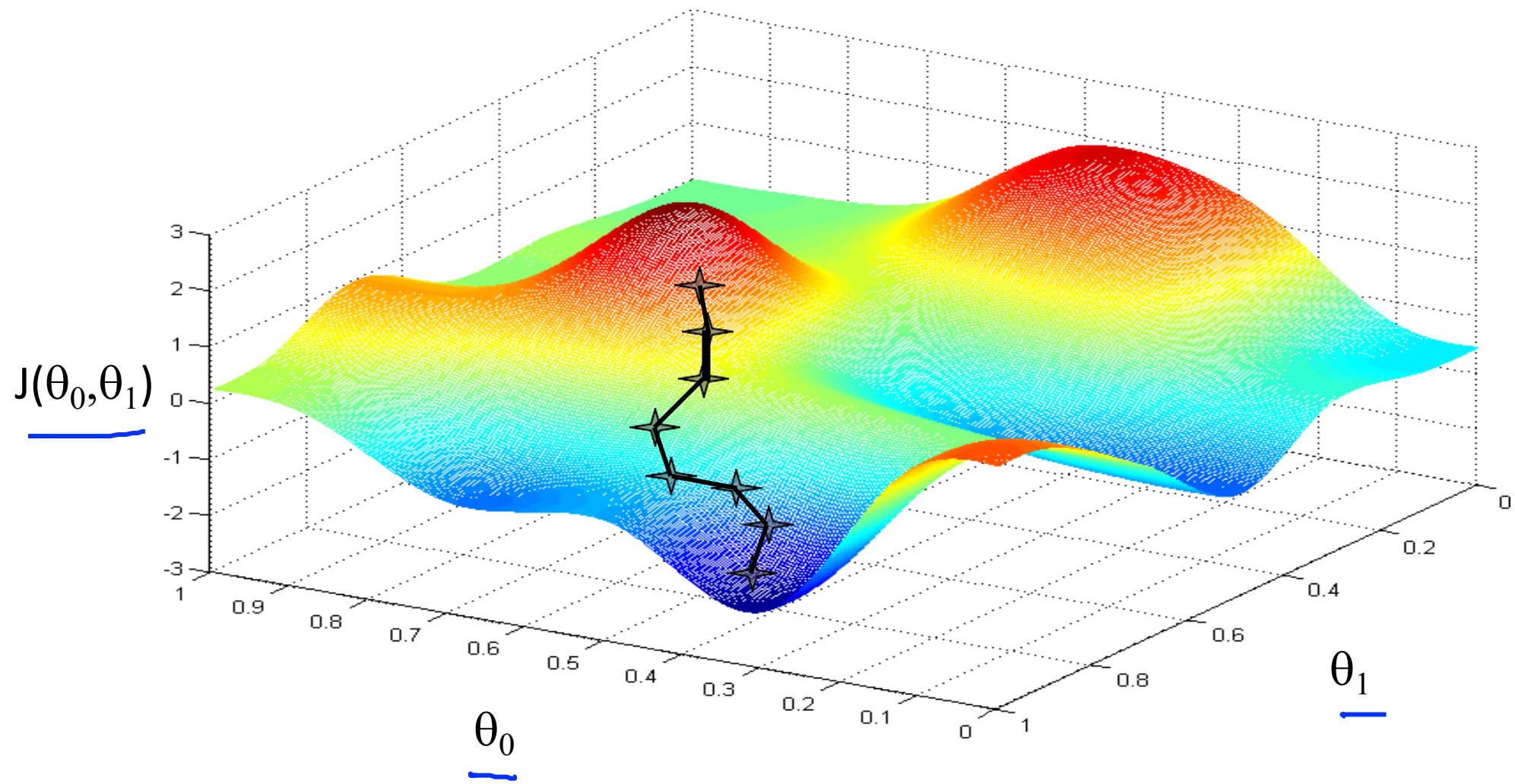
Derivatives in Neural Networks.

Have some function $J(\theta_0, \theta_1)$   $J(\theta_0, \theta_1, \theta_2, \ldots, \theta_n)$
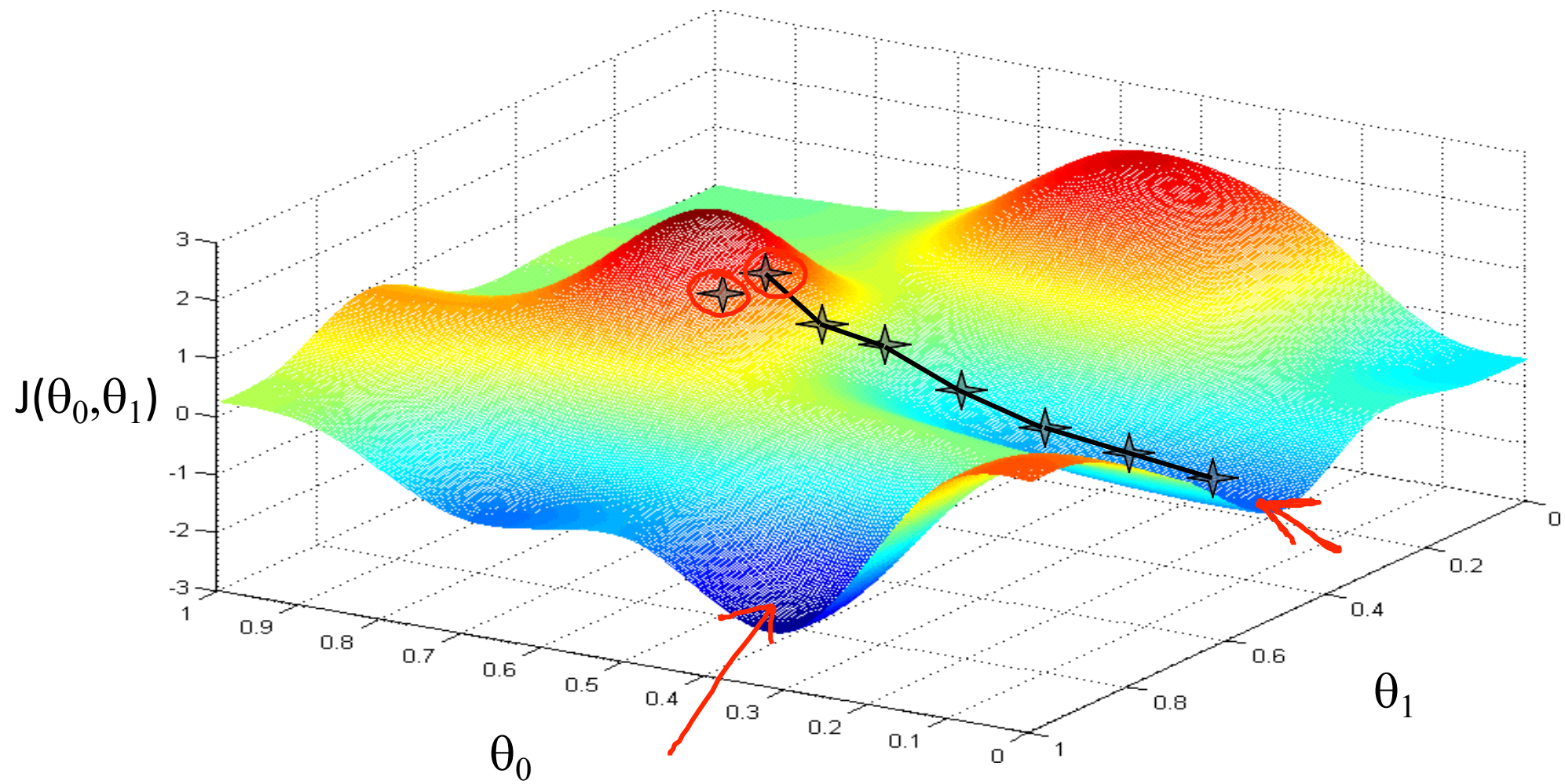
Want $\displaystyle\min_{\theta_0, \theta_1} J(\theta_0, \theta_1)$   $\displaystyle\min_{\theta_0 \ldots \theta_n} J(\theta_0, \ldots, \theta_n)$

**Outline:**

- Start with some $\theta_0, \theta_1$   (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

$J(\theta_0, \theta_1)$

$\theta_0$

$\theta_1$

# Gradient descent algorithm

Assignment

$a := b$

$a := a + 1$

Truth assetion

$a = b$

$a = a + 1$ ✗

repeat until convergence {

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

}

$\theta_0, \theta_1$

learning rate

Simultaneously update $\theta_0$ and $\theta_1$

---

## Correct: Simultaneous update

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

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

$\theta_0 := \text{temp0}$

$\theta_1 := \text{temp1}$
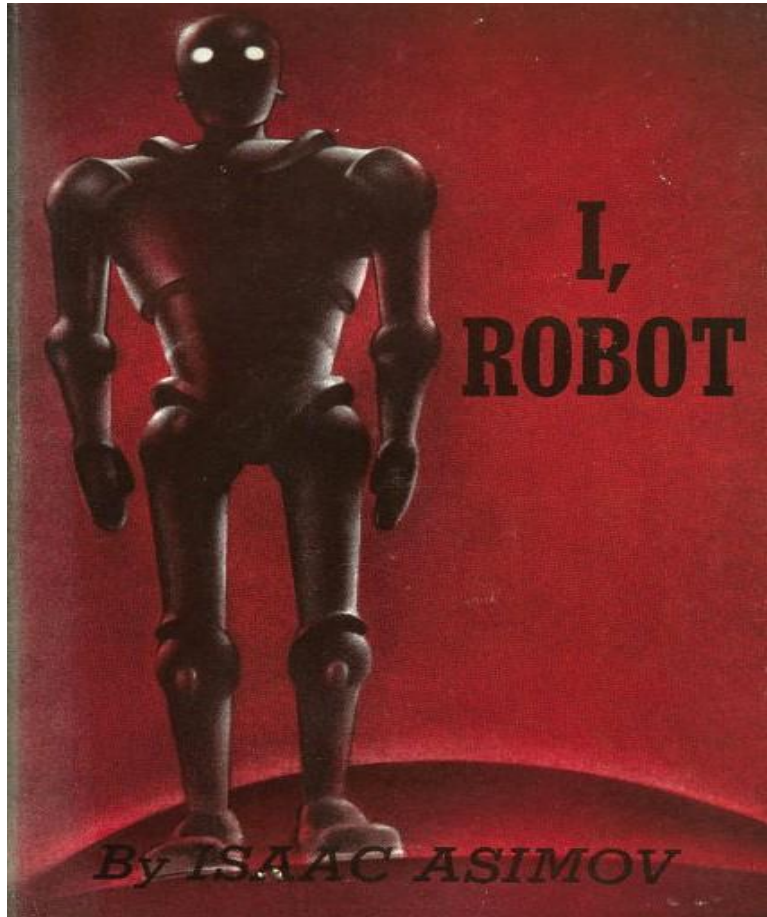
## Incorrect:

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

$\theta_0 := \text{temp0}$

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

$\theta_1 := \text{temp1}$

Machine Learning

Linear regression with one variable

Gradient descent intuition

# What is Gradient Descent?

Gradient Descent is an iterative optimization algorithm for finding the minimum of a function.

In the context of machine learning, this function is usually the cost or loss function, which measures the difference between the model's predictions and the actual data.

The goal is to adjust the model's parameters (e.g., weights in a neural network) to minimize this cost function.

# 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

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

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

$$J(\theta_1)$$

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

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

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

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

$$\leftarrow \theta_1$$

negative slope

$$J(\theta_1)$$

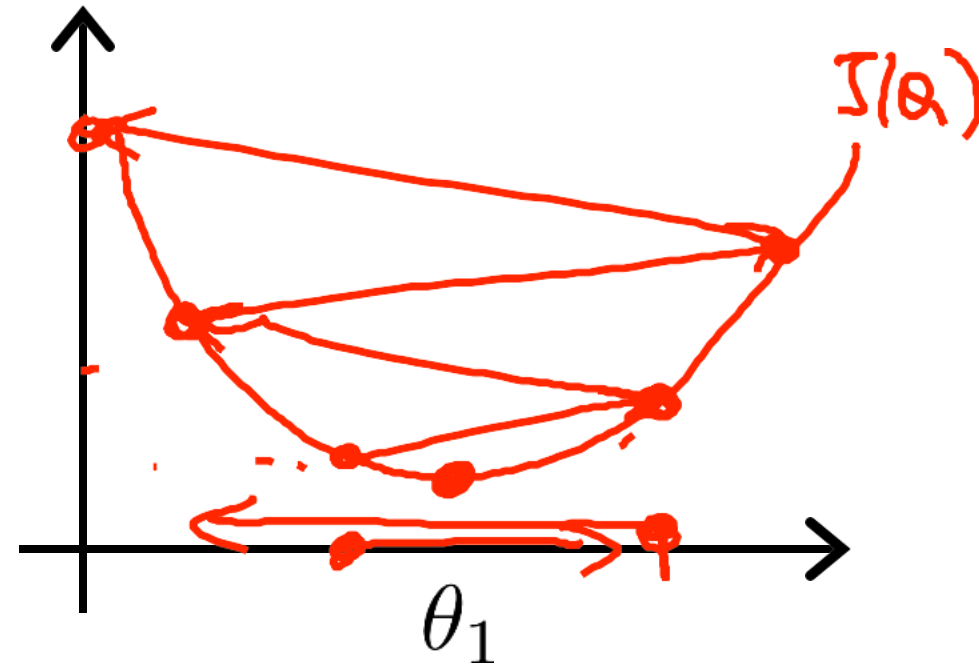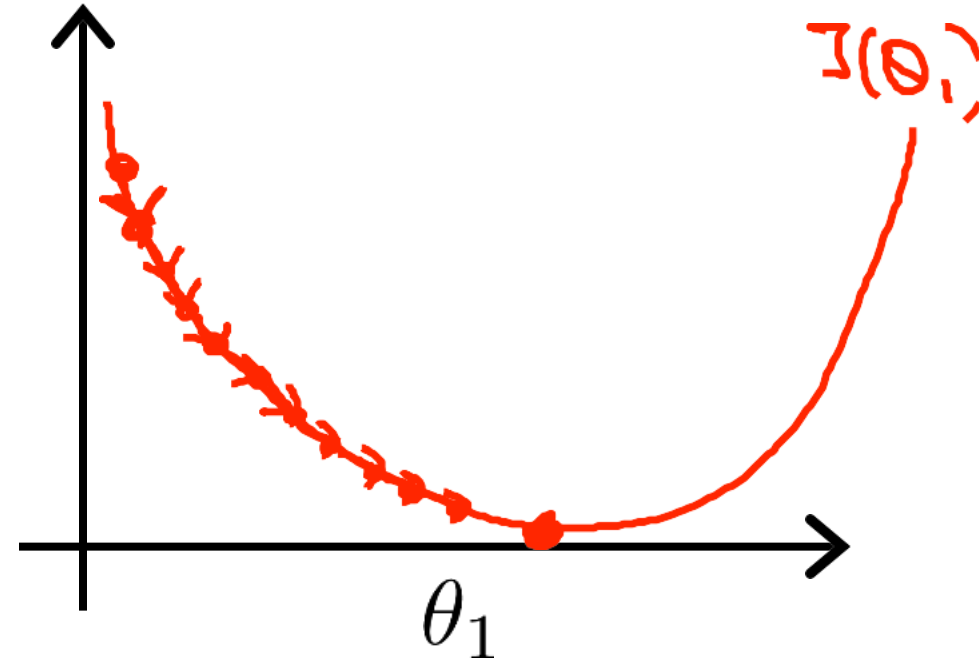$$\theta_1 \rightarrow$$

$$\frac{\frac{d}{d\theta_1} J(\theta_1)}{\leq 0}$$
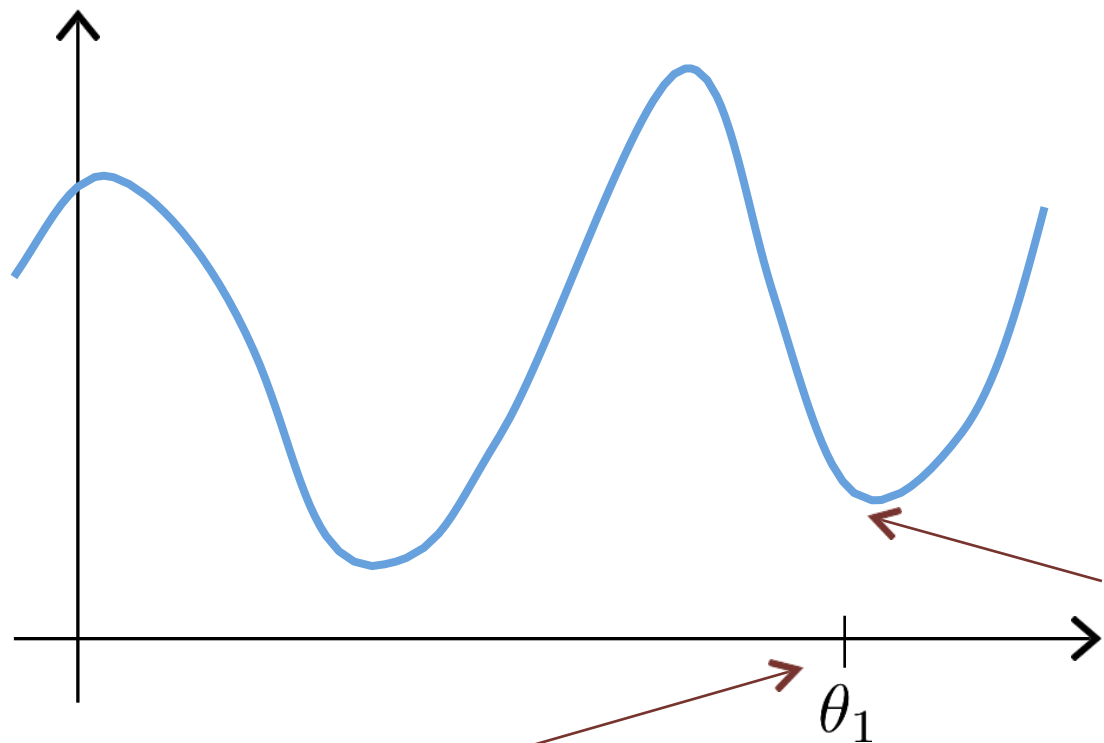
$$\theta_1 := \theta_1 - \alpha (\text{negative number})$$

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

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.
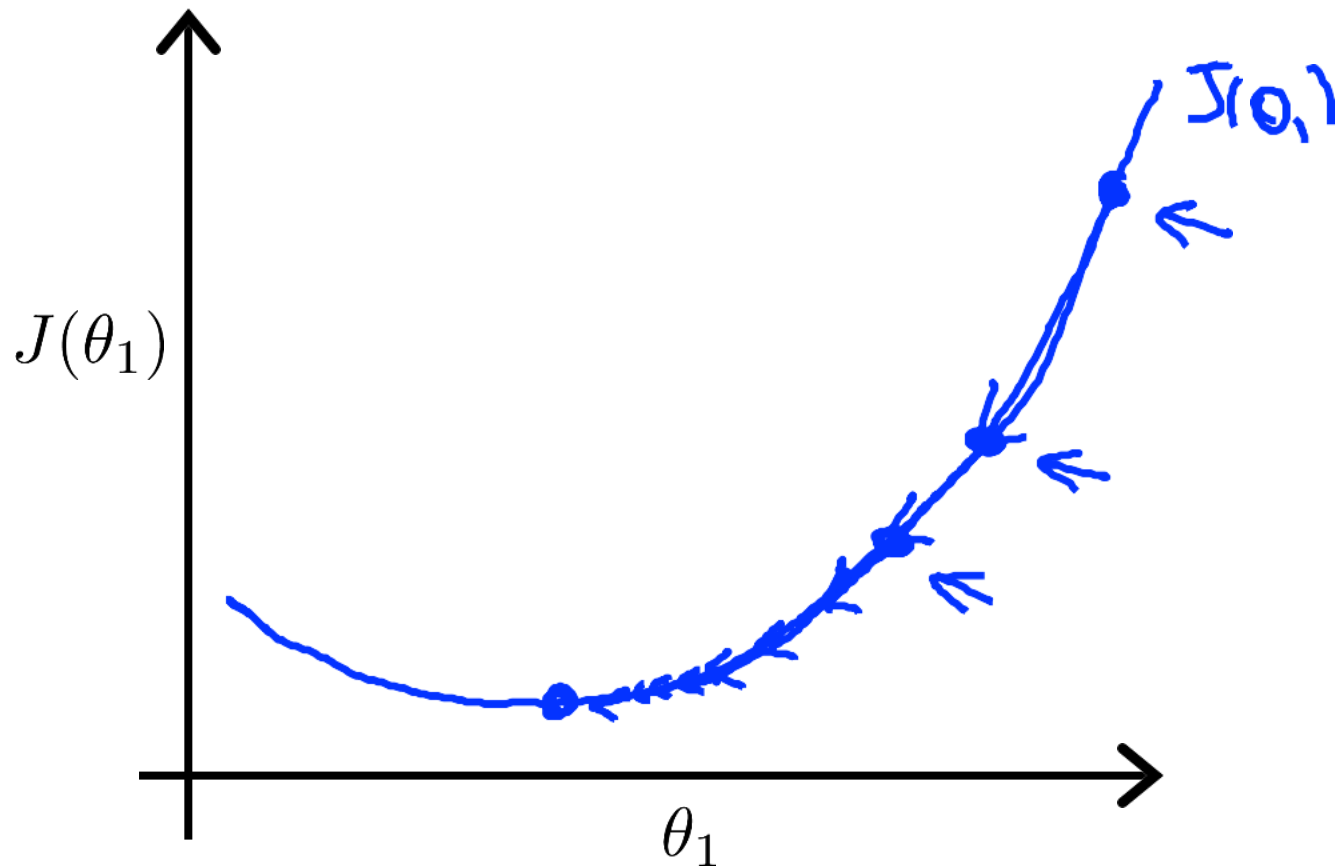
$\theta_1$ at local optima
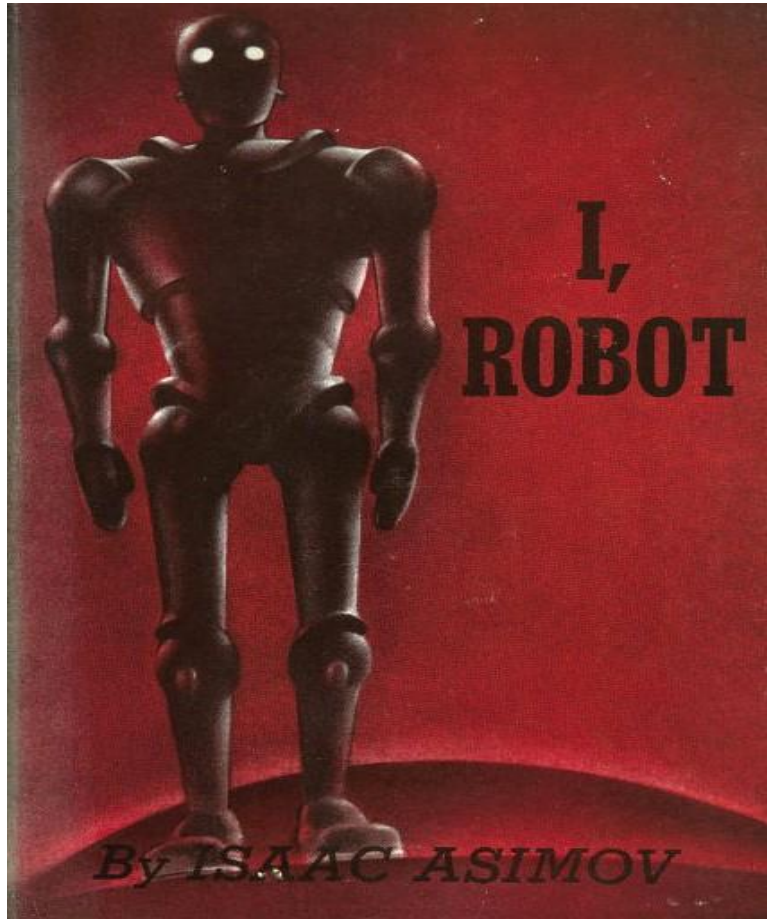
Current value of $\theta_1$

$\theta_1$

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

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.

$J(\theta_1)$

$J(\theta_1)$

$\theta_1$

Machine Learning

Linear regression with one va

Gradient descent for linear regression

## 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

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

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

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

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

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

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

# Gradient descent algorithm

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

repeat until convergence {

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

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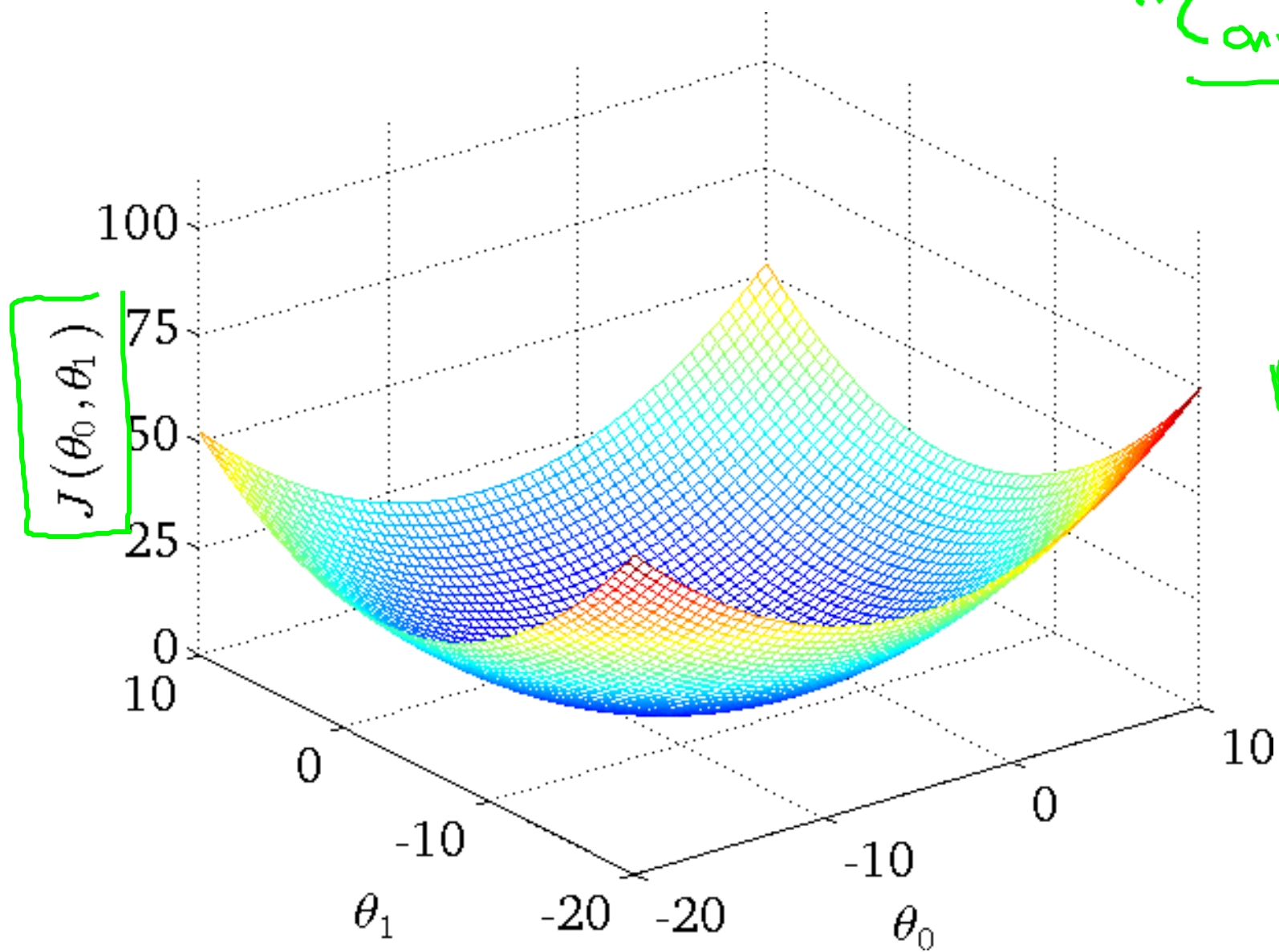
}

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

update $\theta_0$ and $\theta_1$ simultaneously
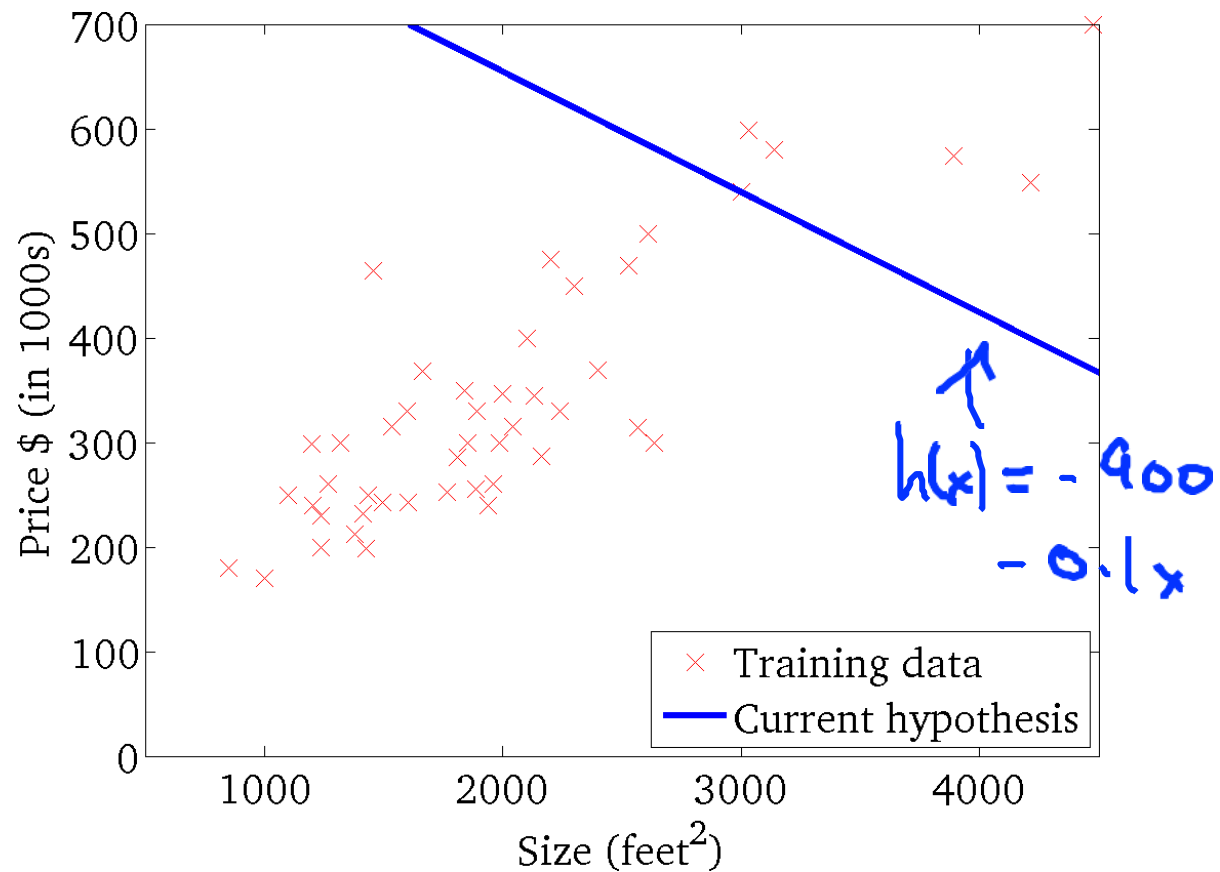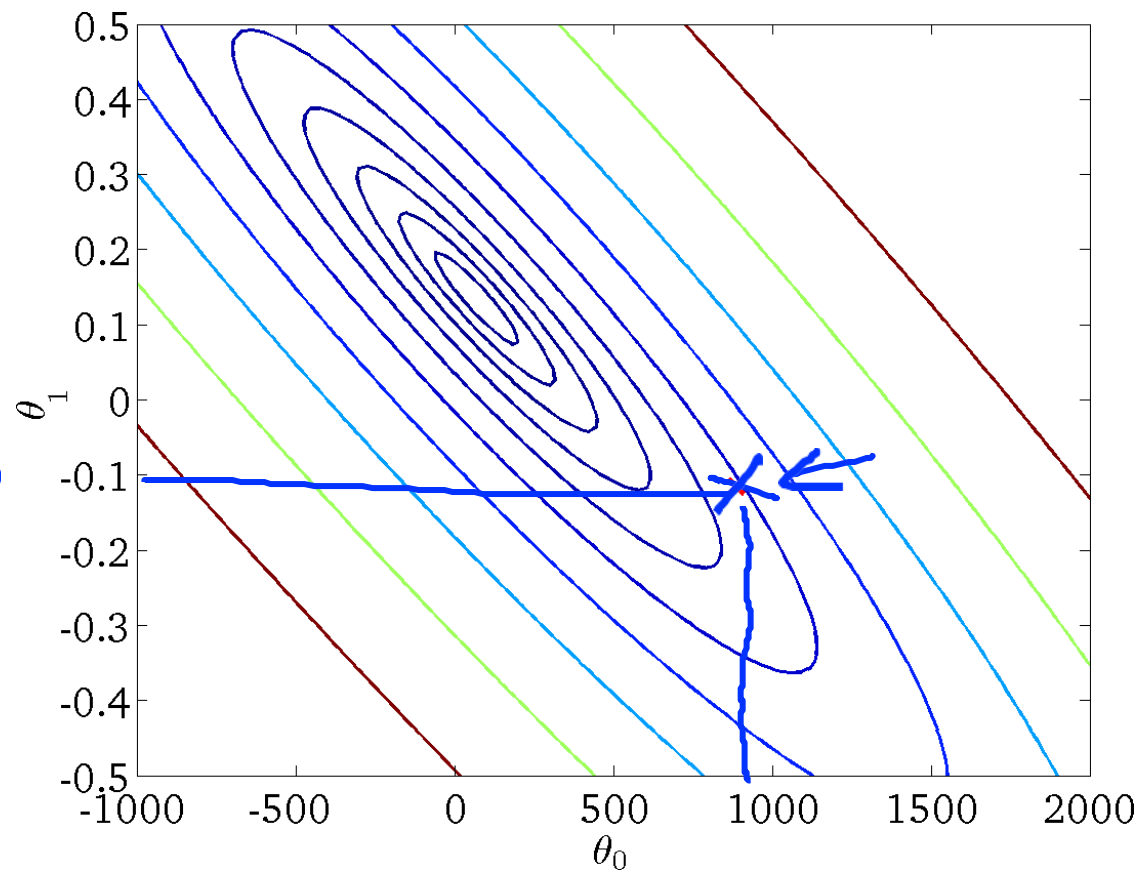
"Convex function"

Bowl-shaped

# $h_\theta(x)$

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

# $J(\theta_0, \theta_1)$

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



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

Legend (left plot):
- × Training data
- ── Current hypothesis

# $h_\theta(x)$

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

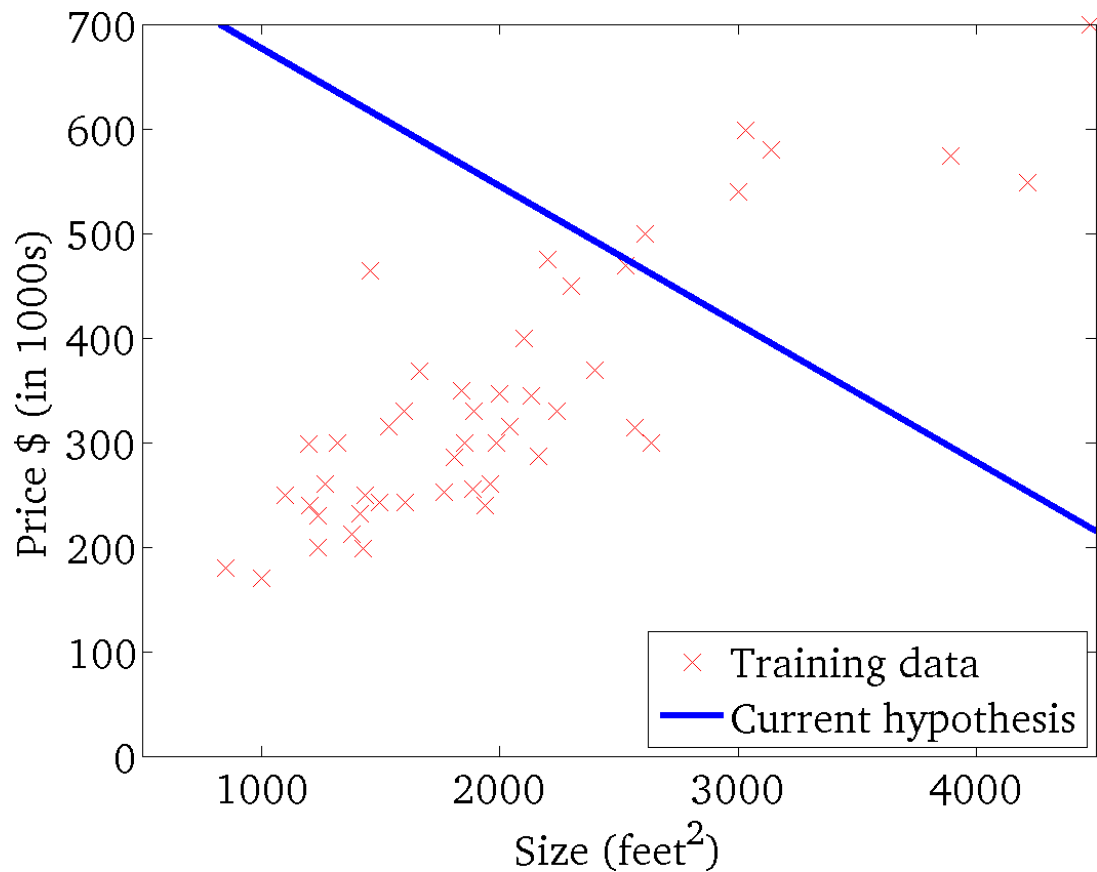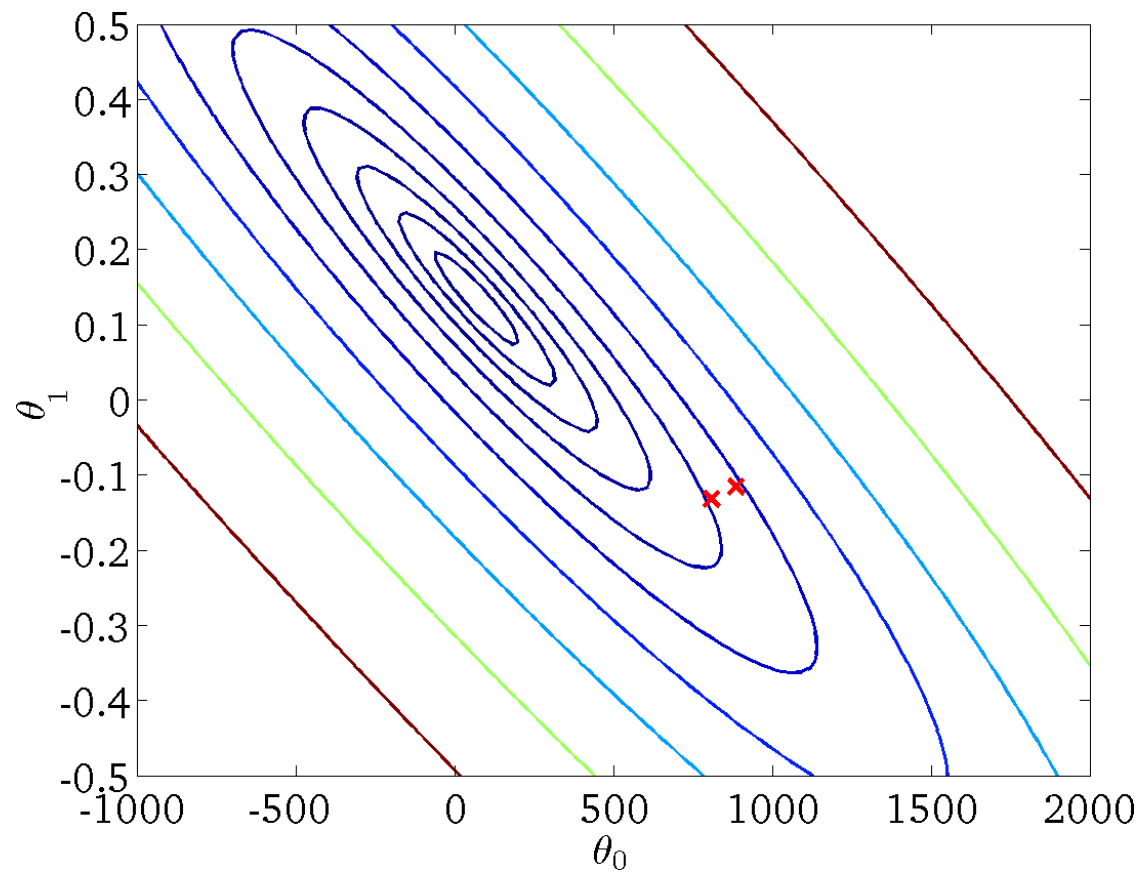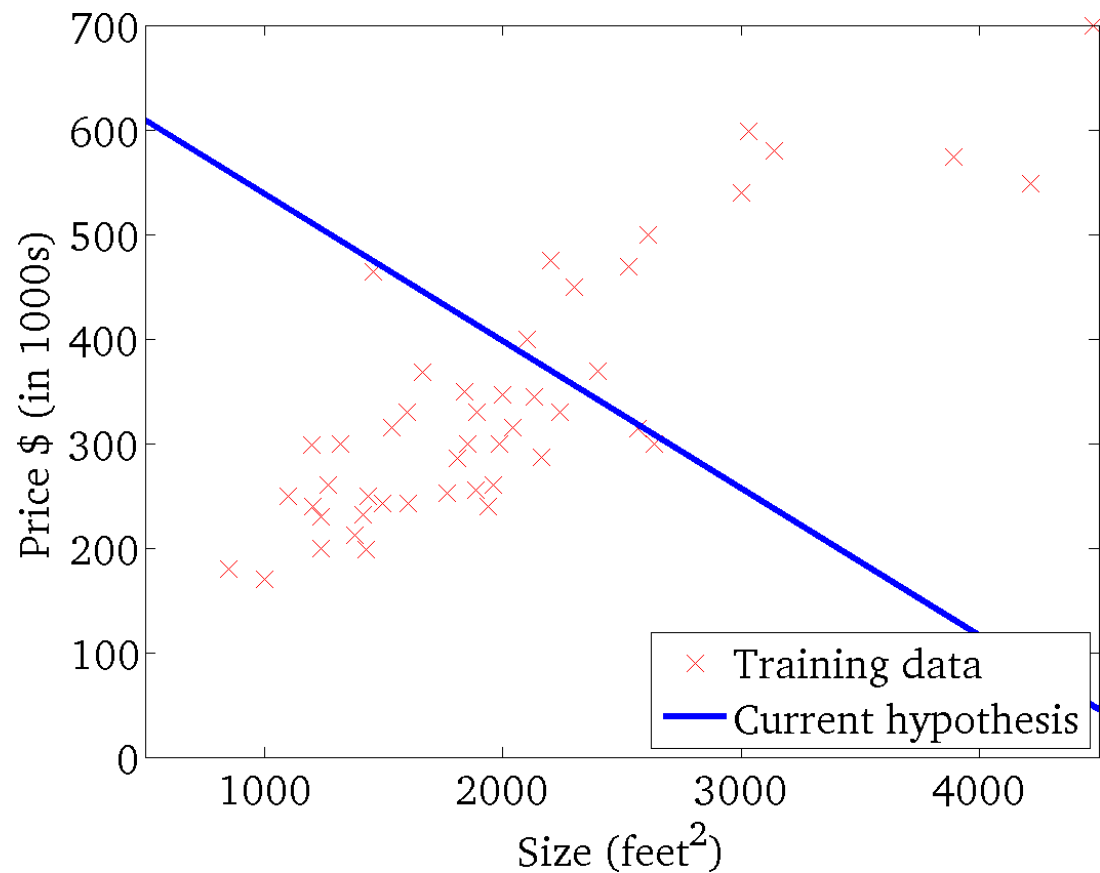# $J(\theta_0, \theta_1)$

(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)

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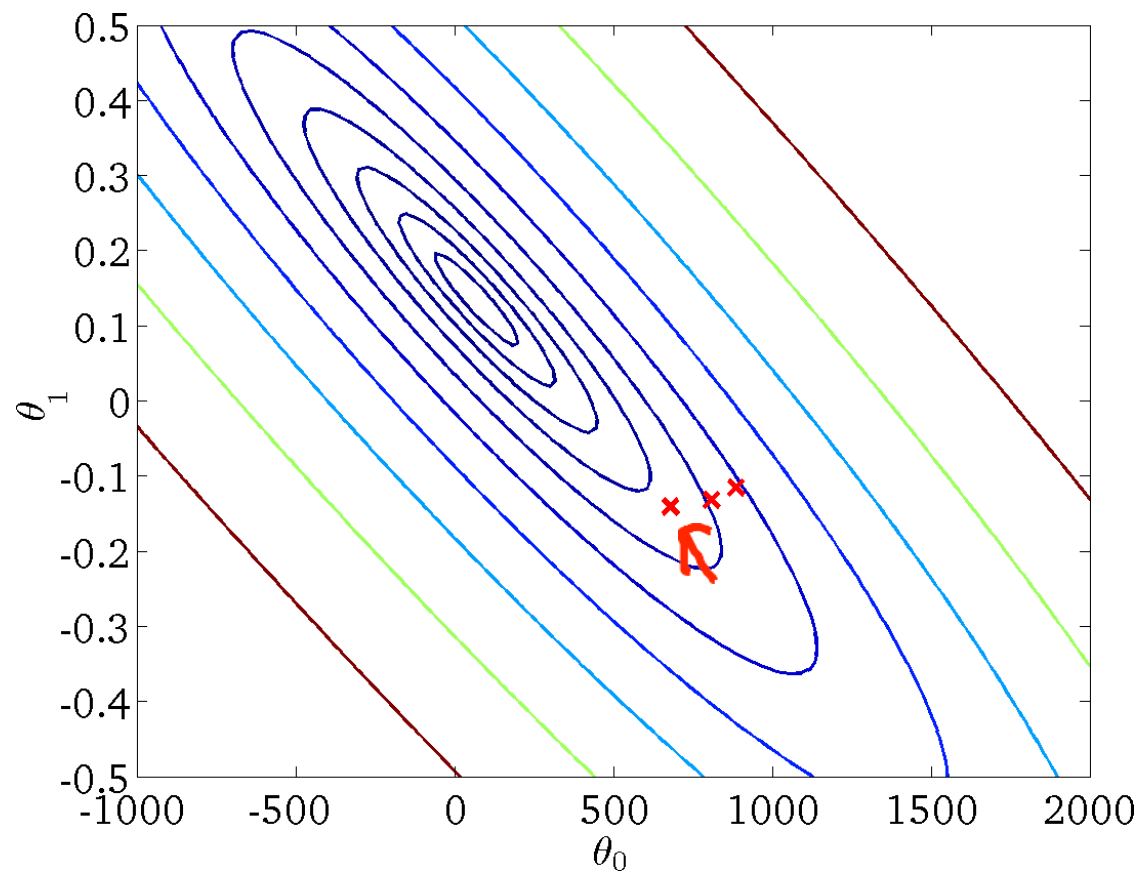
(function of the parameters $\theta_0, \theta_1$)

# $h_\theta(x)$

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

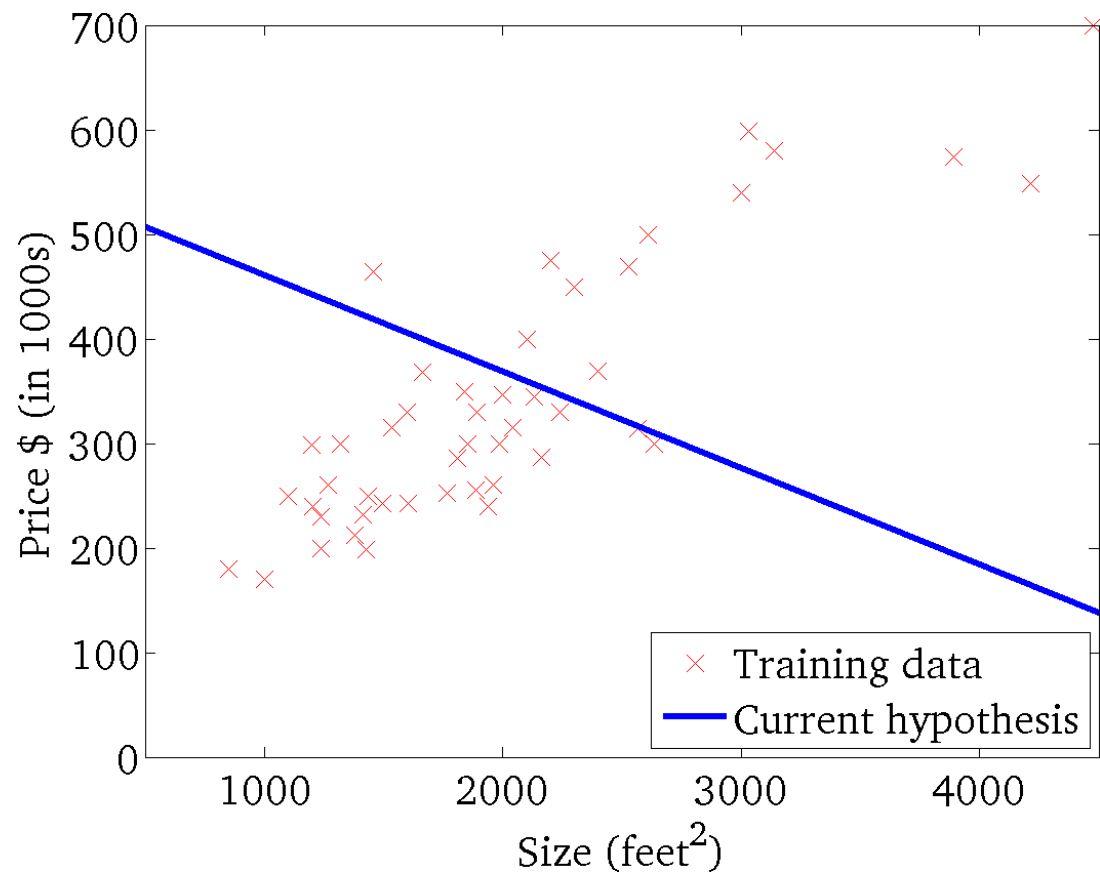# $J(\theta_0, \theta_1)$

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

$h_\theta(x)$

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

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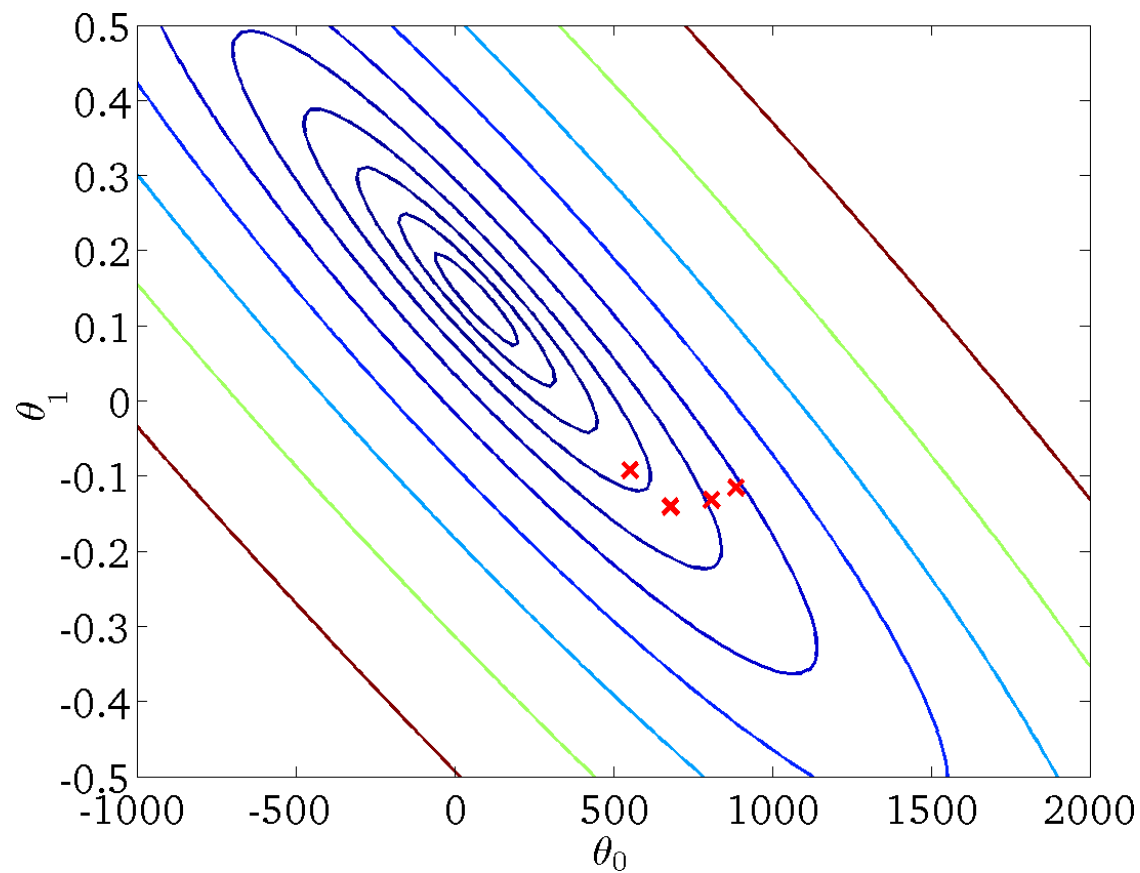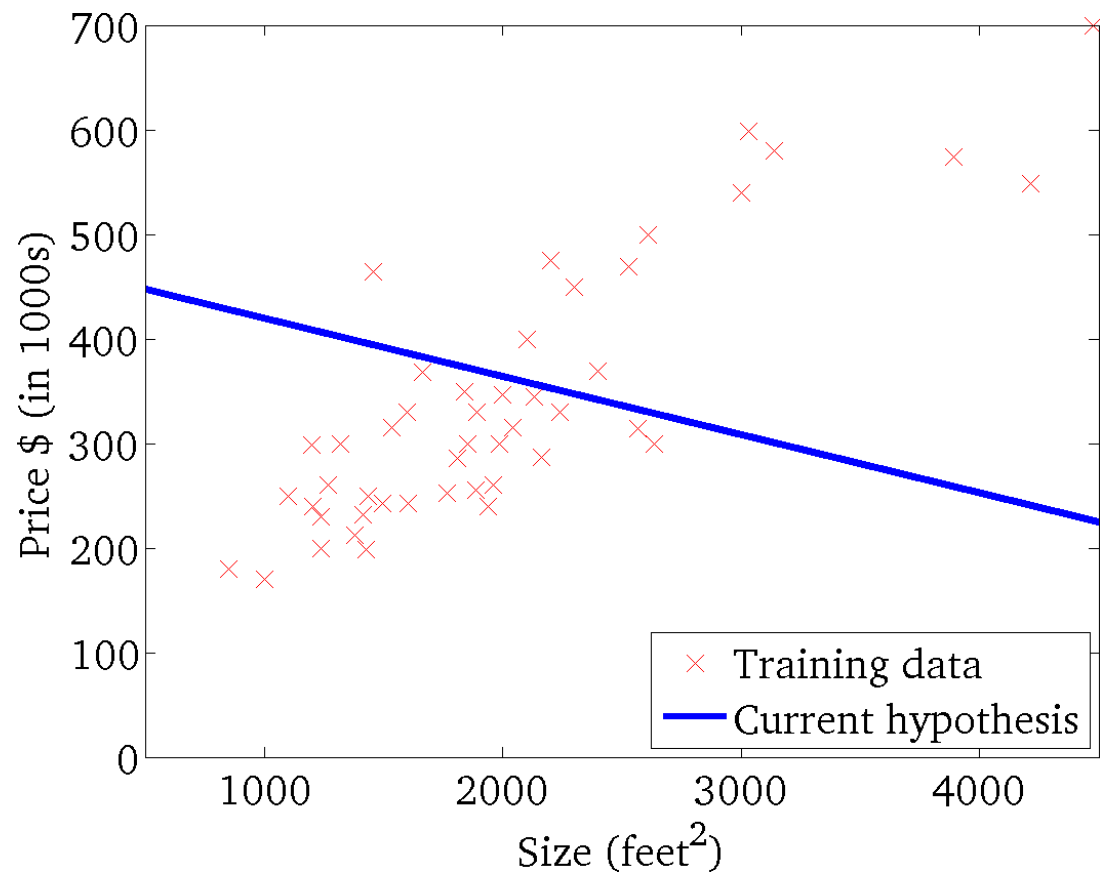
(function of the parameters $\theta_0, \theta_1$)

# $h_\theta(x)$

## (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 $\theta_0, \theta_1$)

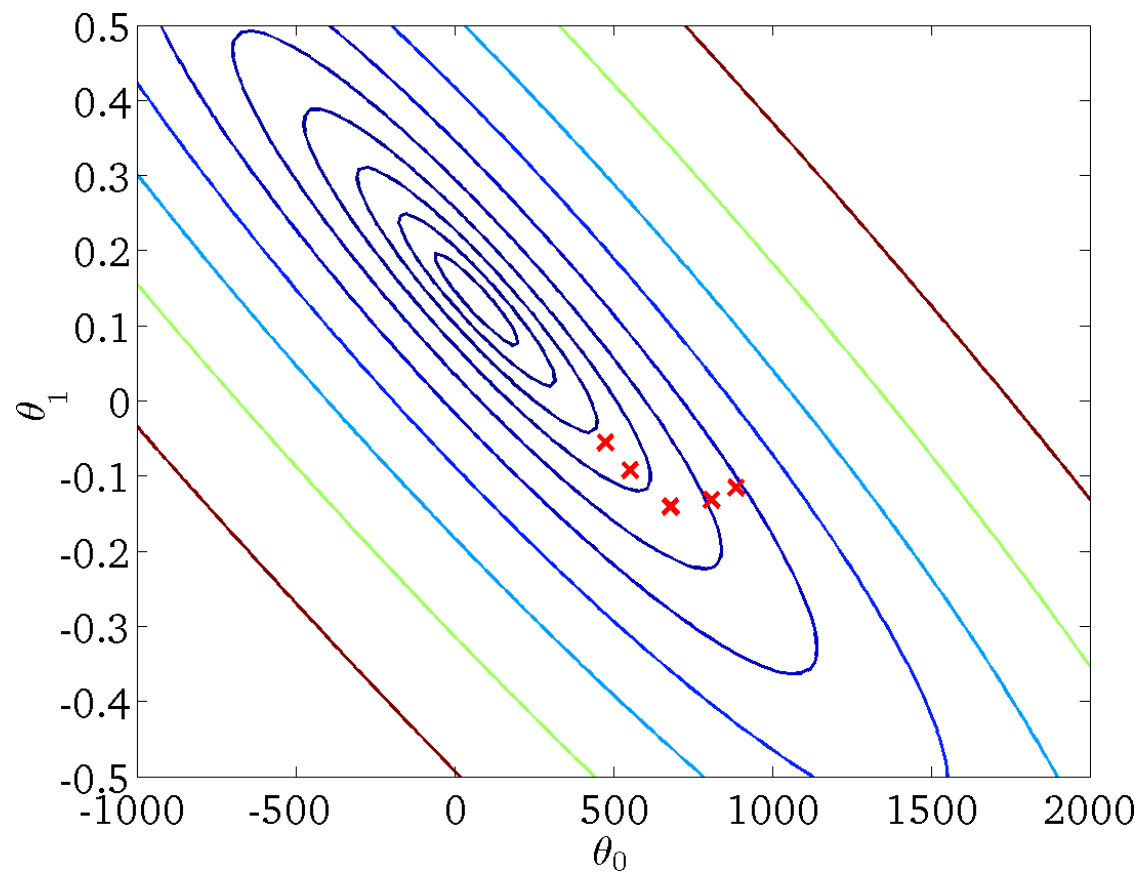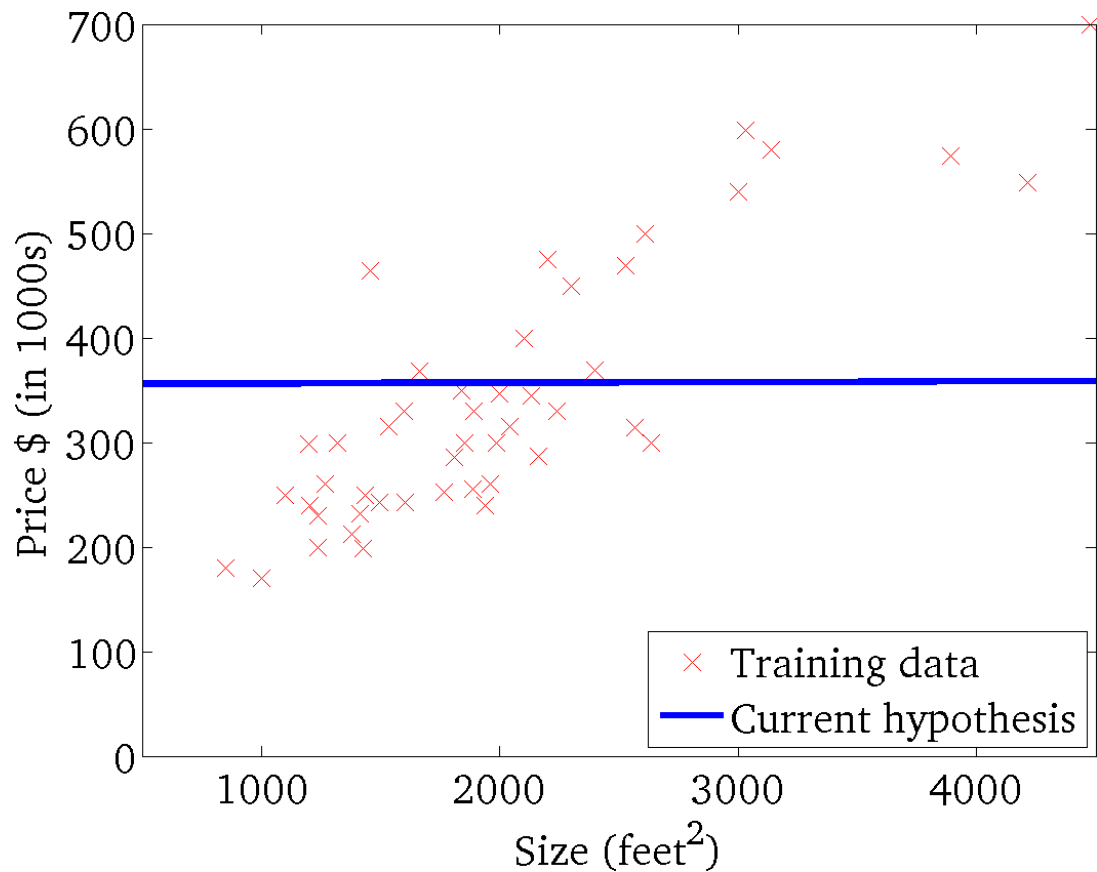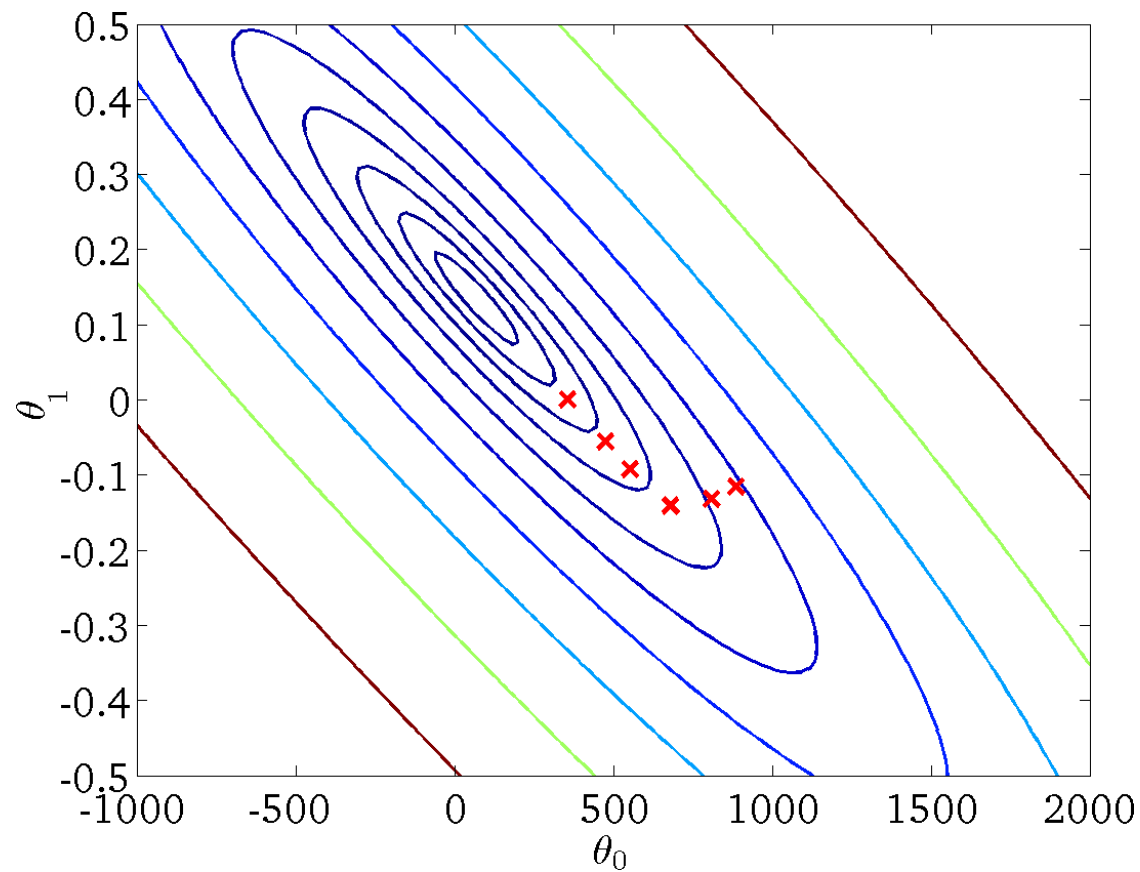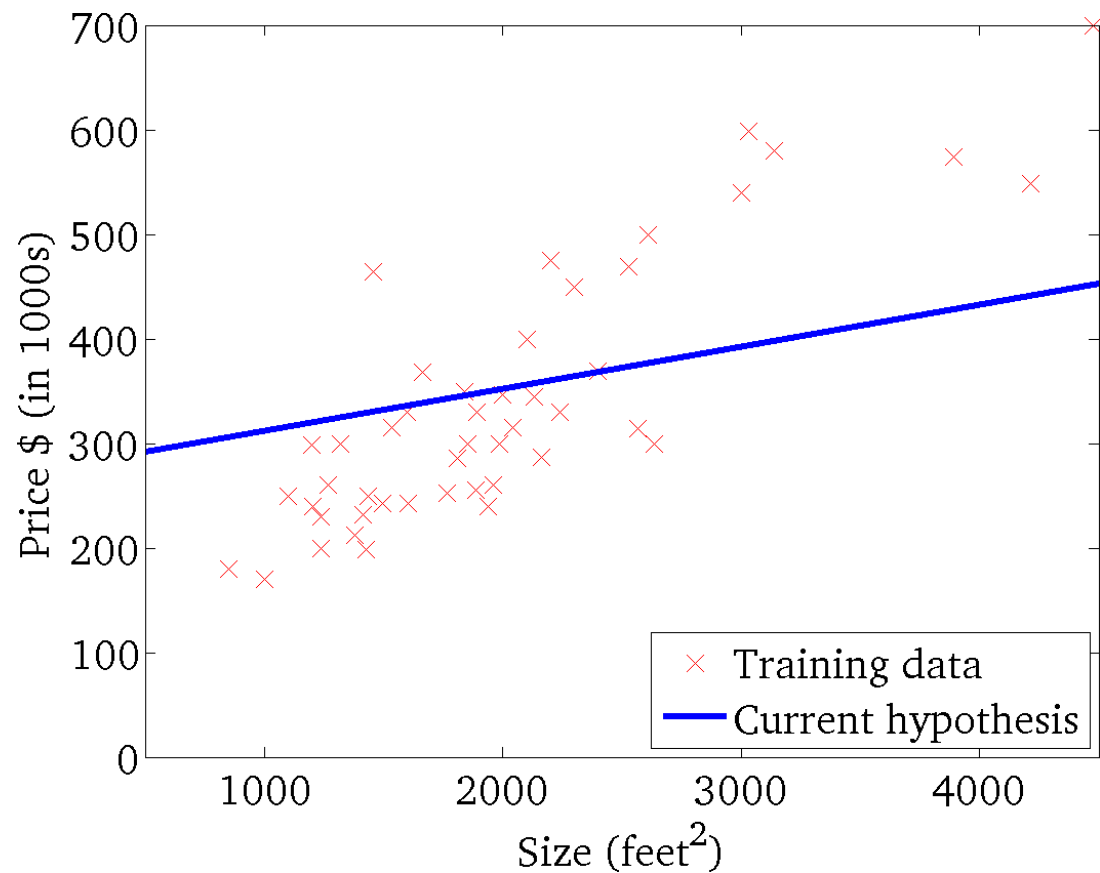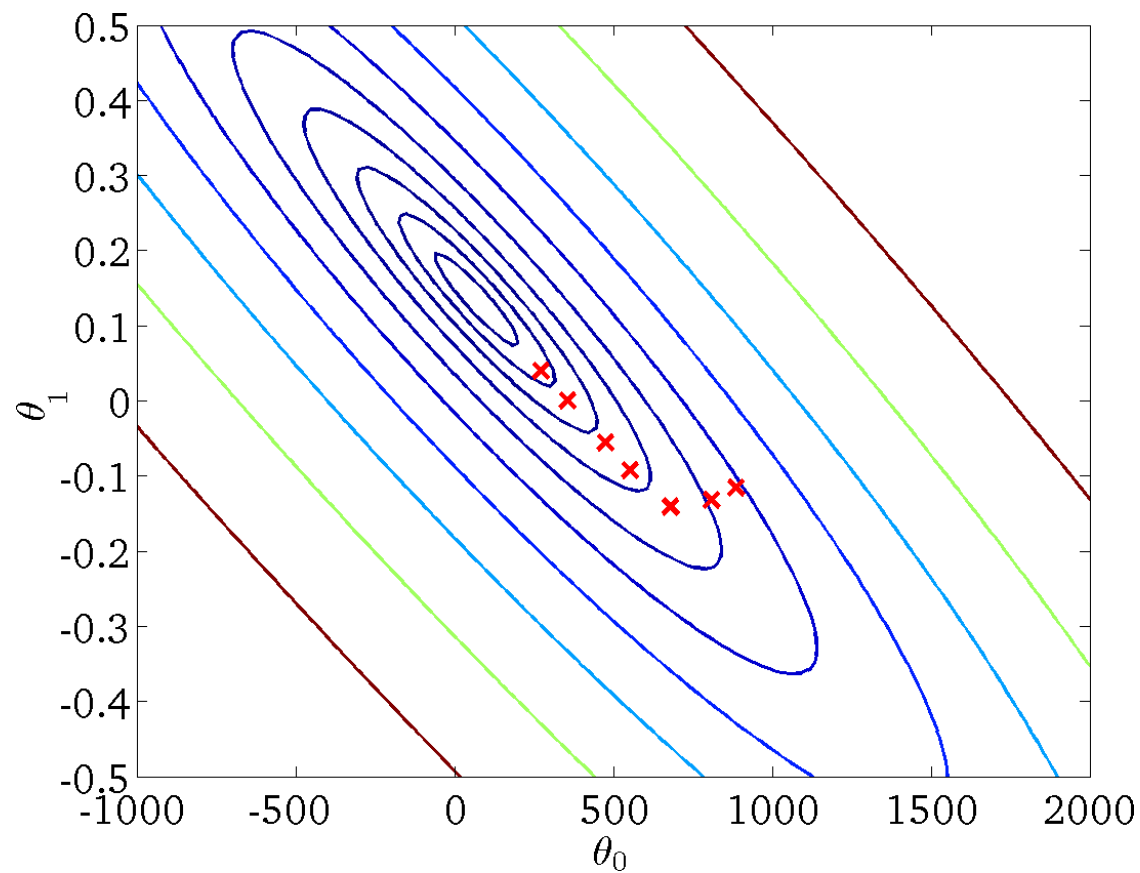$$h_\theta(x)$$

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

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

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

# "Batch" Gradient Descent

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

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

"Batch" Gradient Descent is a variation of the Gradient Descent optimization algorithm, which is used to minimize a cost function. It's a fundamental algorithm in machine learning for training a wide range of models, including linear regression, logistic regression, and neural networks.

# How Does "Batch" Gradient Descent Work?

- **Batch Processing:** In Batch Gradient Descent, the term "batch" refers to the use of the entire training dataset to calculate the gradient of the cost function with respect to the model parameters. This means that for each iteration of the algorithm, the model uses all available data points to perform an update.

- **Gradient Calculation:** The gradient—a vector of partial derivatives—indicates the direction in which the cost function is increasing most rapidly. By moving in the opposite direction (i.e., descending), we aim to find the function's minimum.

- **Parameter Update:** The model parameters are updated in the direction of the negative gradient scaled by the learning rate, a hyperparameter that controls the size of the step. The equation for updating each parameter $\vartheta$ is given by: $\theta = \theta - \eta \cdot \nabla_\theta j(\theta)$ where **$\eta$** is the learning rate and $\nabla_\theta j(\theta)$ is the gradient of the cost function **J** with respect to $\theta$.

# Advantages of Batch Gradient Descent

- **Simplicity:** It's straightforward to implement and understand.

- **Stable Convergence:** Using the entire dataset provides a clear direction toward the minimum, leading to stable convergence in convex or quadratic problems.

- **Efficiency on Small Datasets:** When the dataset is small enough to fit into memory, batch processing can be computationally efficient due to optimized matrix operations.

**Disadvantages of Batch Gradient Descent**

- **Scalability:** For very large datasets, loading the entire dataset into memory can be impractical, leading to computational inefficiencies.

- **Speed:** Each iteration can be slow because it must process the entire dataset before making a single update to the parameters.

- **Convergence to Global Minimum:** In non-convex functions, such as those common in deep learning, Batch Gradient Descent can get stuck in local minima.

**Alternatives to Batch Gradient Descent**

- **Stochastic Gradient Descent (SGD):** Updates parameters for each training example, leading to faster iterations but more variance in the convergence path.

- **Mini-batch Gradient Descent:** Strikes a balance between Batch and SGD by updating parameters after computing the gradient on small subsets of the data. This approach is often used in practice due to its balance between efficiency and stability.

- In summary, Batch Gradient Descent is a powerful but somewhat simplistic tool in the machine learning optimization arsenal. Its use is best suited to situations where datasets are manageable in size and the cost function is relatively simple. For larger datasets or more complex cost functions, its alternatives, such as Mini-batch Gradient Descent, are generally preferred.