# Machine Learning
# 机器学习

## Lecture2:线性回归

李洁

nijanice@163.com

# 学习任务的类型
## Types of learning task

- Supervised learning
  - infer a function from labeled training data.

- Unsupervised learning
  - try to find hidden structure in unlabeled training data
  - clustering

- Reinforcement learning
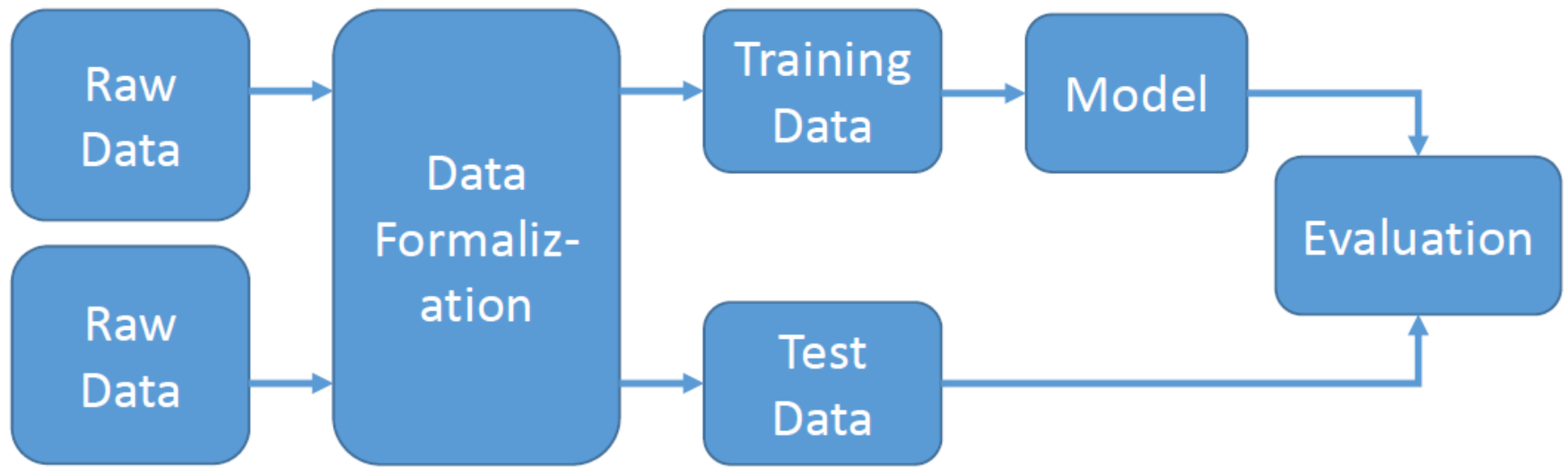  - To learn a policy of taking actions in a dynamic environment and acquire rewards

# 学习任务的类型
## Types of learning task

|  | **Supervised Learning** | **Unsupervised Learning** |
|---|---|---|
| **Discrete** | classification or categorization | clustering |
| **Continuous** | regression | dimensionality reduction |

# 机器学习的一般过程
# Machine Learning Process



- Basic assumption: there exist the same patterns across training and test data

# 监督学习
# Supervised Learning

- Given the training dataset of (data,label) pairs,

$$D = \left\{ \left( x^{(i)}, y^{(i)} \right) \right\}_{i=1.2,\ldots,N}$$

$$x^{(i)} = (x_1^{(i)}, x_2^{(i)}, \ldots x_n^{(i)})^T$$

$y^{(i)} =$ output data(label) of $i^{th}$ training example

 let the machine learn a function from data to label

$$y^{(i)} \approx f_\theta \left( x^{(i)} \right)$$

- Function set $\{ f_\theta (x^{(i)}) \}$ is called hypothesis space
- Learning is referred to as updating the parameter $\theta$ to make the prediction closed to the corresponding label

# 线性模型
# Linear Model

easily understood and implemented, efficient and scalable

- – Linear regression
- – Linear classification

# 线性模型举例
# Linear model example

$$f_{好瓜}(\boldsymbol{x}) = 0.2 \cdot x_{色泽} + 0.5 \cdot x_{根蒂} + 0.3 \cdot x_{敲声} + 1$$



周志华. "机器学习" (西瓜书)

# 线性模型举例
# Linear model example

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

$$f_{好瓜}(\boldsymbol{x}) = 0.2 \cdot x_{色泽} + 0.5 \cdot x_{根蒂} + 0.3 \cdot x_{敲声} + 1$$



周志华．"机器学习"　（西瓜书）

# 线性回归模型
# Linear Regression Model

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

sample $x$

features/variables: $x_1, x_2, \ldots x_n$

# 线性回归模型
# Linear Regression Model

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

sample $x$

features/variables: $x_1, x_2, \ldots x_n$

$x = (x_1, x_2, \ldots x_n)^T$

Feature vector $(x_1, x_2, \ldots x_n)^T$

# 监督学习
# Supervised Learning

- Given the training dataset of (data,label) pairs,

$$D = \left\{ \left( x^{(i)}, y^{(i)} \right) \right\}_{i=1.2,\ldots,N}$$

$x^{(i)}$ = input data(features) of $i^{th}$ training example
$y^{(i)}$ = output data(label) of $i^{th}$ training example

 let the machine learn a function from data to label

$$y^{(i)} \approx f_\theta \left( x^{(i)} \right)$$

- Function set $\{ f_\theta (x^{(i)}) \}$ is called hypothesis space
- Learning is referred to as updating the parameter $\theta$ to make the prediction closed to the corresponding label

# 监督学习
# Supervised Learning

- Given the training dataset of (data,label) pairs,

$$D = \{(x^{(i)}, y^{(i)})\}_{i=1.2,...,N}$$

$$x^{(i)} = (x_1^{(i)}, x_2^{(i)}, ... x_n^{(i)})^T$$

$$y^{(i)} = \text{output data(label) of } i^{th} \text{ training example}$$

let the machine learn a function from data to label

$$y^{(i)} \approx f_\theta(x^{(i)})$$

- Function set $\{f_\theta(x^{(i)})\}$ is called hypothesis space
- Learning is referred to as updating the parameter $\theta$ to make the prediction closed to the corresponding label

# 线性回归模型
## Linear Regression Model

- Given the training dataset of (data,label) pairs,

$$D = \left\{\left(x^{(i)}, y^{(i)}\right)\right\}_{i=1.2,\dots,N}$$

$$x^{(i)} = (x_1^{(i)}, x_2^{(i)}, \dots x_n^{(i)})^T$$

$y^{(i)} =$ output data(label) of $i^{th}$ training example

let the machine learn a function from data to label

$$y^{(i)} \approx f_\theta\left(x^{(i)}\right) \implies f_\theta(x^{(i)}) = \theta_1 x_1^{(i)} + \theta_2 x_2^{(i)} + \cdots + \theta_n x_n^{(i)} + \theta_0$$

- Function set $\{f_\theta(x^{(i)})\}$ is called hypothesis space
- Learning is referred to as updating the parameter $\theta$ to make the prediction closed to the corresponding label

# 线性回归模型
# Linear Regression Model

- Given the training dataset of (data,label) pairs,

$$D = \left\{ \left( x^{(i)}, y^{(i)} \right) \right\}_{i=1.2,\dots,N}$$

$$x^{(i)} = (x_1^{(i)}, x_2^{(i)}, \dots x_n^{(i)})^T$$

$y^{(i)} =$ output data(label) of $i^{th}$ training example

let the machine learn a function from data to label

$$y \approx f_\theta(x) \quad \Rightarrow \quad f_\theta(x) = \theta_1 \, x_1 + \theta_2 \, x_2 + \cdots + \theta_n \, x_n + \theta_0$$

- Function set $\{f_\theta(x^{(i)})\}$ is called hypothesis space
- Learning is referred to as updating the parameter $\theta$ to make the prediction closed to the corresponding label

# 线性回归模型
# Linear Regression Model

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

sample $x$

features/variables: $x_1, x_2, \ldots x_n$

# 线性回归模型
# Linear Regression Model

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

sample $x$

One feature/variable: $x$

# 线性回归模型
# Linear Regression Model

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

sample $x$

One feature/variable: $x$

Linear regression
with one variable

(One-dimensional linear regression)



$f(x)$

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

$x$

Linear Regression

# 线性回归模型
# Linear Regression Model

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

sample $x$

One feature/variable: $x$

Linear regression
with one variable

quadratic regression
with one variable

(One-dimensional regression)

$f(x)$

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

Linear Regression

$x$

$f(x)$

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

Quadratic Regression
(A kind of generalized
linear model)

$x$

# 线性回归模型
# Linear Regression Model

$$f_\theta(x) = \theta_1 x_1 + \theta_2 x_2 + \theta_0$$

sample $x$

Two features/variables: $x_1, x_2$

Linear regression
with two variable
(two-dimensional
linear regression)

$$f_\theta(x) = \theta_0 + \theta_1 x_1 + \theta_2 x_2$$

$f(x)$

# 单变量线性回归
# Linear regression with one variable

Housing Prices
(Portland, OR)

Price
(in 1000s
of dollars)



Size (feet²)

Supervised Learning

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

Regression Problem

Predict real-valued output

# 单变量线性回归
## Linear regression with one variable

Training set of
housing prices
(Portland, OR)

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

Notation:

$N$ = Number of training examples

$x$ = "input" variable / features

$y$ = "output" variable / "target" variable

# 单变量线性回归
## Linear regression with one variable



How do we represent  f?

# 单变量线性回归
# Linear regression with one variable

Training Set

Learning Algorithm

Size of
house → $f$ → Estimated
price

**How do we represent  f?**

$$f_\theta(x) = \theta_0 + \theta_1\, x$$

Linear regression with one variable.
Univariate(one variable ) linear
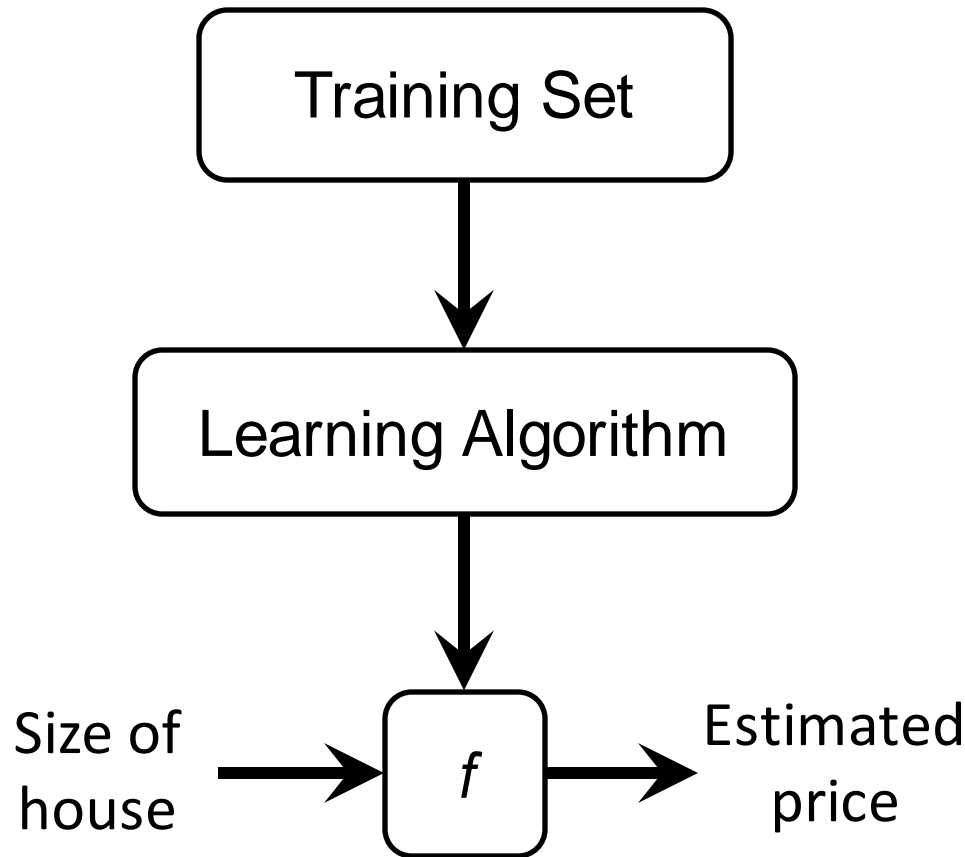regression.

# 单变量线性回归
# Linear regression with one variable

Training Set

↓

Learning Algorithm

↓

Size of house → $f$ → Estimated price

**How do we represent f?**

$$f_\theta(x) = \theta_0 + \theta_1\, x$$

$\theta_0, \theta_1$ : Parameters

Linear regression with one variable.
Univariate(one variable ) linear regression.

# 单变量线性回归
# Linear regression with one variable

Training Set

Learning Algorithm

Size of house → $f$ → Estimated price

**How do we represent f?**

$$f_\theta(x) = \theta_0 + \theta_1\, x$$

$\theta_0, \theta_1$: Parameters

**How to choose** $\theta_0, \theta_1$ **?**

Linear regression with one variable.
Univariate(one variable ) linear regression.

# 单变量线性回归
# Linear regression with one variable



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

Training Set

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

# 单变量线性回归
## Linear regression with one variable

Training Set

$\downarrow$

Learning Algorithm

$\downarrow$

Size of house $\rightarrow$ $f$ $\rightarrow$ Estimated price

Hypothesis:

$$f_\theta(x) = \theta_0 + \theta_1\, x$$

Parameters:

$$\theta_0, \theta_1$$

Cost Function:

$$J(\theta_0, \theta_1) = \frac{1}{2N} \sum_{i=1}^{N} (f_\theta(x^{(i)}) - y^{(i)})^2$$

Goal:

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

# 单变量线性回归
# Linear regression with one variable

Hypothesis:

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

Parameters:

$$\theta_0, \theta_1$$

Cost  Function:

$$J(\theta_0, \theta_1) = \frac{1}{2N} \sum_{i=1}^{N} (f_\theta(x^{(i)}) - y^{(i)})^2$$

Goal: $\displaystyle\operatorname*{minimize}_{\theta_0, \theta_1} J(\theta_0, \theta_1)$

# 单变量线性回归
# Linear regression with one variable

Hypothesis:

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

Parameters:

$$\theta_0, \theta_1$$

Cost Function:

$$J(\theta_0, \theta_1) = \frac{1}{2N} \sum_{i=1}^{N} (f_\theta(x^{(i)}) - y^{(i)})^2$$

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

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

$$\theta_1$$

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

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

# 单变量线性回归
# Linear regression with one variable

$$f_\theta(x)$$

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



$$f_\theta(x^i) = \theta_1 x^{(i)}$$

$$J(\theta_1)$$

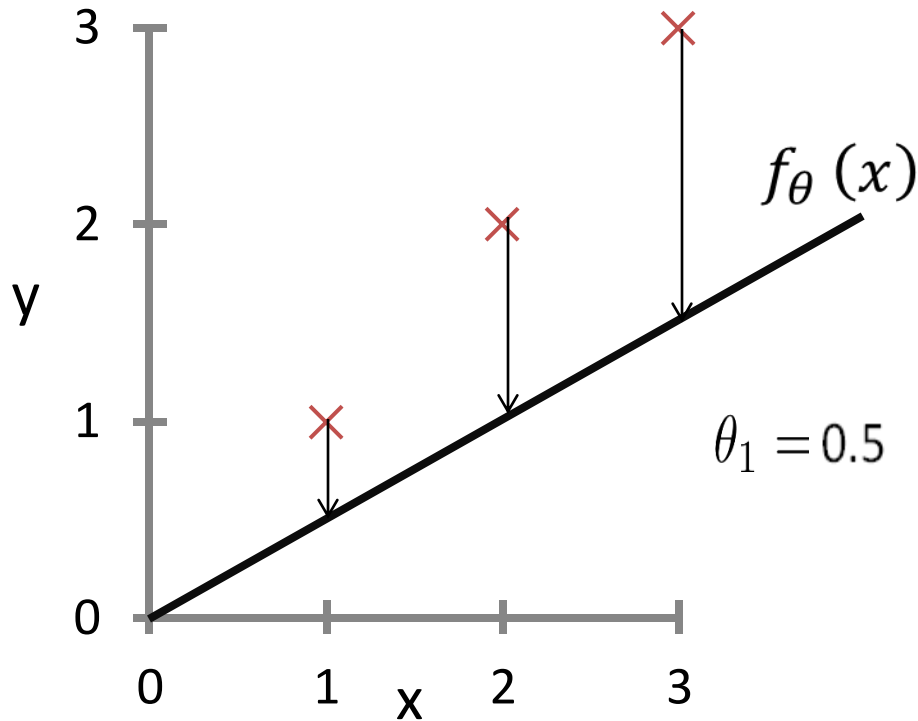(function of the parameter $\theta_1$)



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

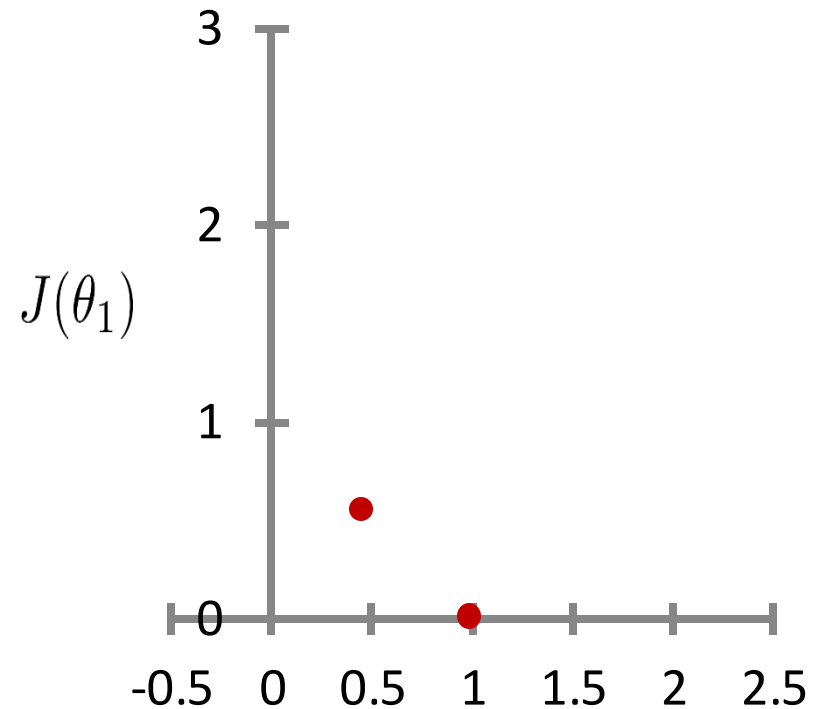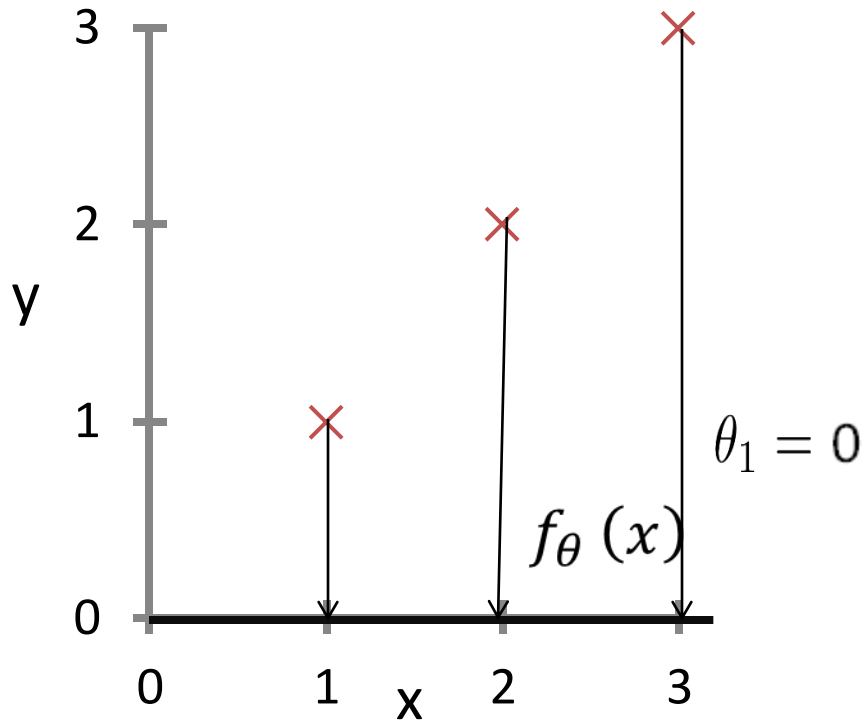# 单变量线性回归
# Linear regression with one variable

$f_\theta(x)$

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

$J(\theta_1)$

(function of the parameter $\theta_1$)



$f_\theta(x)$

$\theta_1 = 1$

$f_\theta(x^i) = \theta_1 x^{(i)}$

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

# 单变量线性回归
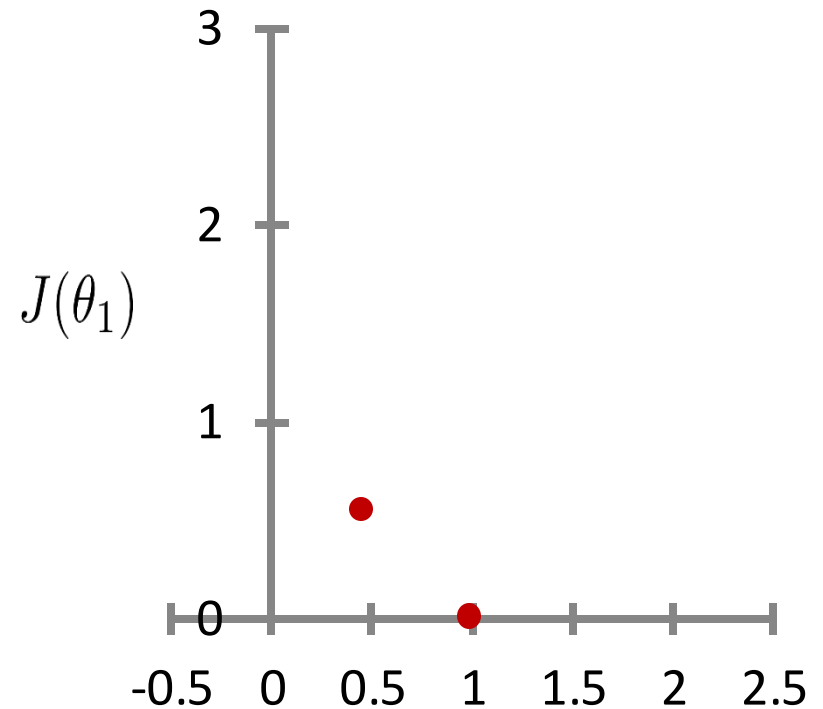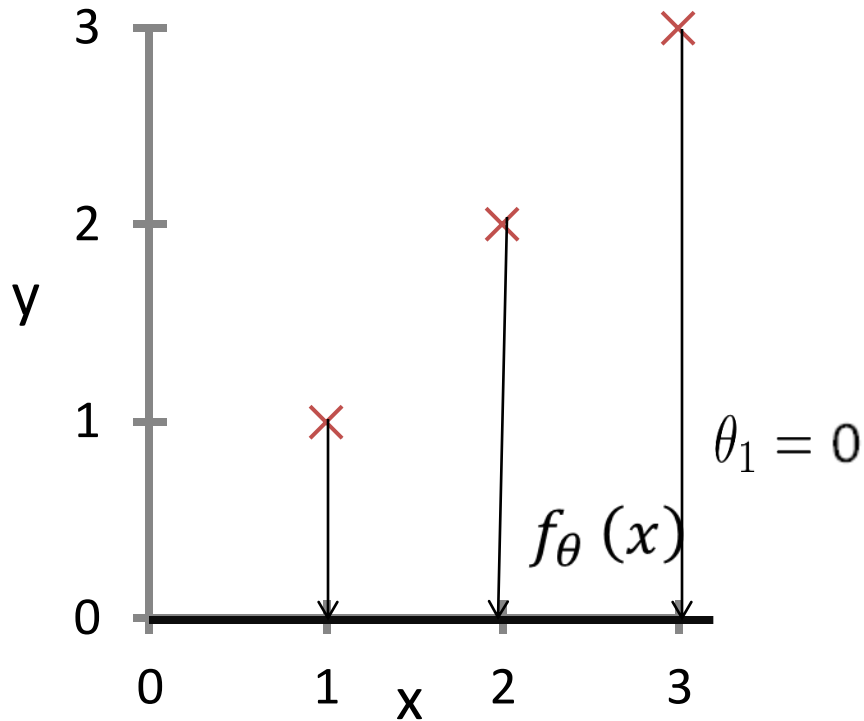# Linear regression with one variable

$$f_\theta(x)$$

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

$$J(\theta_1)$$

(function of the parameter $\theta_1$)



$f_\theta(x)$

$\theta_1 = 0.5$

$$f_\theta(x^i) = \theta_1 x^{(i)}$$

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

# 单变量线性回归
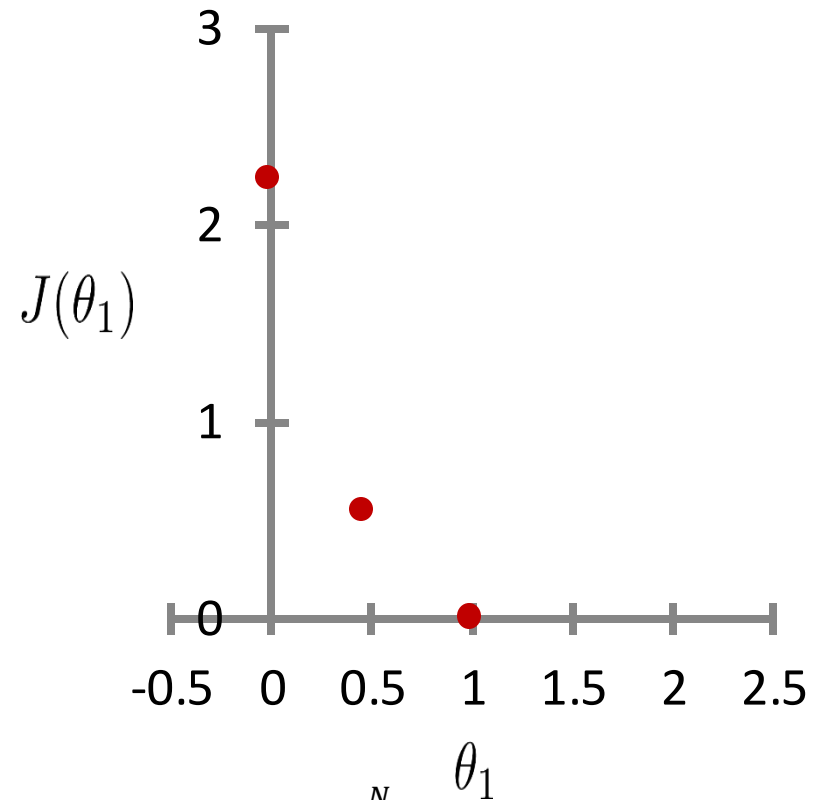# Linear regression with one variable

$$f_\theta(x)$$

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

$$J(\theta_1)$$

(function of the parameter $\theta_1$)



$f_\theta(x)$

$\theta_1 = 0.5$
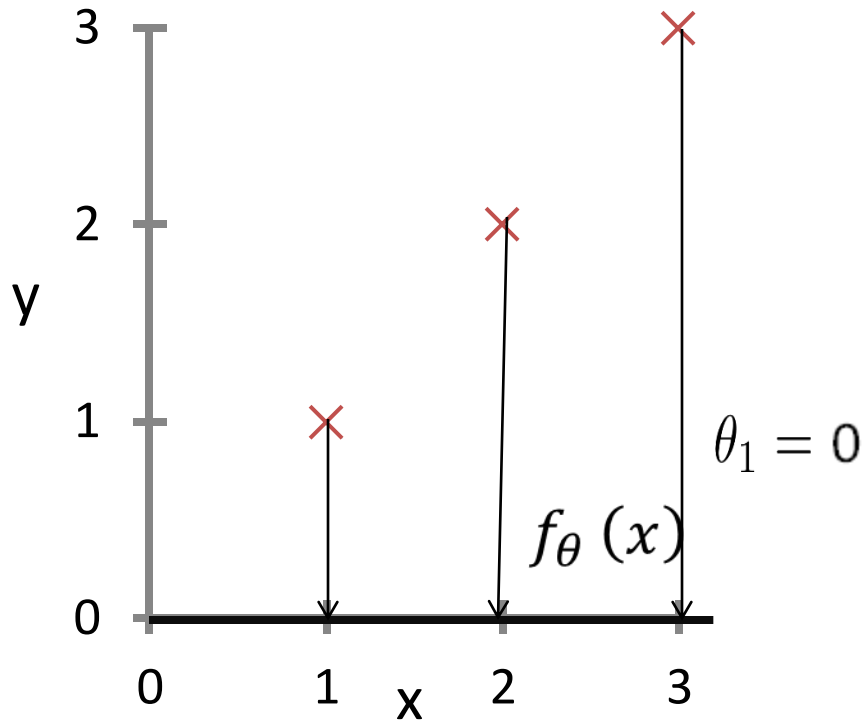
$$f_\theta(x^i) = \theta_1 x^{(i)}$$

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

# 单变量线性回归
# Linear regression with one variable

$f_\theta(x)$

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

$J(\theta_1)$

(function of the parameter $\theta_1$)



$\theta_1 = 0$

$f_\theta(x)$

$$f_\theta(x^i) = \theta_1 x^{(i)}$$

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

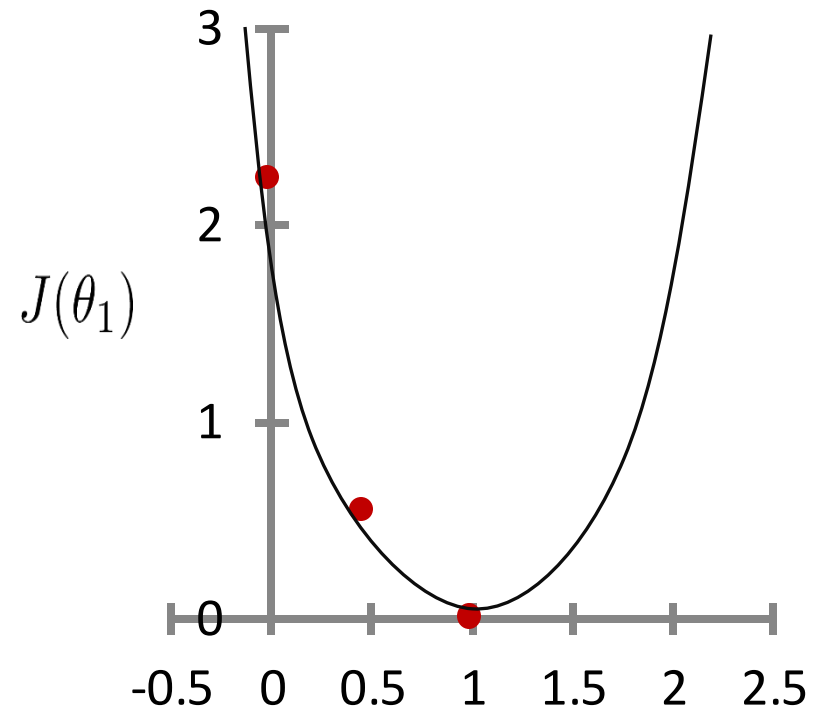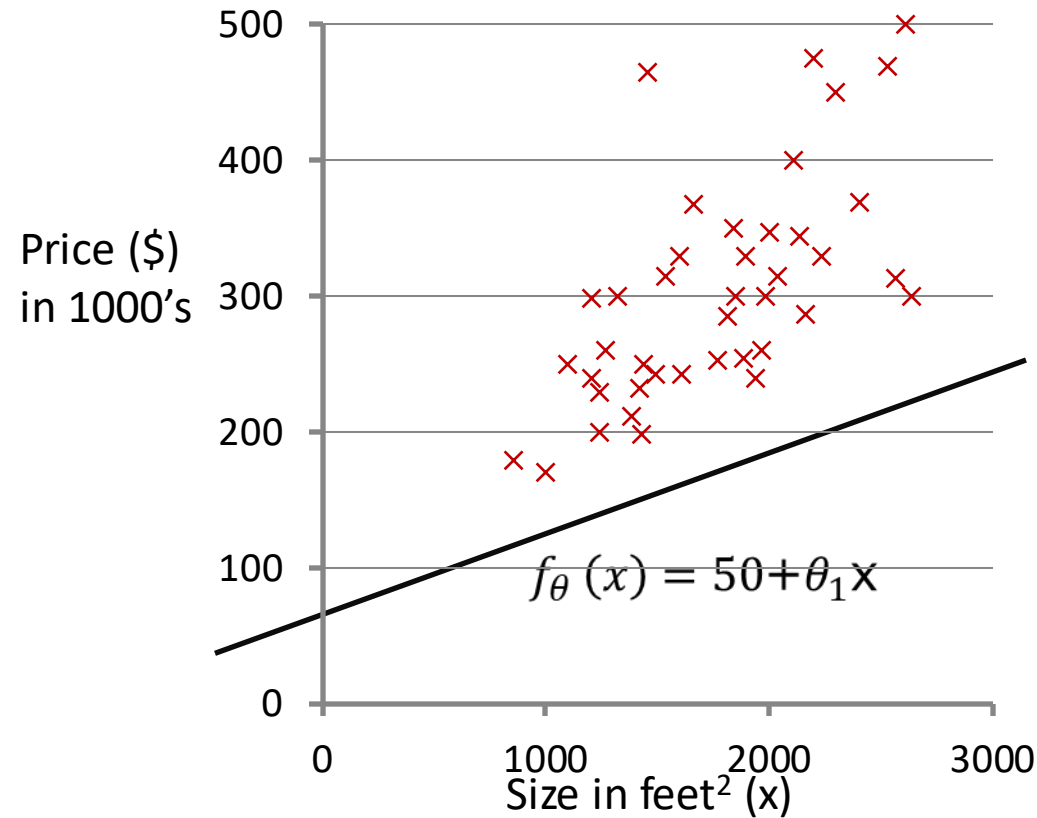# 单变量线性回归
# Linear regression with one variable

$f_\theta(x)$

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

$J(\theta_1)$

(function of the parameter $\theta_1$)



$\theta_1 = 0$

$f_\theta(x)$

$$f_\theta(x^i) = \theta_1 x^{(i)}$$

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

# 单变量线性回归
# Linear regression with one variable

$f_\theta(x)$

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

$J(\theta_1)$

(function of the parameter $\theta_1$)



$f_\theta(x^i) = \theta_1 x^{(i)}$

$J(\theta_1) = \dfrac{1}{2N} \sum_{i=1}^{N} (f_\theta(x^{(i)}) - y^{(i)})^2$

# 单变量线性回归
# Linear regression with one variable

$f_\theta(x)$

$J(\theta_0, \theta_1)$



$f_\theta(x) = 50 + \theta_1 x$

Price ($) in 1000's

Size in feet$^2$ (x)
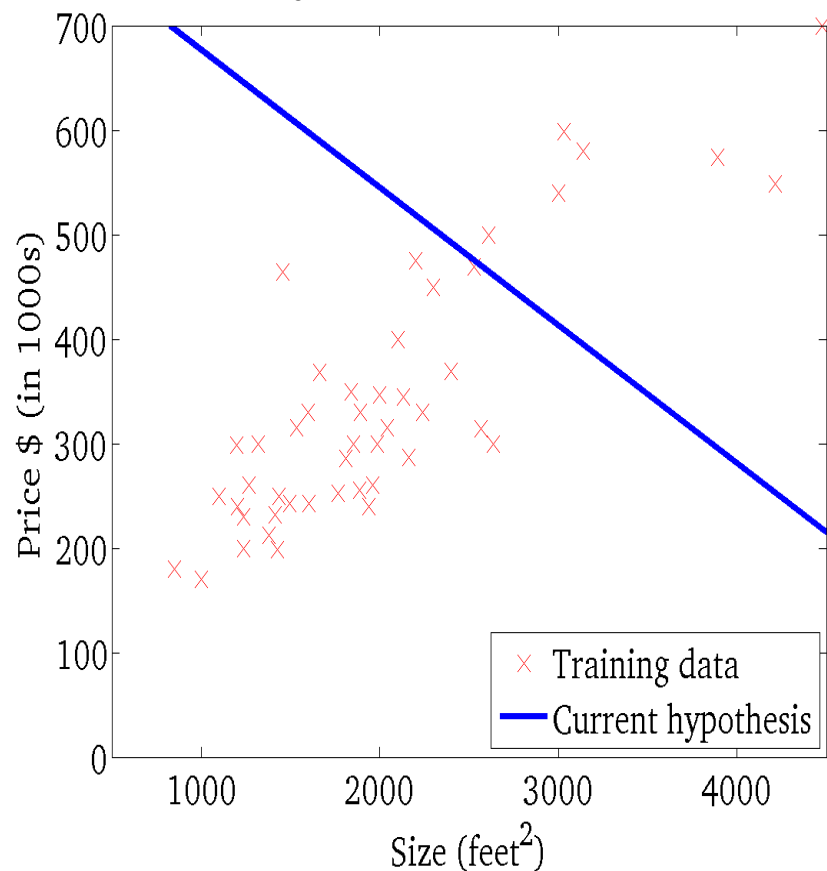
$f_\theta(x) = \theta_0 + \theta_1 x$

$$J(\theta_0, \theta_1) = \frac{1}{2N} \sum_{i=1}^{N} (f_\theta(x^{(i)}) - y^{(i)})^2$$

# 单变量线性回归
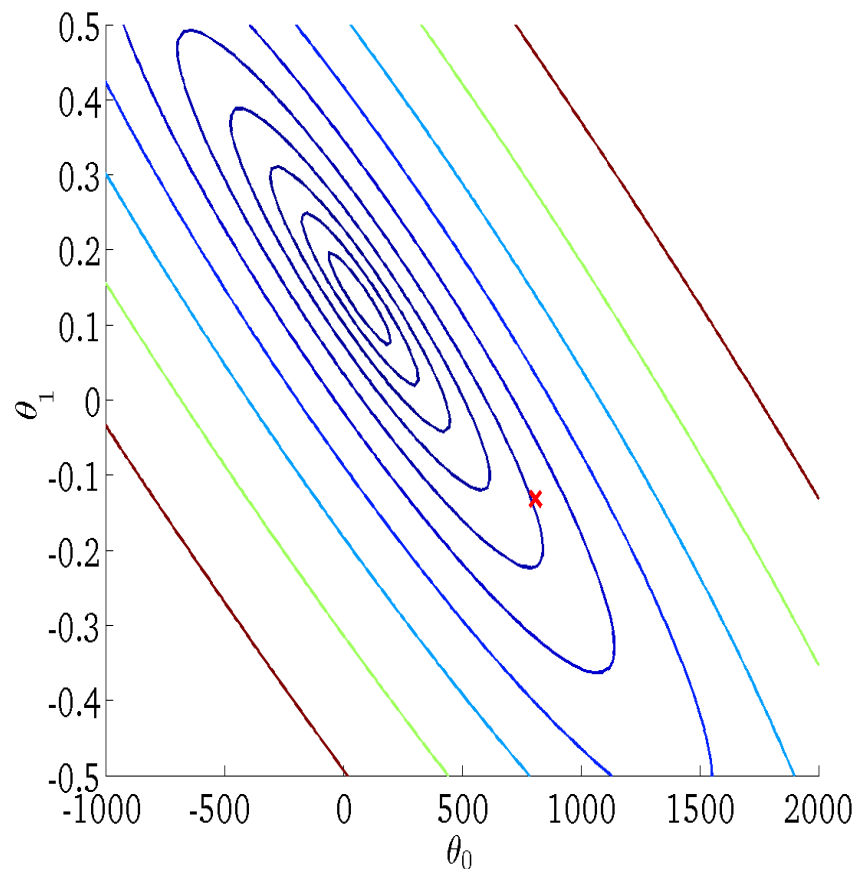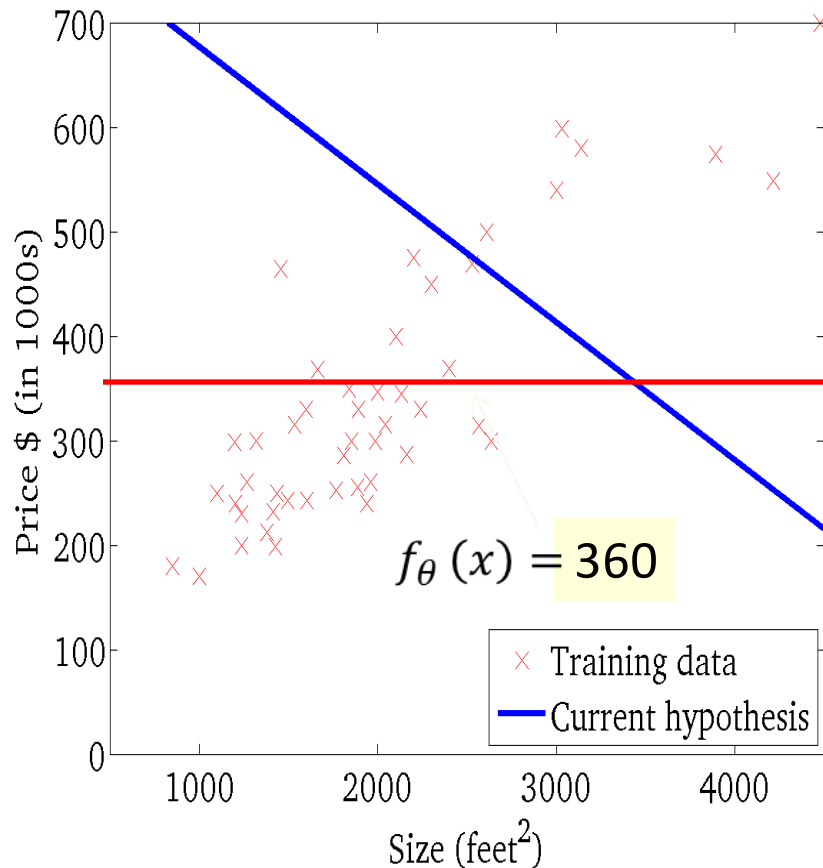# Linear regression with one variable

$$f_\theta(x)$$

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

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

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



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

$$J(\theta_0, \theta_1) = \frac{1}{2N} \sum_{i=1}^{N} (f_\theta(x^{(i)}) - y^{(i)})^2$$

# 单变量线性回归
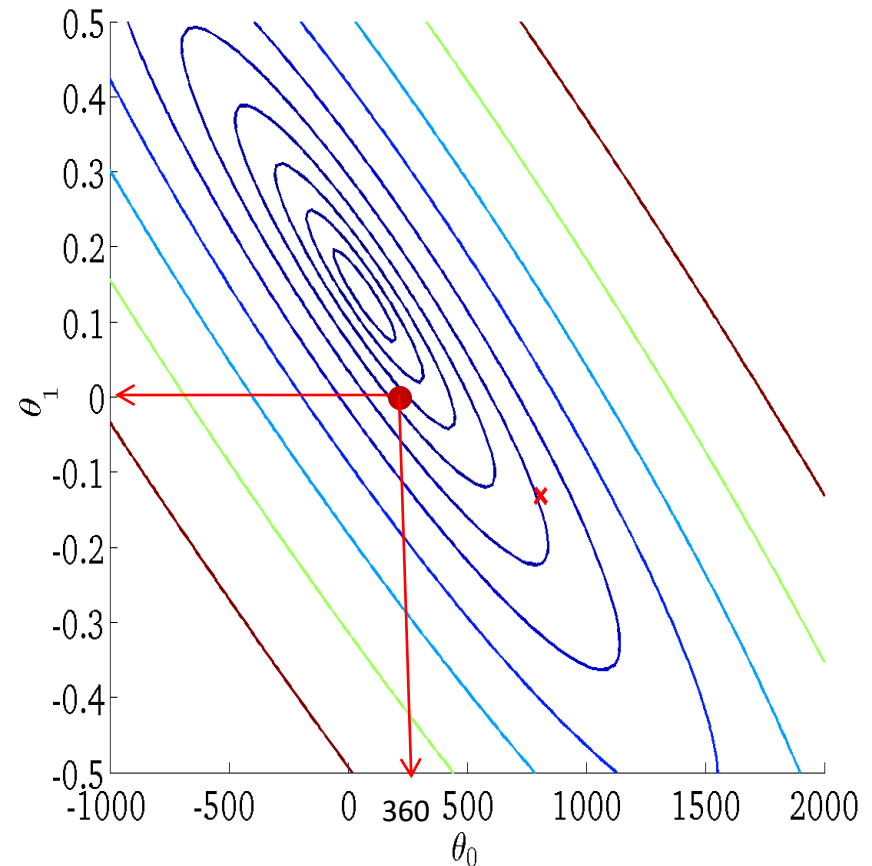# Linear regression with one variable

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

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



$f_\theta(x) = 360$

Training data ×

Current hypothesis ——

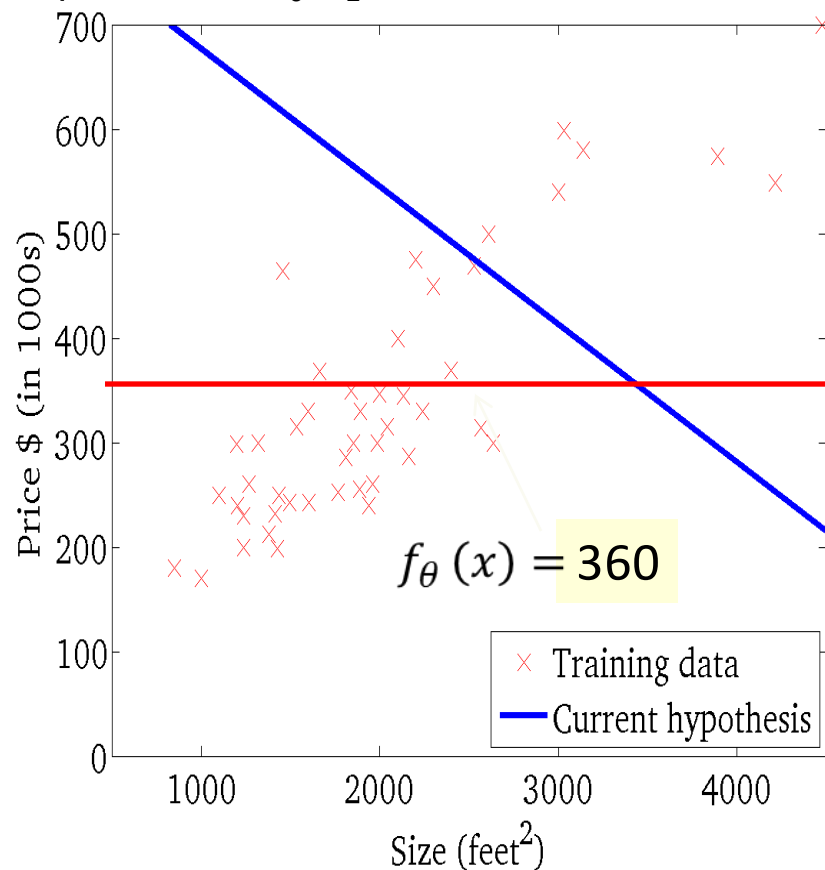$$f_\theta(x) = \theta_0 + \theta_1 x$$

$$J(\theta_0, \theta_1) = \frac{1}{2N} \sum_{i=1}^{N} (f_\theta(x^{(i)}) - y^{(i)})^2$$

# 单变量线性回归
# Linear regression with one variable
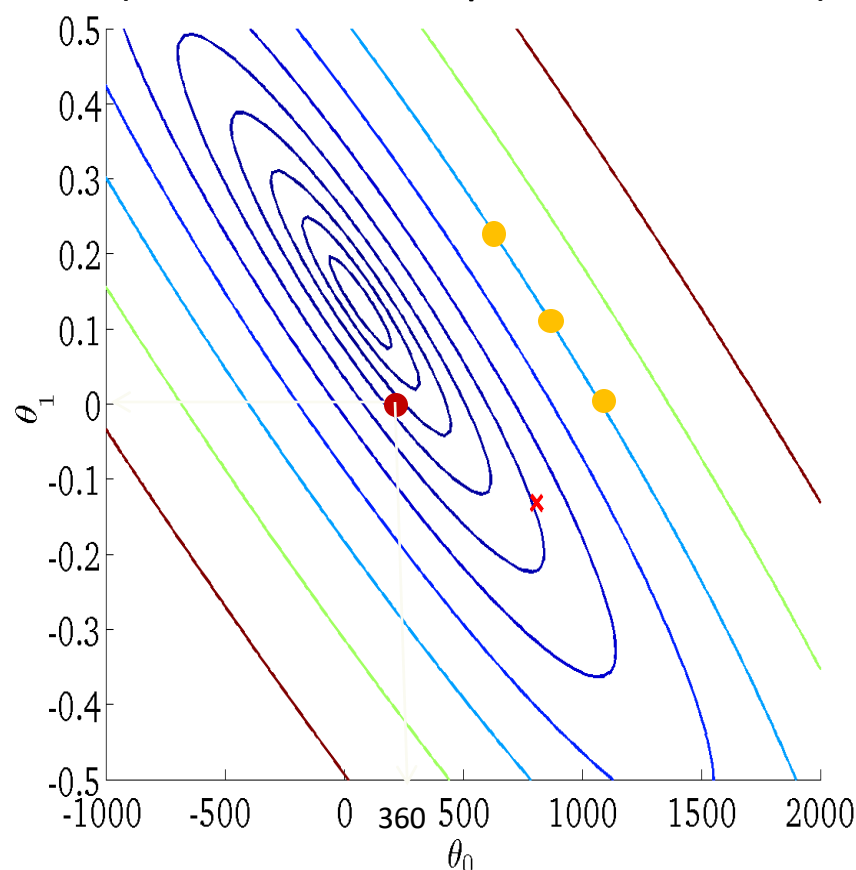
$f_\theta(x)$ · · · · · · · · · · · · · · · · · · · · · · · $J(\theta_0, \theta_1)$

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



$f_\theta(x) = 360$

$$f_\theta(x) = \theta_0 + \theta_1 x \qquad\qquad J(\theta_0, \theta_1) = \frac{1}{2N} \sum_{i=1}^{N} (f_\theta(x^{(i)}) - y^{(i)})^2$$
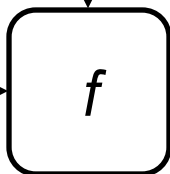
# 单变量线性回归
# Linear regression with one variable

Training Set

Learning Algorithm

Size of house $\rightarrow$ $f$ $\rightarrow$ Estimated price

- 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

Hypothesis:

$$f_\theta(x) = \theta_0 + \theta_1\,x$$

Parameters:

$$\theta_0, \theta_1$$

Cost Function:

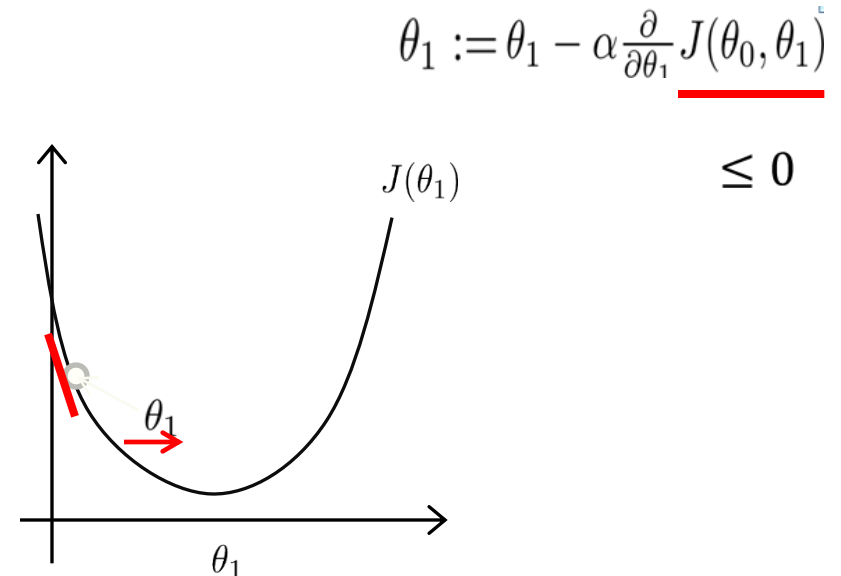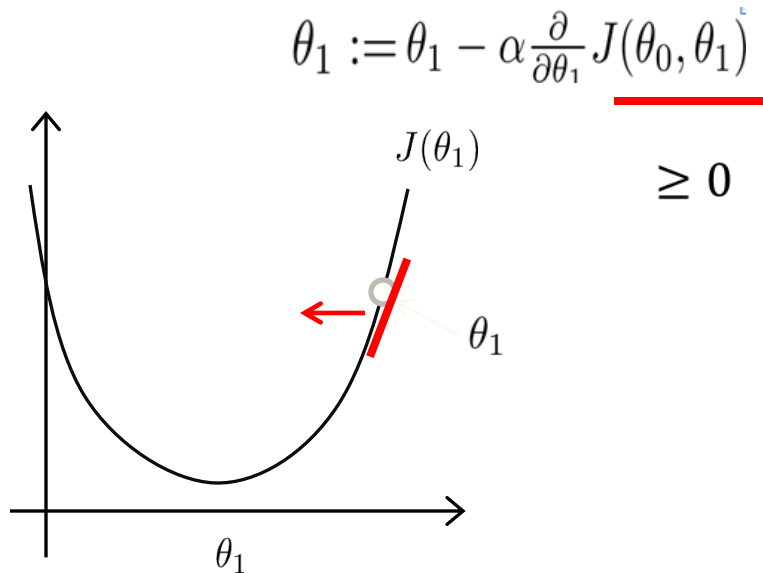$$J(\theta_0, \theta_1) = \frac{1}{2N}\sum_{i=1}^{N}(f_\theta(x^{(i)}) - y^{(i)})^2$$

Goal:

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

# 梯度下降法
## Gradient descent algorithm

repeat until convergence {

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

}

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

$\geq 0$

$J(\theta_1)$

$\theta_1$

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

$\leq 0$

$J(\theta_1)$

$\theta_1$

# 梯度下降法
## 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 $j = 1$)

}

# 梯度下降法
## 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 $j = 1$)

}

---

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

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

$\theta_0 := \text{temp0}$

$\theta_1 := \text{temp1}$

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

$\theta_0 := \text{temp0}$

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

$\theta_1 := \text{temp1}$

# 梯度下降法
# 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 $j = 1$)

}

---

**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:**

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

$\theta_0 := \text{temp0}$

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

$\theta_1 := \text{temp1}$

# 梯度下降法
# Gradient descent algorithm

### Gradient descent algorithm

repeat until convergence {

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

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

}

### Linear Regression Model

$$f_\theta(x) = \theta_0 + \theta_1 \, x$$

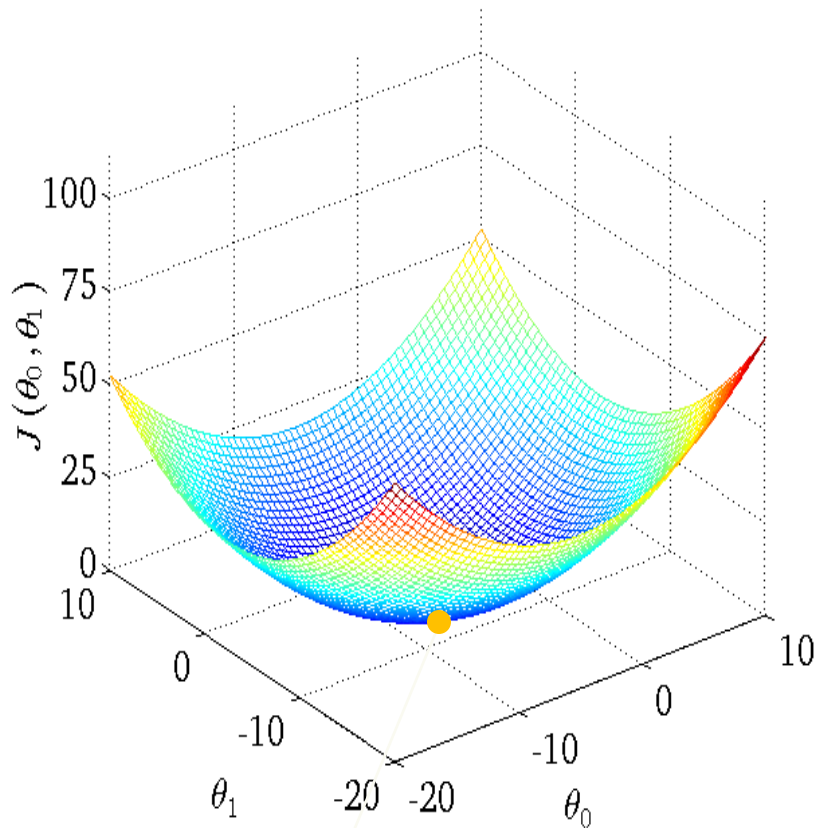$$J(\theta_0, \theta_1) = \frac{1}{2N} \sum_{i=1}^{N} (f_\theta(x^{(i)}) - y^{(i)})^2$$

### Repeat until convergece

$$\theta_0 := \theta_0 - a \frac{1}{N} \sum_{i=1}^{N} (f_\theta(x^{(i)}) - y^{(i)})$$

$$\theta_1 := \theta_1 - a \frac{1}{N} \sum_{i=1}^{N} (f_\theta(x^{(i)}) - y^{(i)}) x^{(i)}$$

update $\theta_0$ and $\theta_1$ simultaneously

# 凸函数
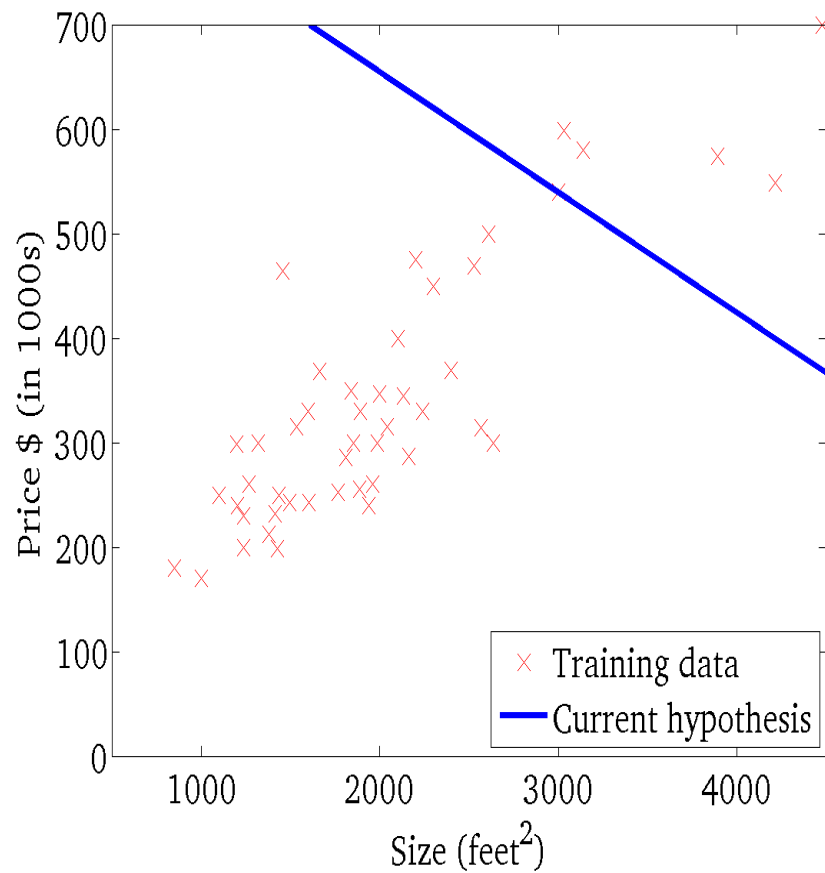# Bowled shape Convex Function



$$J(\theta_0, \theta_1) = \frac{1}{2N} \sum_{i=1}^{N} (f_\theta(x^{(i)}) - y^{(i)})^2$$

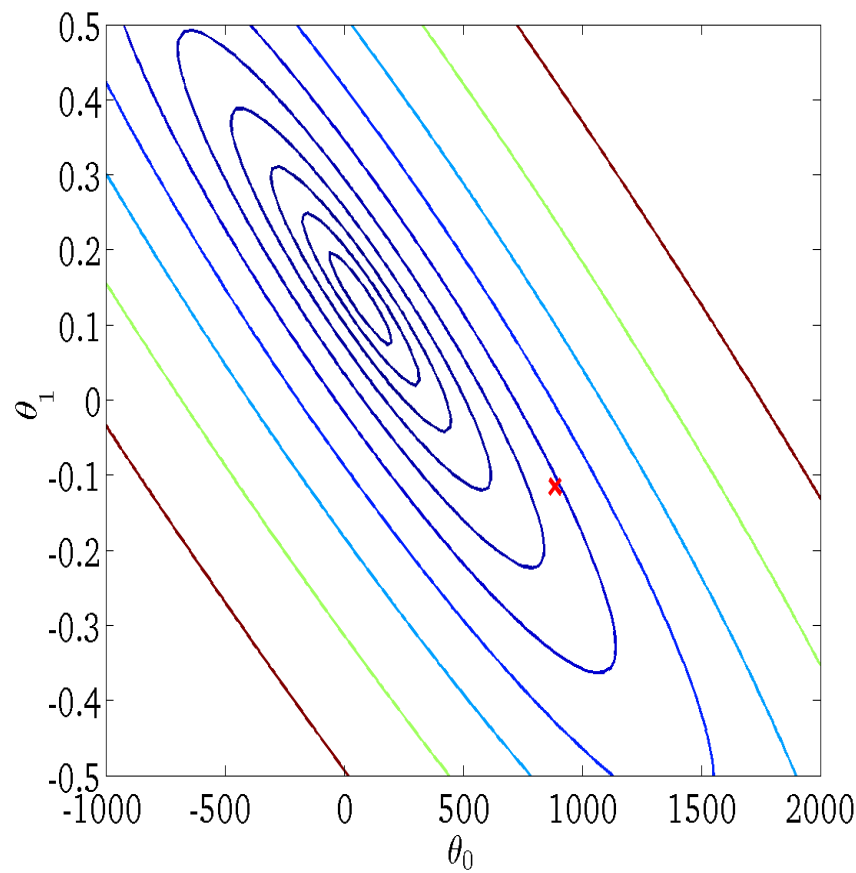Different initial lead to the same optimum

Unique Minimum

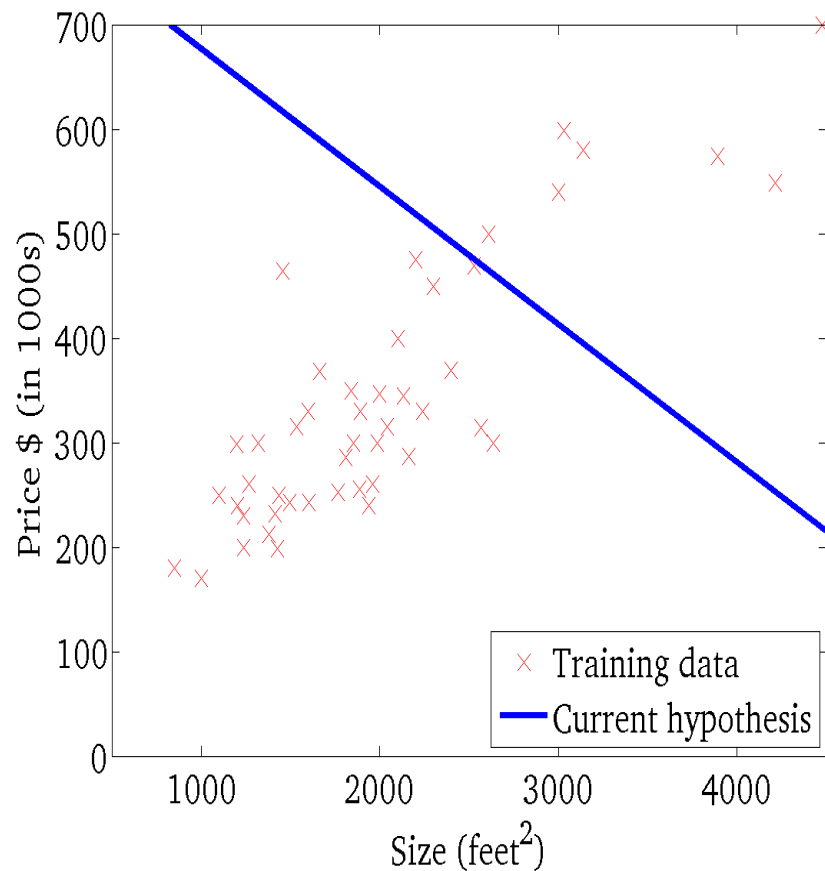# 搜索过程
# Search Procedure



$$f_{\theta}(x)$$

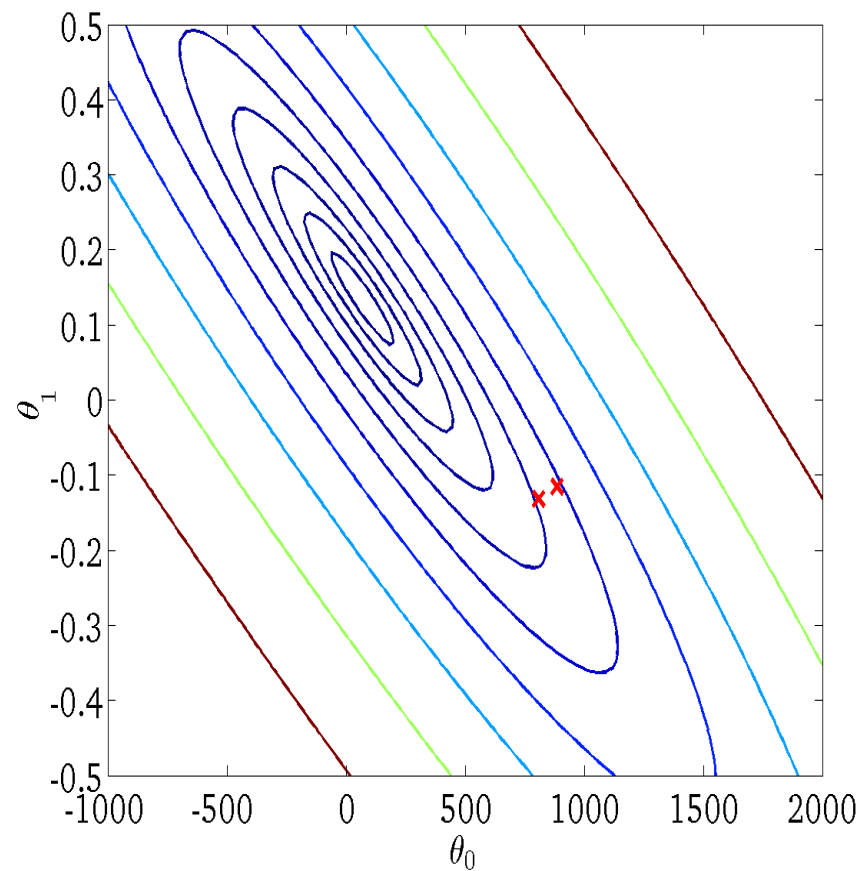(for fixed $\theta_0, \theta_1$ , this is a function of x)

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

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

# 搜索过程
# Search Procedure



$$f_\theta\left(x\right)$$

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

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

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

# 搜索过程
# Search Procedure



$$f_\theta\left(x\right)$$

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

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

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

# 搜索过程
# Search Procedure



$f_{\boldsymbol{\theta}}\left(x\right)$

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

$J(\theta_0, \theta_1)$

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

# 搜索过程
# Search Procedure
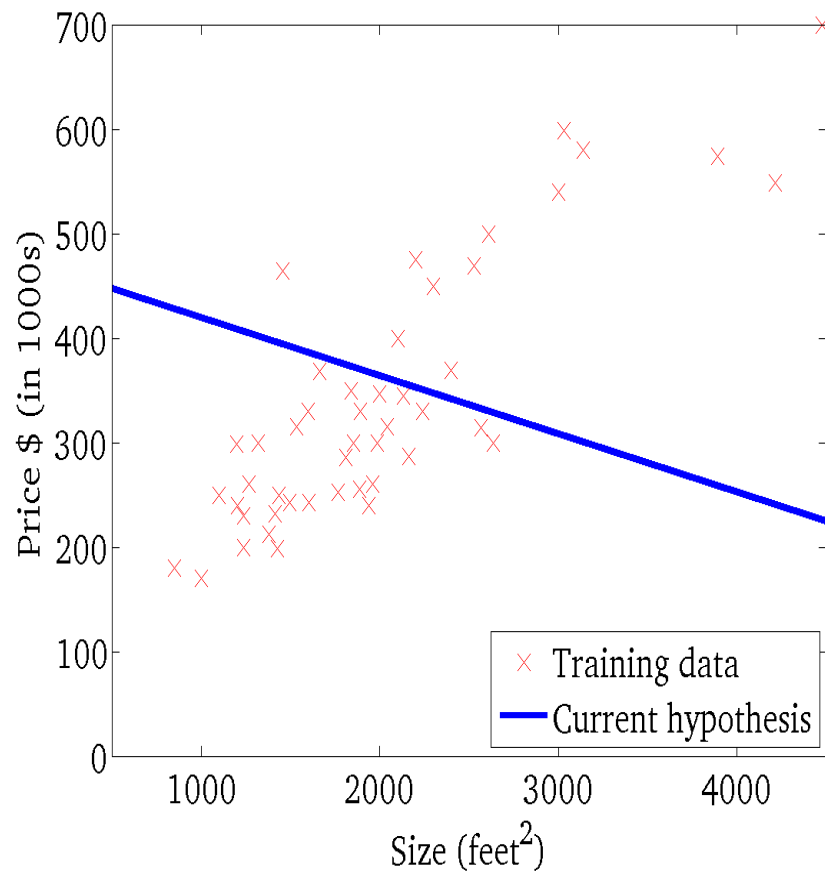


$$f_{\theta}(x)$$
(for fixed $\theta_0, \theta_1$ , this is a function of x)

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

# 搜索过程
# Search Procedure



$f_{\theta}(x)$

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

$J(\theta_0, \theta_1)$

(function of the parameters $\theta_0, \theta_1$ )
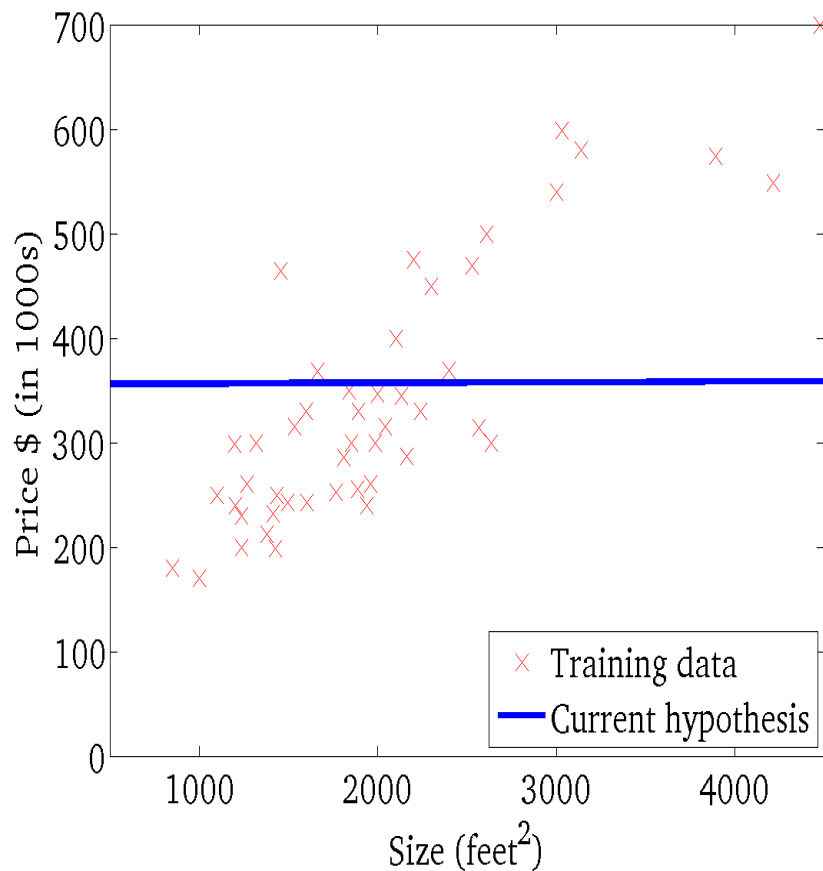
# 搜索过程
# Search Procedure



$$f_{\theta}(x)$$

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

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

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

# 搜索过程
# Search Procedure



$$f_{\theta}\left(x\right)$$
(for fixed $\theta_0, \theta_1$ , this is a function of x)

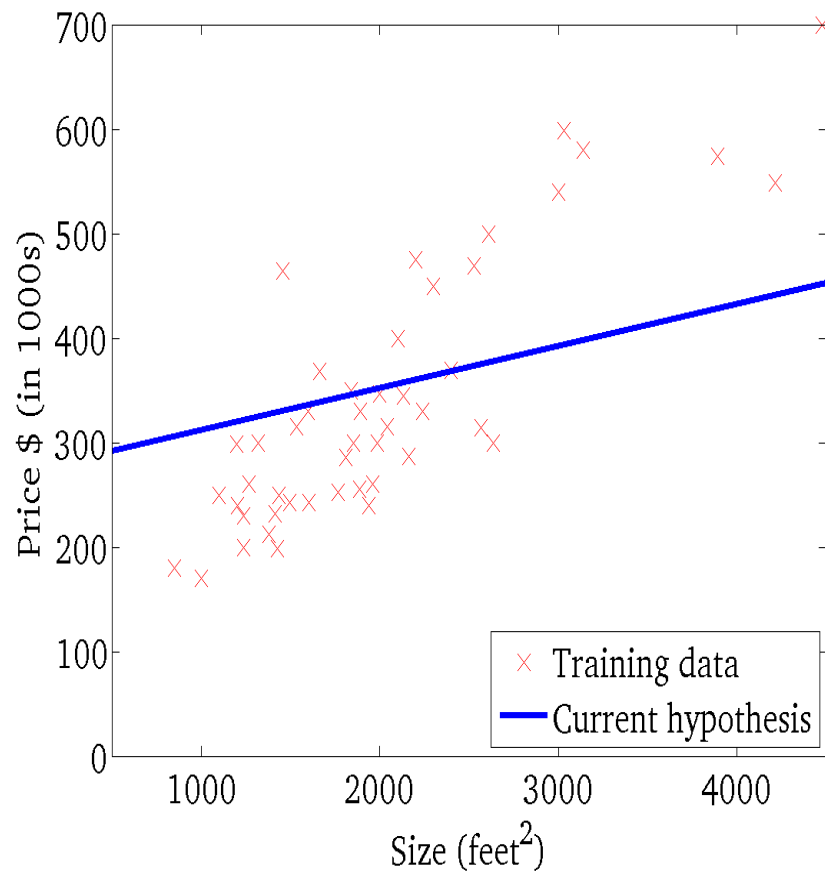$$J(\theta_0, \theta_1)$$
(function of the parameters $\theta_0, \theta_1$ )

# 搜索过程
# Search Procedure



$$f_\theta\,(x)$$

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

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

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

# 搜索过程
# Search Procedure



- Choose an initial value for $\theta$
- Update $\theta$ iteratively with the data
- Until we research a minimum

# 搜索过程
# Search Procedure



- Choose a new initial value for $\theta$
- Update $\theta$ iteratively with the data
- Until we research a minimum

# 搜索过程
# Search Procedure



$J(\theta_0, \theta_1)$

$\theta_0$

$\theta_1$

- Choose a new initial value for $\theta$
- Update $\theta$ iteratively with the data
- Until we research a minimum

In linear regression, the loss function L is convex. Different initial lead to the same optimum.

# 批量梯度下降
# Batch Gradient descent

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

Repeat until convergence

$$\theta_0 := \theta_0 - a\frac{1}{N}\sum_{i=1}^{N}(f_\theta(x^{(i)}) - y^{(i)})$$

$$\theta_1 := \theta_1 - a\frac{1}{N}\sum_{i=1}^{N}(f_\theta(x^{(i)}) - y^{(i)})x^{(i)}$$


loss w.r.t. parameters round = 200

# 随机梯度下降
# Stochastic Gradient descent

"stochastic": Each step of gradient descent uses single  training example.

Repeat until convergence

$$\theta_0 := \theta_0 - a\frac{1}{N}\sum_{i=1}^{N}(f_\theta(x^{(i)}) - y^{(i)})$$

$$\theta_1 := \theta_1 - a\frac{1}{N}\sum_{i=1}^{N}(f_\theta(x^{(i)}) - y^{(i)})x^{(i)}$$


loss w.r.t. parameters round = 50 case = 2

Compare with BGD
- Faster learning
- Uncertainty or fluctuation in learning

# 小批量梯度下降
# Mini-Batch Gradient descent

- A combination of batch GD and stochastic GD

- Split the whole dataset into $K$ mini-batches

$$\{1, 2, 3, \ldots, K\}$$

- For each mini-batch $k$, perform one-step BGD toward

$$J^k(\theta) := \frac{1}{2N_k} \sum_{i=1}^{N_k} (f_\theta(x^{(i)}) - y^{(i)})^2$$

- Good learning stability (BGD)
- Good convergence rate (SGD)

- Update $\theta_{\text{new}} = \theta_{\text{old}} - \eta \frac{\partial J^{(k)}(\theta)}{\partial \theta}$ for each mini-batch

$$\theta_0 := \theta_0 - a\frac{1}{N_k}\sum_{i=1}^{N_k}(f_\theta(x^{(i)}) - y^{(i)})$$

$$\theta_1 := \theta_1 - a\frac{1}{N_k}\sum_{i=1}^{N_k}(f_\theta(x^{(i)}) - y^{(i)})x^{(i)}$$

# 学习率选择
## Choose learning rate

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.



Initial weight

$J(\theta)$

Gradient

Global cost minimum

$J_{min}(\theta)$

$\theta$

$J(\theta)$

$\theta$

To see if gradient descent is working, print out for each or every $J(\theta)$ several iterations. If $J(\theta)$ does not drop properly, adjust the learning rate!

# 学习率选择
# Choose learning rate



Loss

Very Large

small

Large

Just make

Loss

Loss

No. of parameters updates

To see if gradient descent is working, print out for each or every $J(\theta)$ several iterations. If $J(\theta)$ does not drop properly, adjust the learning rate!

# 多变量线性回归
## Linear regression with multiple variable

| Size (feet²) | Price ($1000) |
|---|---|
| 2104 | 460 |
| 1416 | 232 |
| 1534 | 315 |
| 852 | 178 |
| … | … |

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

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

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

# 单变量线性回归
# Linear regression with one variable

Training Set

Learning Algorithm

Size of house → $f$ → Estimated price

- 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

Hypothesis:

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

Parameters:

$$\theta_0, \theta_1$$

Cost Function:

$$J(\theta_0, \theta_1) = \frac{1}{2N} \sum_{i=1}^{N} (f_\theta(x^{(i)}) - y^{(i)})^2$$

Goal:

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

# 多变量线性回归
# Linear regression with multiple variable

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

Parameters: $\theta_0, \theta_1, \ldots, \theta_n$

Cost function: $J(\theta_0, \theta_1, \ldots \theta_n) = \dfrac{1}{2N} \sum_{i=1}^{N} (f_\theta(x^{(i)}) - y^{(i)})^2$

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.

Gradient descent:

Repeat $\{$

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

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

# 多变量线性回归
## Linear regression with multiple variable

Previously (n=1):

Repeat $\Big\{$

$$\theta_0 := \theta_0 - a\frac{1}{N}\underbrace{\sum_{i=1}^{N}(f_\theta(x^{(i)}) - y^{(i)})}_{\frac{\partial}{\partial\theta_0}J(\theta)}$$

$$\theta_1 := \theta_1 - a\frac{1}{N}\sum_{i=1}^{N}(f_\theta(x^{(i)}) - y^{(i)})x^{(i)}$$

(simultaneously update $\theta_0, \theta_1$)

$\Big\}$

New algorithm $(n \geq 1)$ :

Repeat $\Big\{$

$$\theta_j := \theta_j - a\frac{1}{N}\sum_{i=1}^{N}(f_\theta(x^{(i)}) - y^{(i)})x_j^{(i)}$$

(simultaneously update $\theta_j$ for

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

$\Big\}$

---

$$\theta_0 := \theta_0 - a\frac{1}{N}\sum_{i=1}^{N}(f_\theta(x^{(i)}) - y^{(i)})x_0^{(i)}$$
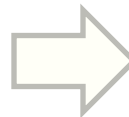
$$\theta_1 := \theta_1 - a\frac{1}{N}\sum_{i=1}^{N}(f_\theta(x^{(i)}) - y^{(i)})x_1^{(i)}$$

...

# 多变量线性回归
## Linear regression with multiple variable

| Size (feet²) | Price ($1000) |
|---|---|
| 2104 | 460 |
| 1416 | 232 |
| 1534 | 315 |
| 852 | 178 |
| … | … |

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

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

$$f_\theta(x) = \theta_0 + \theta_1 x_1 + \theta_2 x_2 + \cdots + \theta_n x_n$$
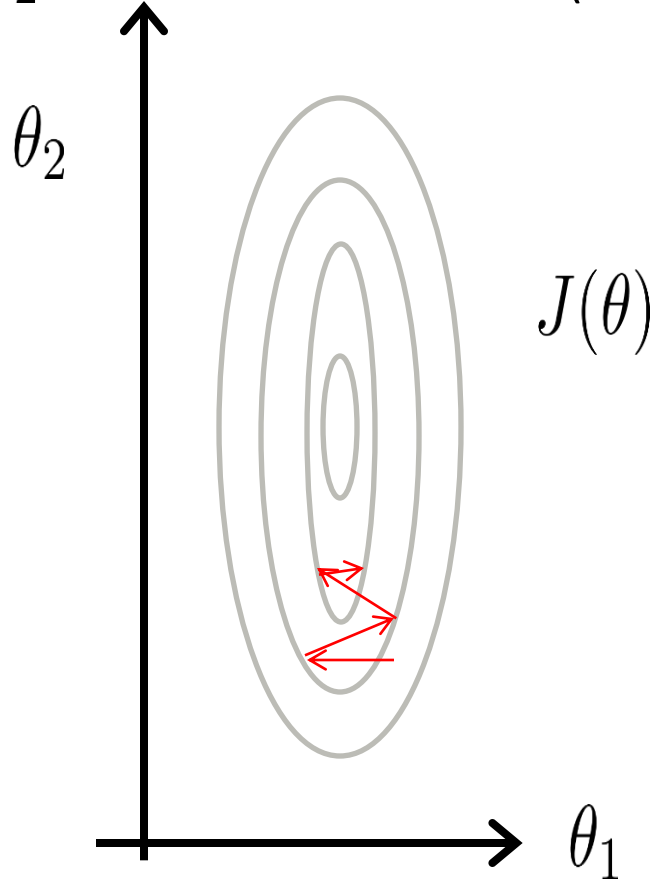
# 特征归一化
# Feature Scaling

$$f_\theta(x) = \theta_0 + \theta_1 \, x_1 + \theta_2 \, x_2$$

E.g. $x_1$ = size (0-2000 feet$^2$)

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

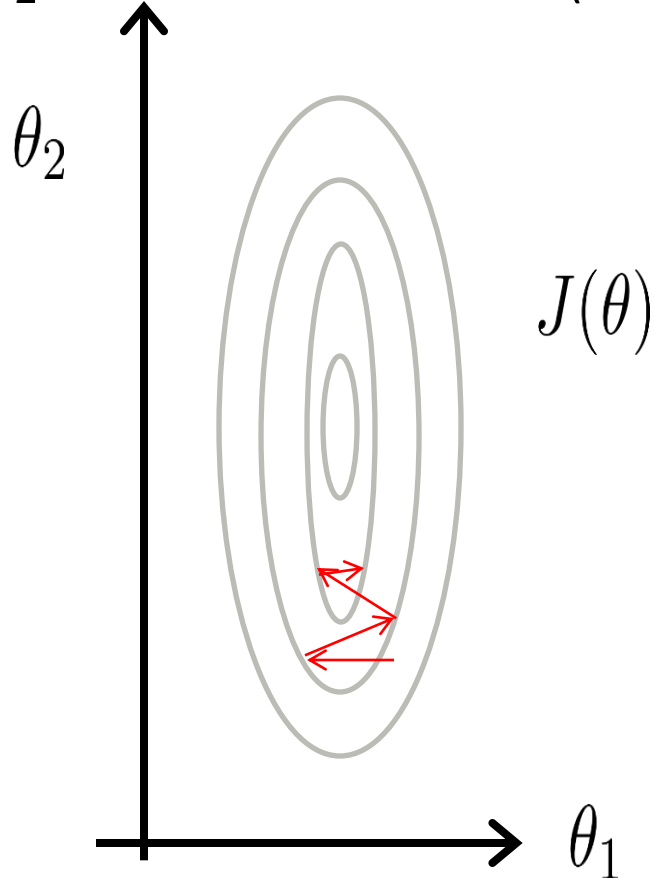# 特征归一化
# Feature Scaling

Idea: Make sure features are on a similar scale.

E.g. $x_1$ = size (0-2000 feet²)

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

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

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

# 特征归一化
# Feature Scaling

Get every feature into approximately a similar scale.

**Mean Normalization**

$$x' = \frac{x - mean(x)}{max(x) - min(x)}$$

**Standardization**

$$x' = \frac{x - mean(x)}{std(x)} \qquad std(x) = \sqrt{\frac{\sum(x - mean(x))^2}{n}}$$

e.g. Replace $x_i$ with $x_i - \mu_i$ to make features have approximately zero mean.

E.g.

$$x_1 = \frac{size - 1000}{2000}$$

$$x_2 = \frac{\#bedrooms - 2}{4}$$

$$-0.5 \le x_1 \le 0.5, -0.5 \le x_2 \le 0.5$$

# 多变量线性回归
## Linear regression with multiple variable

Previously (n=1):

Repeat $\{$

$$\theta_0 := \theta_0 - a\frac{1}{N}\underbrace{\sum_{i=1}^{N}(f_\theta(x^{(i)}) - y^{(i)})}_{\frac{\partial}{\partial\theta_0}J(\theta)}$$

$$\theta_1 := \theta_1 - a\frac{1}{N}\sum_{i=1}^{N}(f_\theta(x^{(i)}) - y^{(i)})x^{(i)}$$

(simultaneously update $\theta_0, \theta_1$)

$\}$

New algorithm $(n \geq 1)$ :

Repeat $\{$

$$\theta_j := \theta_j - a\frac{1}{N}\sum_{i=1}^{N}(f_\theta(x^{(i)}) - y^{(i)})x_j^{(i)}$$

(simultaneously update $\theta_j$ for

$\}$  $\qquad j = 0, \ldots, n$  )

$$\theta_0 := \theta_0 - a\frac{1}{N}\sum_{i=1}^{N}(f_\theta(x^{(i)}) - y^{(i)})x_0^{(i)}$$

$$\theta_1 := \theta_1 - a\frac{1}{N}\sum_{i=1}^{N}(f_\theta(x^{(i)}) - y^{(i)})x_1^{(i)}$$

$\ldots$

# 自适应的学习率
# Adaptive Learning Rates

## Adagrad

Divide the learning rate of each parameter by the root mean square of its previous derivatives

$$\theta^{(t+1)} := \theta^{(t)} - \frac{a}{\sqrt{\sum_{i=0}^{t}(g^{(i)})^2}} g^{(t)}$$

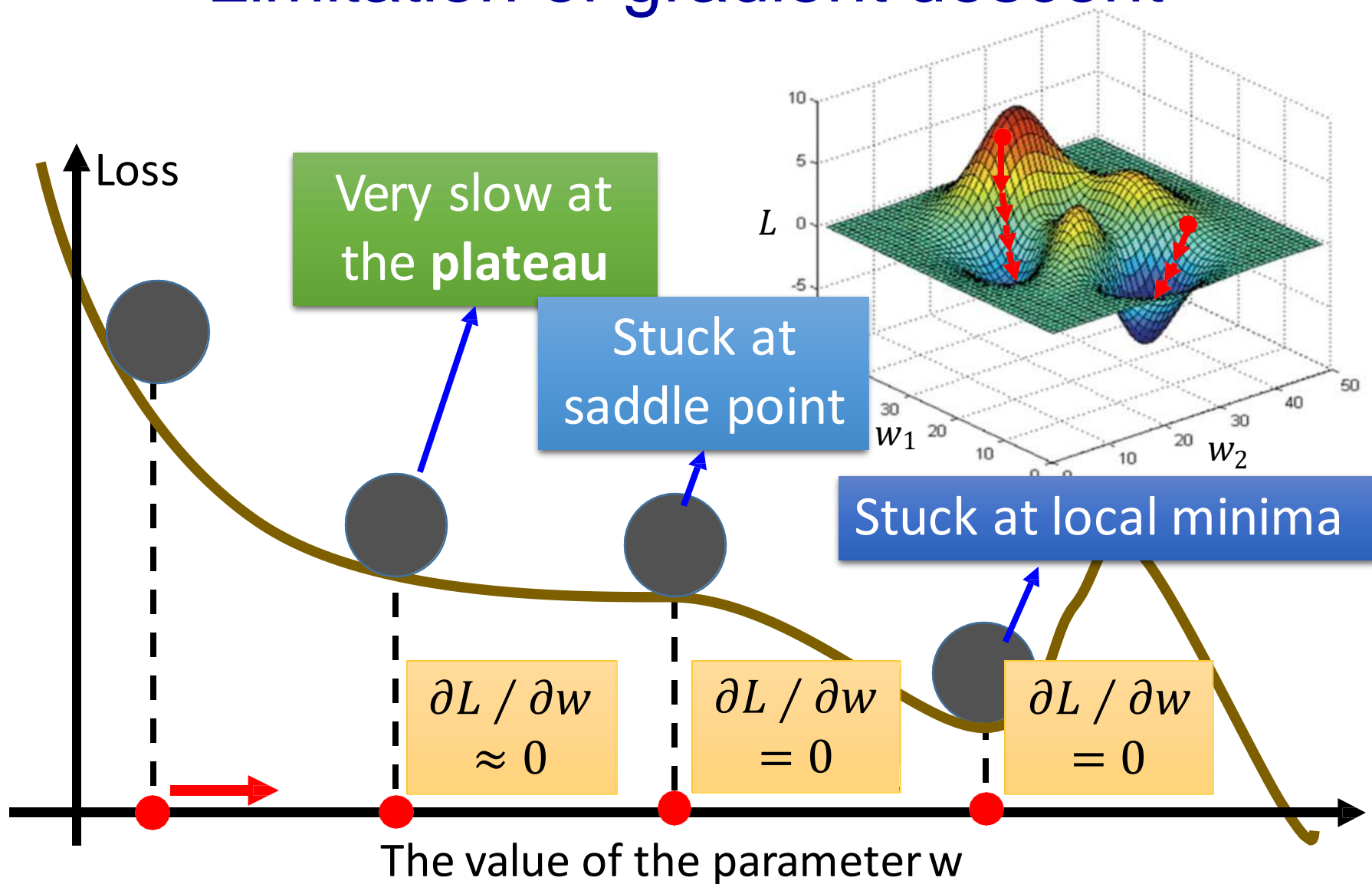$$g^{(t)} = \frac{\partial J(\theta^{(t)})}{\partial \theta}$$

$\theta_2$

$J(\theta)$

$\theta_1$

adapts the learning rate to the parameters, performing larger updates for infrequent and smaller updates for frequent parameters

# 梯度的限制
# Limitation of gradient descent

# 梯度下降优化算法
# Gradient descent optimization algorithms

- **Momentum**

  helps accelerate SGD in the relevant direction and dampens oscillations

- **Adagrad**

  （Adaptive Gradient）

  adapts the learning rate to the parameters, performing larger updates for infrequent and smaller updates for frequent parameters

- **RMSProp**

  （Root Mean Square propagation）

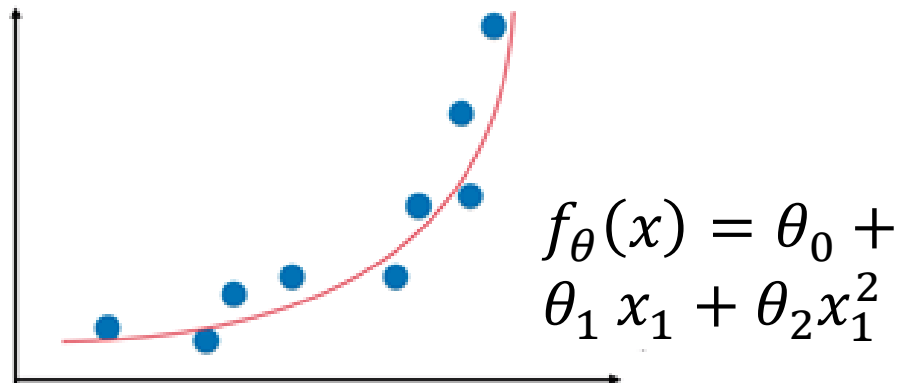  divides the learning rate by an exponentially decaying average of squared gradients

- **Adam**

  （Adaptive Moment Estimation）

  stores an exponentially decaying average of past squared gradients like RMSprop, also keeps an exponentially decaying average of past gradients, similar to momentum

# 多项式回归
# Polynomial regression



$$f_\theta(x) = \theta_0 + \theta_1 x_1$$

$$f_\theta(x) = \theta_0 + \theta_1 x_1 + \theta_2 x_1^2$$

# 多项式回归
# Polynomial regression



$$f_\theta(x) = \theta_0 + \theta_1 x_1$$

$$f_\theta(x) = \theta_0 + \theta_1 x_1 + \theta_2 x_1^2$$

# 多项式回归
# Polynomial regression

$$f_\theta(x) = \theta_0 + \theta_1 x_1$$

$$f_\theta(x) = \theta_0 + \theta_1 x_1 + \theta_2 x_1^2$$

$$f_\theta(x) = \theta_0 + \theta_1 x_1 + \theta_2 x_2$$

Polynomial feature

# 多项式回归
# Polynomial regression



$$f_\theta(x) = \theta_0 + \theta_1 x_1$$

$$f_\theta(x) = \theta_0 + \theta_1 x_1 + \theta_2 x_1^2$$

$$f_\theta(x) = \theta_0 + \theta_1 x_1 + \theta_2 x_1^2 + \theta_2 x_1^3 + \theta_2 x_1^4 + \theta_2 x_1^5 + \cdots$$

⬇

$$f_\theta(x) = \theta_0 + \theta_1 x_1 + \theta_2 x_2 + \theta_2 x_3 + \theta_2 x_4 + \theta_2 x_5 + \dots$$

# 多项式回归
# Polynomial regression



$$f_\theta(x) = \theta_0 + \theta_1 x_1$$

$$f_\theta(x) = \theta_0 + \theta_1 x_1 + \theta_2 x_1^2$$

$$f_\theta(x) = \theta_0 + \theta_1 x_1 + \theta_2 x_1^2 + \theta_2 x_1^3 + \theta_2 x_1^4 + \theta_2 x_1^5 + \cdots$$

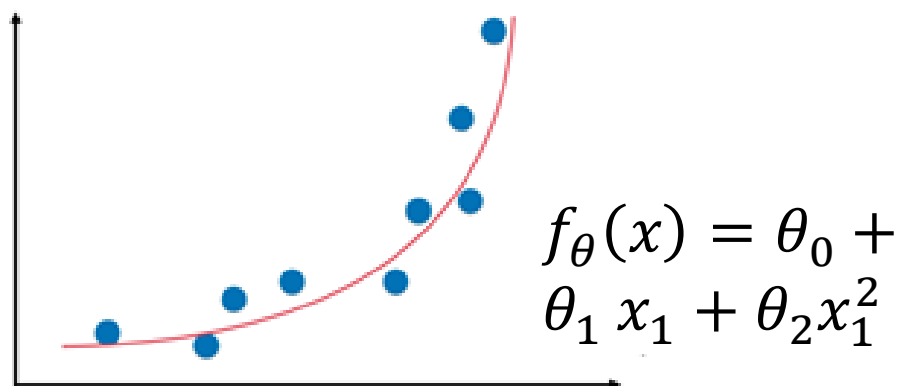$$f_\theta(x) = \theta_0 + \theta_1 x_1 + \theta_2 x_2 + \theta_2 x_3 + \theta_2 x_4 + \theta_2 x_5 + \ldots$$

✓ able to model all sorts of relationships
X easy to overfit

# 最小二乘方线性回归
# Least square linear regression



$J(\theta)$

# 最小二乘方线性回归
# Least square linear regression

## Gradient Descent



$J(\theta)$

## Normal equation



$J(\theta)$

$\theta$

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

(for every $j$)

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

# 最小二乘法求解
# Least square method

| $x_0$ | 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 |

$$X = \begin{bmatrix} 1 & 2104 & 5 & 1 & 45 \\ 1 & 1416 & 3 & 2 & 40 \\ 1 & 1534 & 3 & 2 & 30 \\ 1 & 852 & 2 & 1 & 36 \end{bmatrix} \qquad y = \begin{bmatrix} 460 \\ 232 \\ 315 \\ 178 \end{bmatrix} \qquad \theta = \begin{bmatrix} \theta_0 \\ \theta_1 \\ \theta_2 \\ \theta_3 \\ \theta_4 \end{bmatrix}$$

$$J(\theta) = \frac{1}{2N} \sum_{i=1}^{N} \left( f_\theta(x^{(i)}) - y^{(i)} \right)^2 = \frac{1}{2}(X\theta - y)^2 = \frac{1}{2}(X\theta - y)^T(X\theta - y)$$

# 最小二乘法
## Least square method

- Objective $\min\limits_{\theta} J(\theta)$

$$J(\theta) = \frac{1}{2N} \sum_{i=1}^{N} \left(f_\theta\left(x^{(i)}\right) - y^{(i)}\right)^2 = \frac{1}{2}(X\theta - y)^2 = \frac{1}{2}(X\theta - y)^T(X\theta - y)$$

- Gradient

$$\frac{\partial J(\theta)}{\partial \theta} = \frac{1}{2}\frac{\partial}{\partial \theta}\left(\theta^T X^T X\theta - y^T X\theta - \theta^T X^T y + y^T y\right)$$

$$= \frac{1}{2}\frac{\partial}{\partial \theta}\left(\theta^T X^T X\theta - 2\theta^T X^T y + y^T y\right)$$

$$= X^T X\theta - X^T y$$

- Solution

$$\frac{\partial J(\theta)}{\partial \theta} = 0 \Rightarrow X^T X\theta - X^T y = 0 \Rightarrow \theta = (X^T X)^{-1} X^T y$$

# 正规方程求解
# Normal equation method

| $x_0$ | Size (feet$^2$) $x_1$ | 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 |

$$X = \begin{bmatrix} 1 & 2104 & 5 & 1 & 45 \\ 1 & 1416 & 3 & 2 & 40 \\ 1 & 1534 & 3 & 2 & 30 \\ 1 & 852 & 2 & 1 & 36 \end{bmatrix} \qquad y = \begin{bmatrix} 460 \\ 232 \\ 315 \\ 178 \end{bmatrix} \qquad \theta = \begin{bmatrix} \theta_0 \\ \theta_1 \\ \theta_2 \\ \theta_3 \\ \theta_4 \end{bmatrix}$$

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

# 最小二乗法
## Least square method

| $x_0$ | Size (feet$^2$) $x_1$ | 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 |

$$X = \begin{bmatrix} 1 & 2104 & 5 & 1 & 45 \\ 1 & 1416 & 3 & 2 & 40 \\ 1 & 1534 & 3 & 2 & 30 \\ 1 & 852 & 2 & 1 & 36 \end{bmatrix} \qquad y = \begin{bmatrix} 460 \\ 232 \\ 315 \\ 178 \end{bmatrix} \qquad \theta = \begin{bmatrix} \theta_0 \\ \theta_1 \\ \theta_2 \\ \theta_3 \\ \theta_4 \end{bmatrix}$$

$$\theta = (X^T X)^{-1} X^T y \quad \leftarrow \quad O(\text{number of features})^3$$

# 梯度下降 VS 最小二乘法
# Gradient descent VS Least square method

Gradient Descent

- Need to choose α .
- Needs many iterations.

- Works well even when the number of features is large.

Least square method

- No need to choose α .
- Don't need to iterate.

- Need to compute $(X^T X)^{-1}$
- Slow if the number of features is very large （>10000） .
- only applicable to linear models
- Sometimes cannot be directly calculated(if $X^T X$ is non-invertible).

# 评价标准
## Evaluation indices

MSE(Mean Squared Error)
均方误差

$$\frac{1}{N}\sum_{i=1}^{N}(y^{(i)} - f(x^{(i)}))^2$$

RMSE（Root Mean Squared Error）
均方根误差

$$\sqrt{\frac{1}{N}\sum_{i=1}^{N}(y^{(i)} - f(x^{(i)}))^2}$$

MAE (Mean absolute Error)
平均绝对误差

$$\frac{1}{N}\sum_{i=1}^{N}|(y^{(i)} - f(x^{(i)}))^2|$$

R-Squared（r2score）R方/决定系数

$$= 1 - \frac{\sum_{i=1}^{N}(y^{(i)} - f(x^{(i)}))^2}{\sum_{i=1}^{N}(y^{(i)} - \bar{y})^2}$$

$$= 1 - \frac{MSE}{Var}$$

# 思考题

多变量线性回归相比单变量回归，采用标准的梯度下降求解会有什么问题及可能的解决方法?