

Machine Learning



Amin Golzari Oskouei

a.golzari@azaruniv.ac.ir

a.golzari@tabrizu.ac.ir

<https://github.com/Amin-Golzari-Oskouei>

Azərbaycan Şahid Mədani Universiteti
2023

Supervised learning: Regression



Table of Contents

3

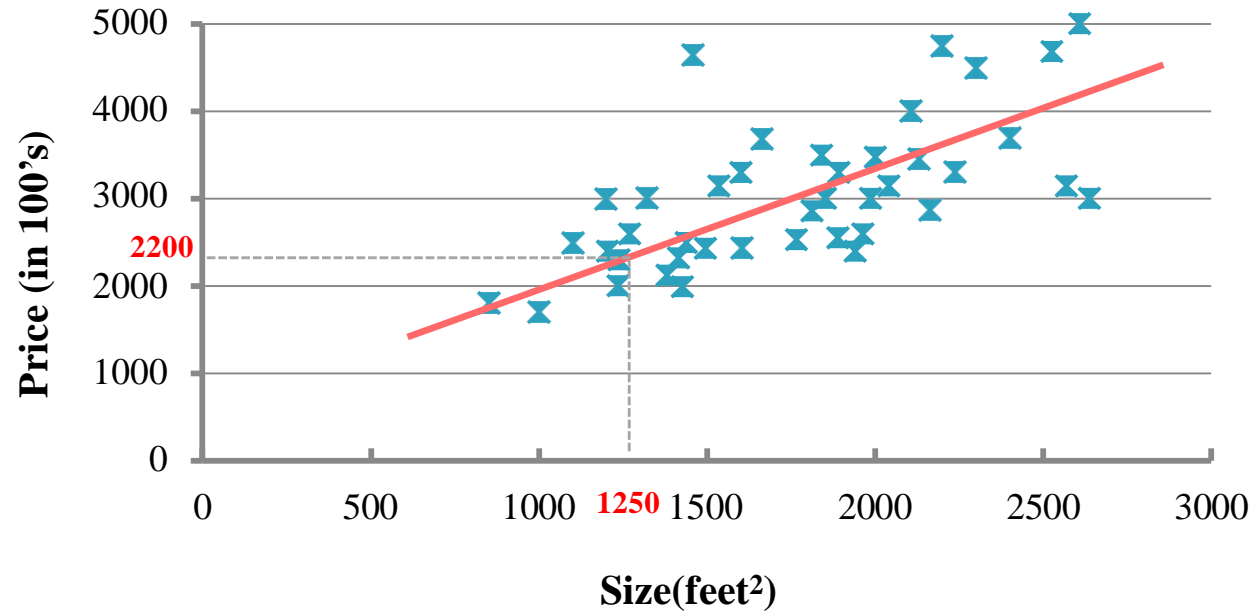
- Regression
 - Univariate and multivariate linear regression
- Gradient descent
- Normal equation
- Regression with local weighting
- Probabilistic interpretation of regression
- Maximum likelihood estimation.

A horizontal decorative bar at the top of the slide, consisting of a red rectangular section on the left and a teal rectangular section on the right.

Univariate linear regression

Home pricing

5



Supervised learning

For each Train set example, the correct answer is given (having Lable)

Regression

Predicting quantities with continuous values (such as the price of a house)

symbolization

6

Train Set

$m = 47$

Area (square feet) (x)	Price (in 1000 dollars) (y)
2104	460
1416	232
1534	315
852	178
....

□ Symbols

m = Number of training samples

x = Input variable, features

y = Output variable, target variable

(x,y) : an trainset sample

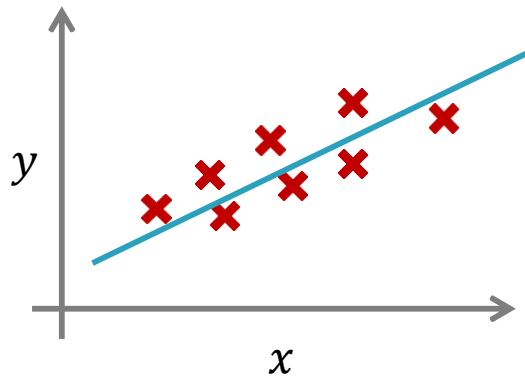
(x^i, y^i) : i -th sample of train set

Model representation

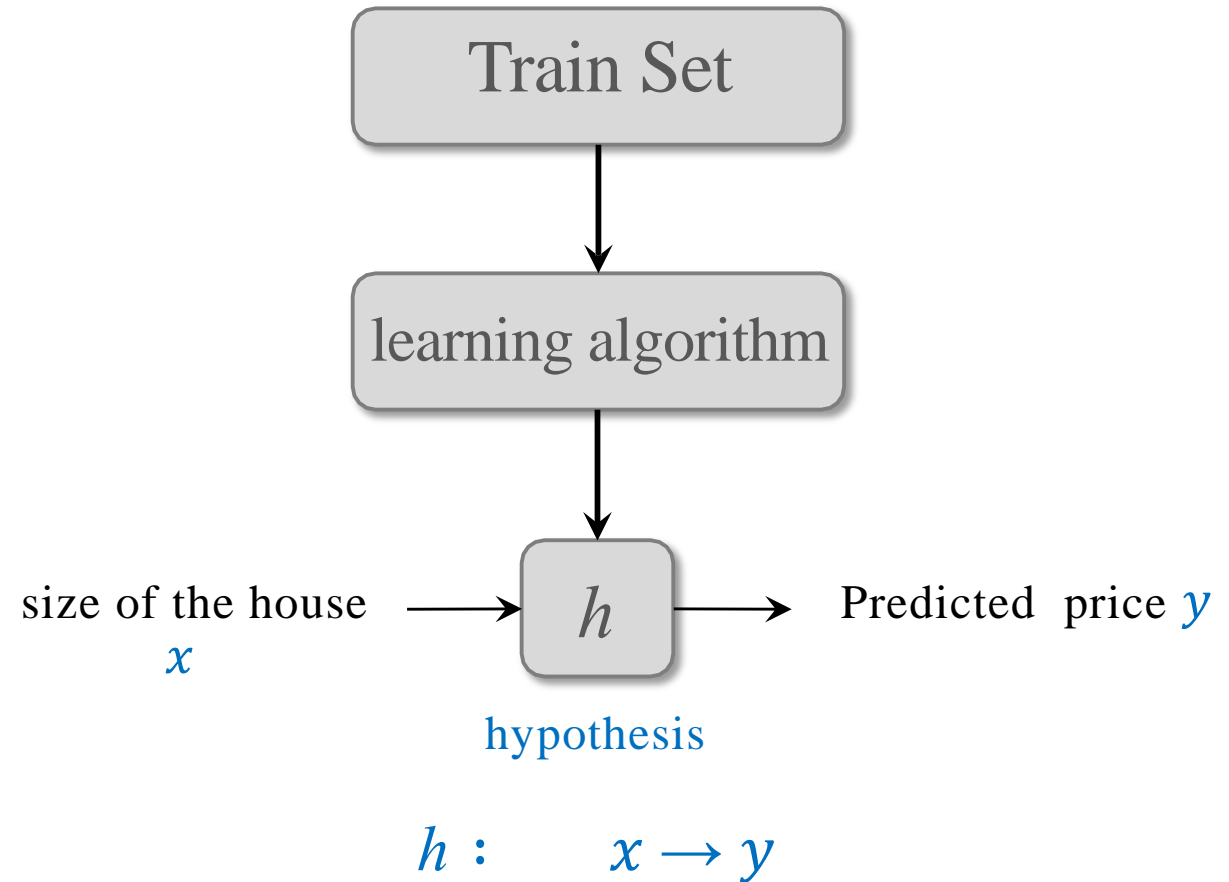
7

Show hypothesis h

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



Univariate linear regression





Cost function

Hypothesis evaluation

9

Train Set

$m = 47$

Area (square feet) (x)	Price (in 1000 dollars) (y)
2104	460
1416	232
1534	315
852	178
....
$h_{\theta}(x) = \theta_0 + \theta_1 x$	
(θ_0, θ_1)	

Hypothesis :

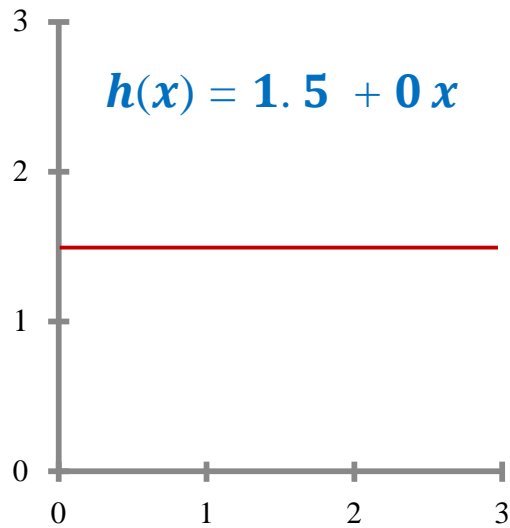
Parameters :

Question: How to choose the value of the parameters?

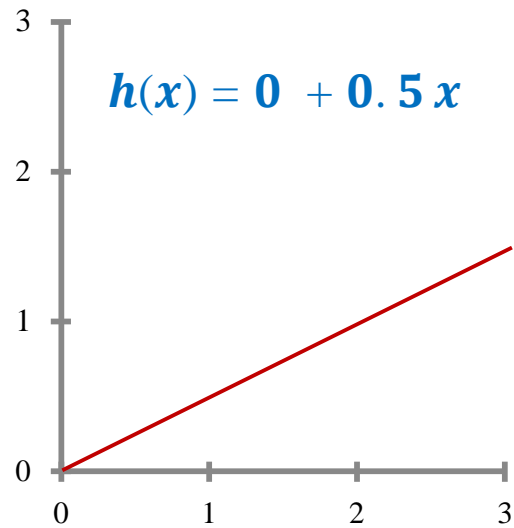
Hypothesis evaluation

10

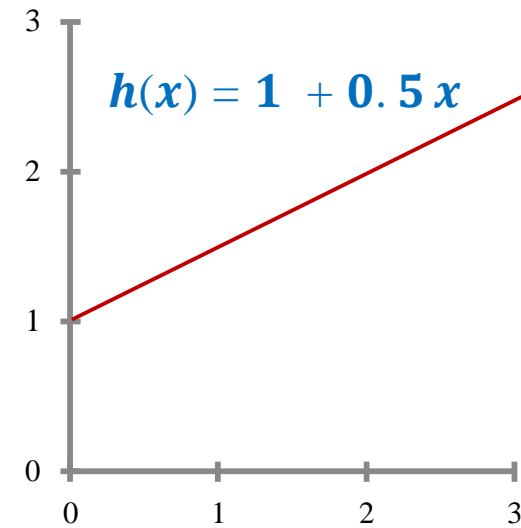
$$h_{\theta}(x) = \theta_0 + \theta_1 x$$



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



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



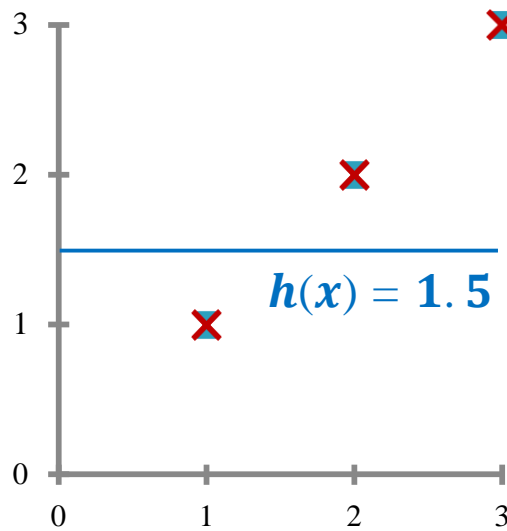
$$\theta_0 = 1$$
$$\theta_1 = 0.5$$

Hypothesis evaluation

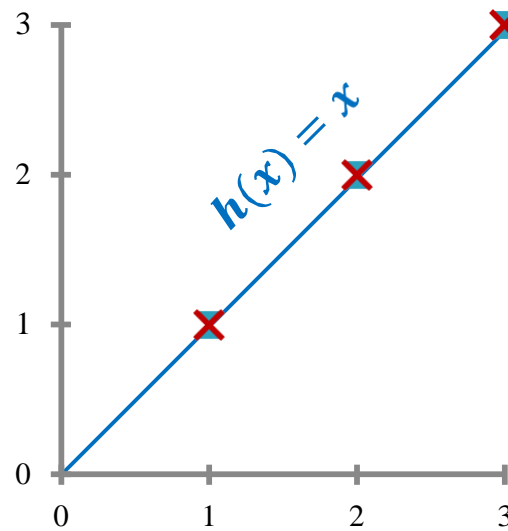
11

Question: Which hypothesis is better?

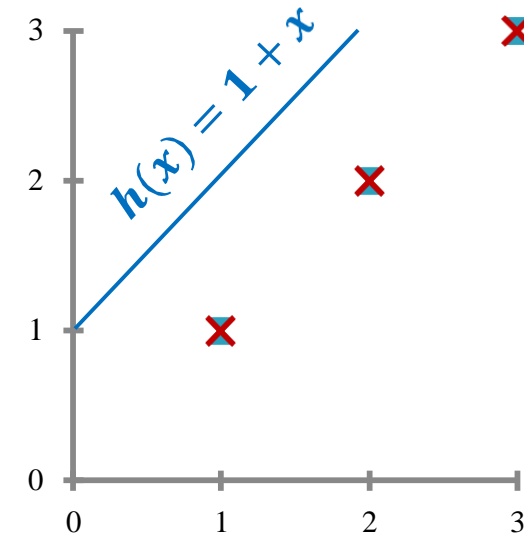
$$h_{\theta}(x) = \theta_0 + \theta_1 x$$



$$\begin{aligned}\theta_0 &= 1.5 \\ \theta_1 &= 0\end{aligned}$$



$$\begin{aligned}\theta_0 &= 0 \\ \theta_1 &= 1\end{aligned}$$



$$\begin{aligned}\theta_0 &= 1 \\ \theta_1 &= 1\end{aligned}$$

Cost function

12

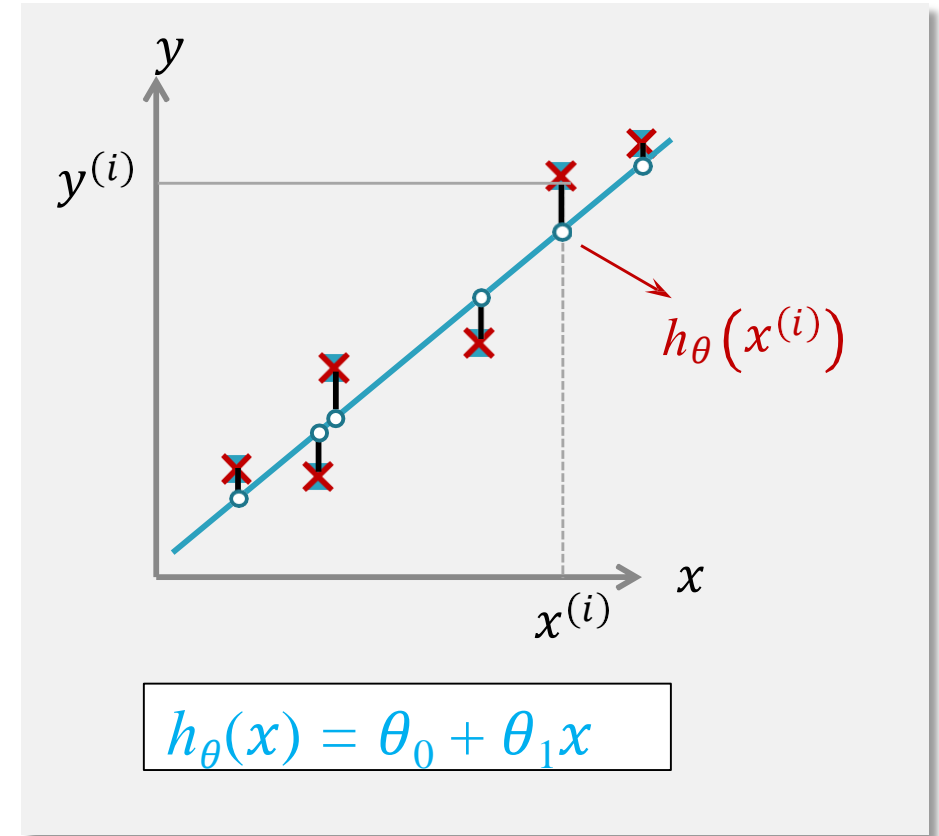
□ **Ideas.** Choosing the parameters so that for each training sample such as x, y , the value of $h_{\theta}(x)$ is as close as possible to the value of y .

□ **cost function**, sum of squared error

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

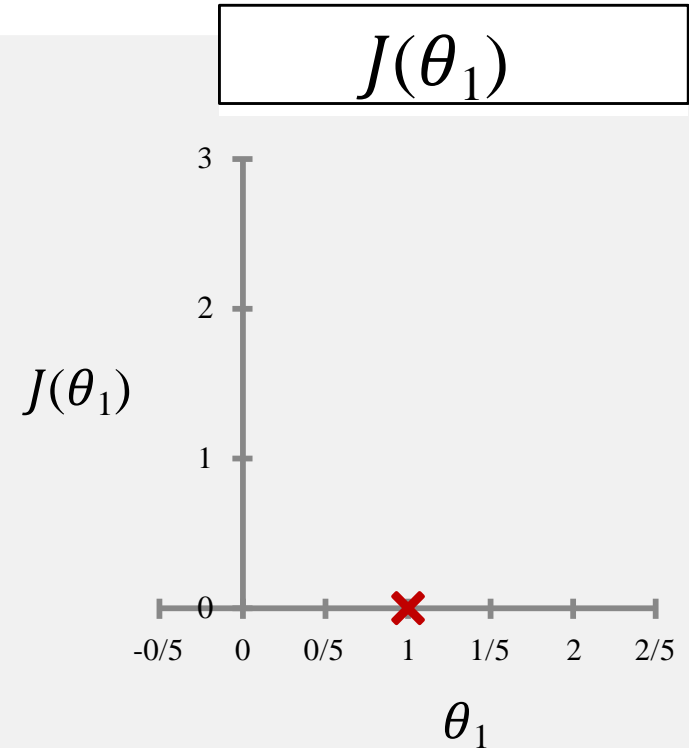
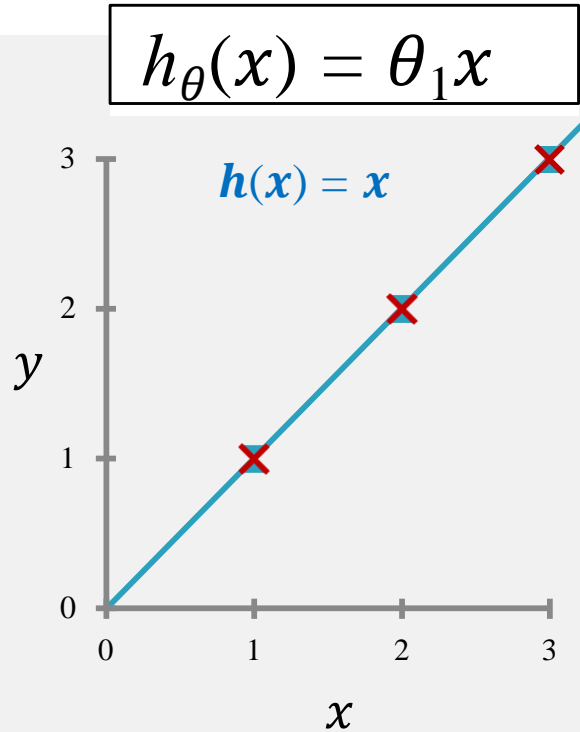
□ **purpose**

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



Simplified cost function

13

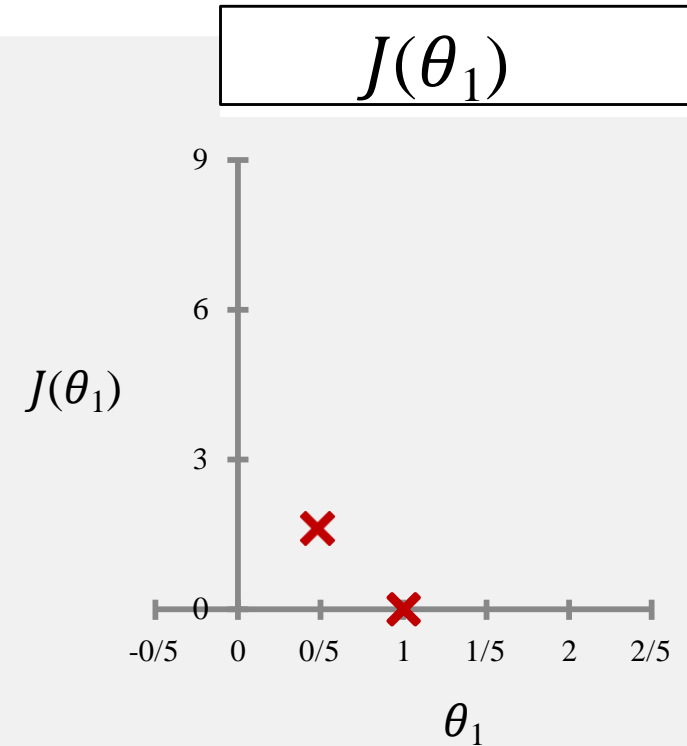
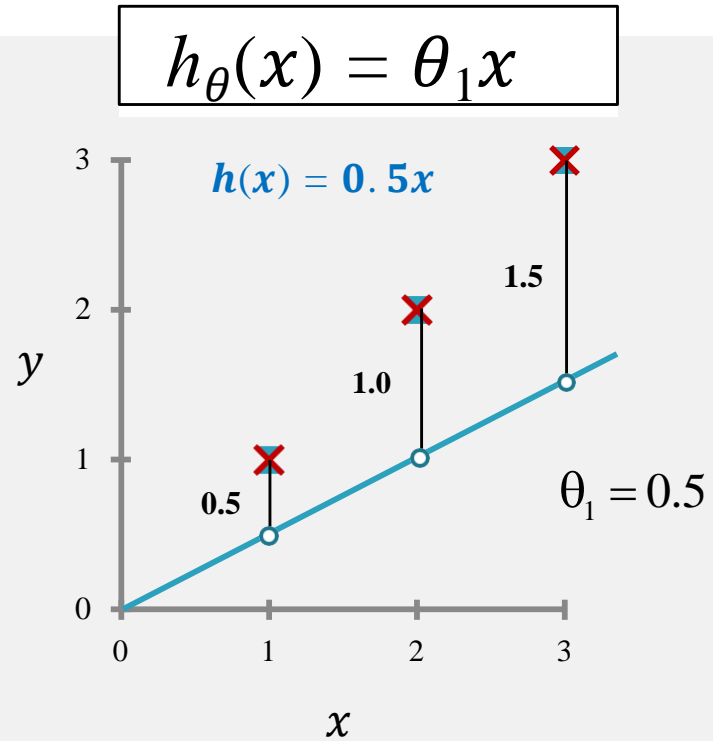


$$\begin{aligned} J(\theta_0, \theta_1) &= \frac{1}{2} \sum_{i=1}^m (h_{\theta}(x^{(i)}) - y^{(i)})^2 \\ &= \frac{1}{2} \sum_{i=1}^m (x^{(i)} - y^{(i)})^2 = \frac{1}{2} (0^2 + 0^2 + 0^2) = 0 \end{aligned}$$

$$J(1) = 0$$

Cost function

14

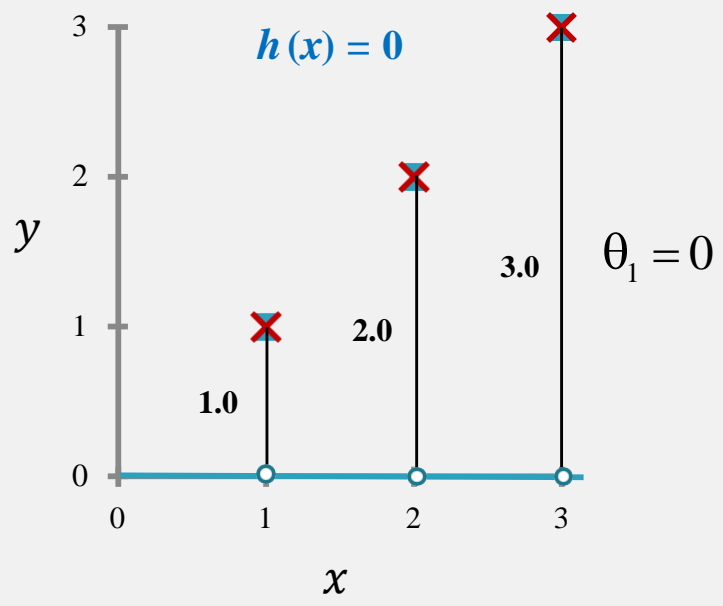


$$J(0.5) = \frac{1}{2} (0.5^2 + 1^2 + 1.5^2) = \frac{1}{2} (3.5) = 1.75$$

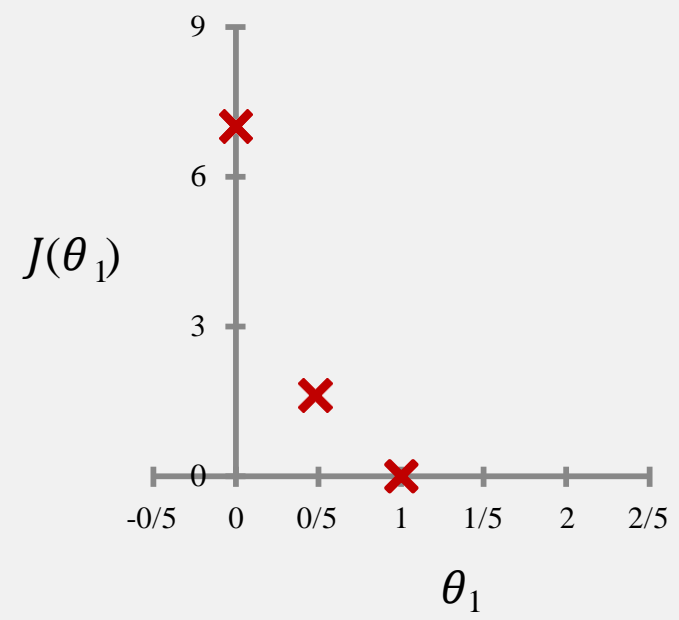
$$J(0.5) = 1.75$$

Cost function

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



$$J(\theta_1)$$



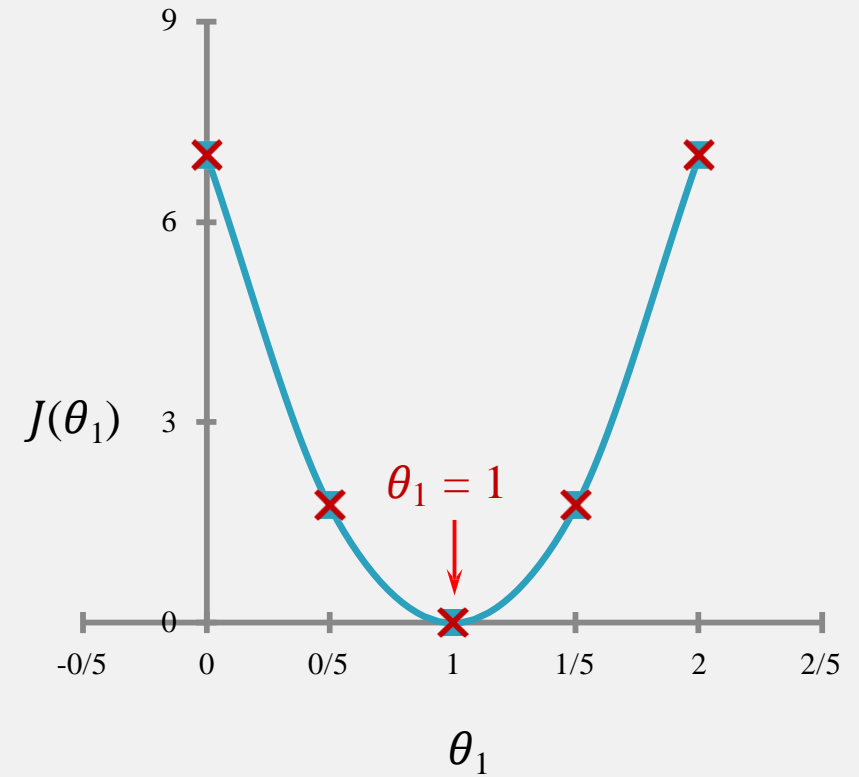
$$J(0) = \frac{1}{2} (1^2 + 2^2 + 3^2) = \frac{1}{2} (14) = 7.0$$

$$J(0) = 7.0$$

Cost function

16

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



Univariate linear regression

17

□ hypothesis

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

□ parameters

$$\theta_0, \theta_1$$

□ Cost functions

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

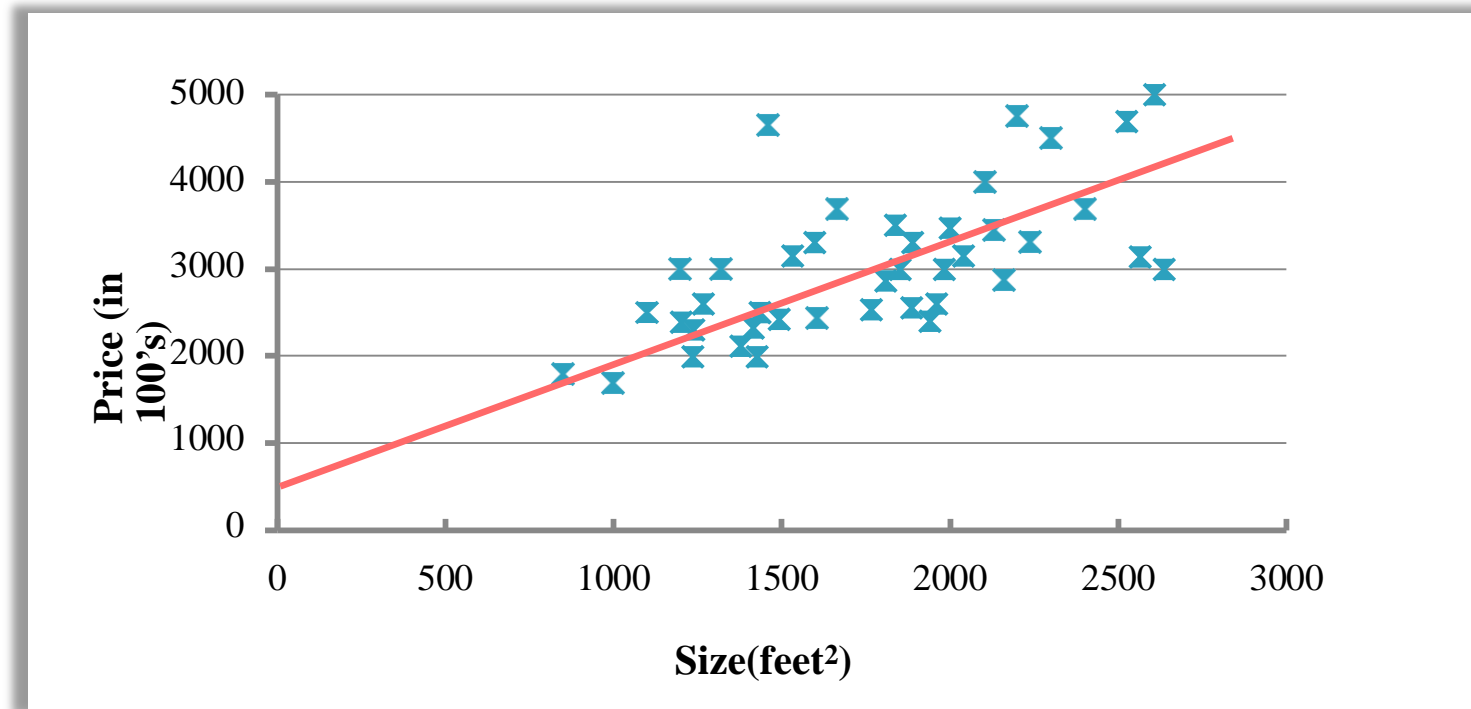
□ purpose

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

Example: home pricing

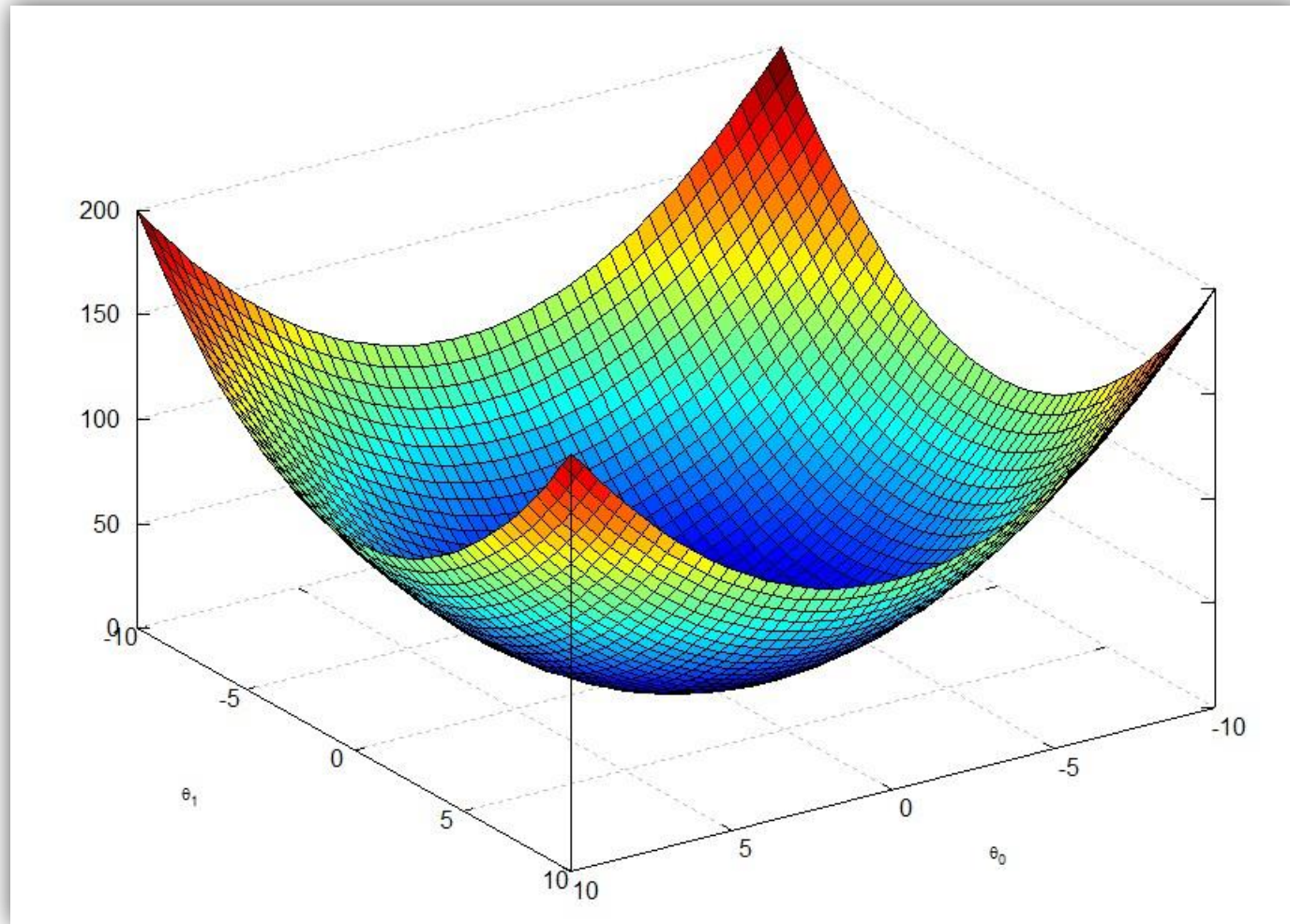
18

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



Cost function

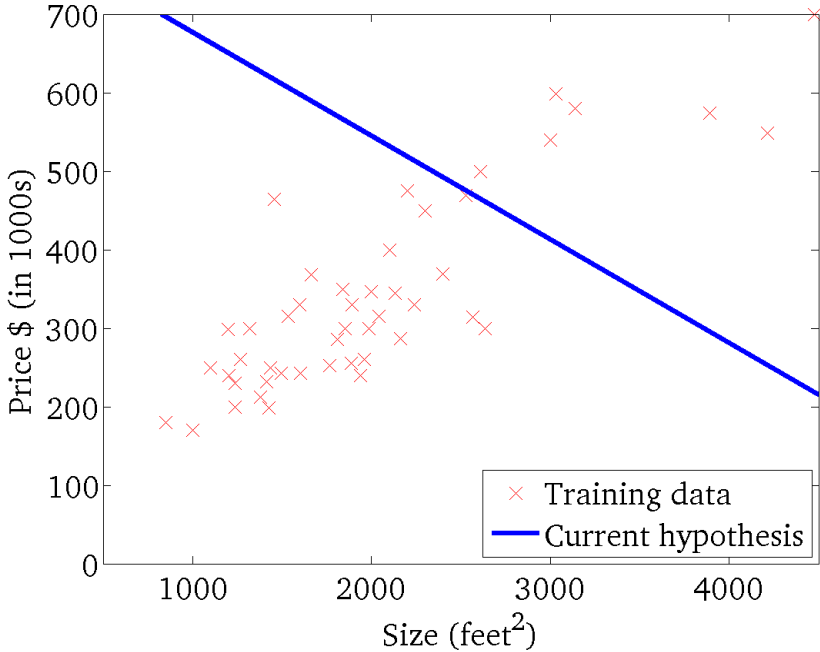
19



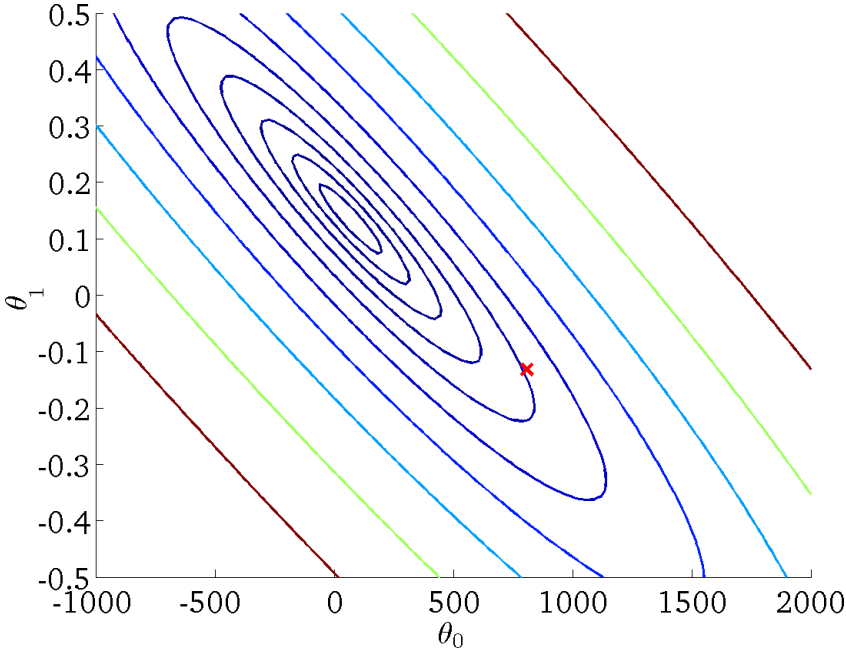
```
x = -10:0.5:10;  
y = -10:0.5:10;  
[X Y] = meshgrid(x,y);  
Z = X.^2 + Y.^2;  
surf(X,Y,Z);
```

Cost function: contour curve

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

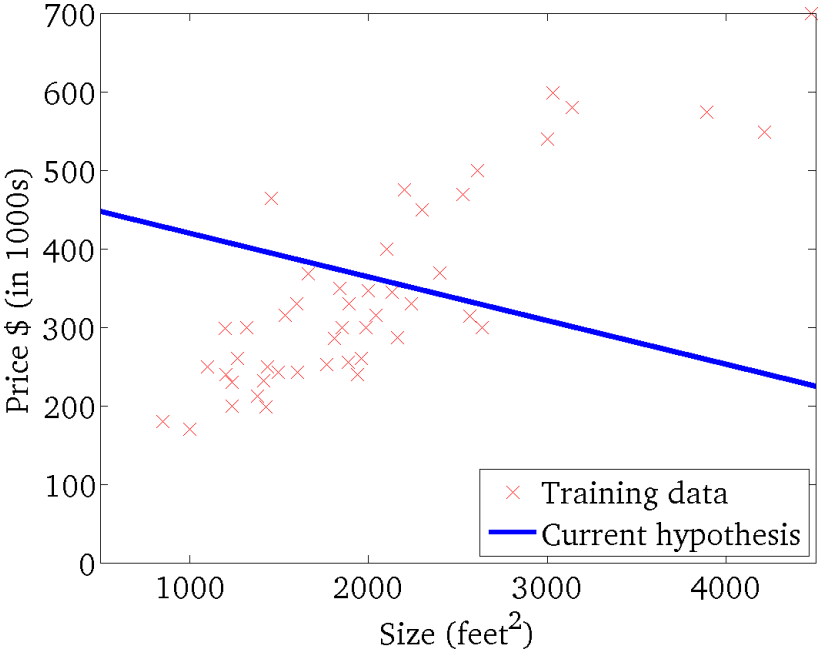


$$J(\theta_0, \theta_1)$$

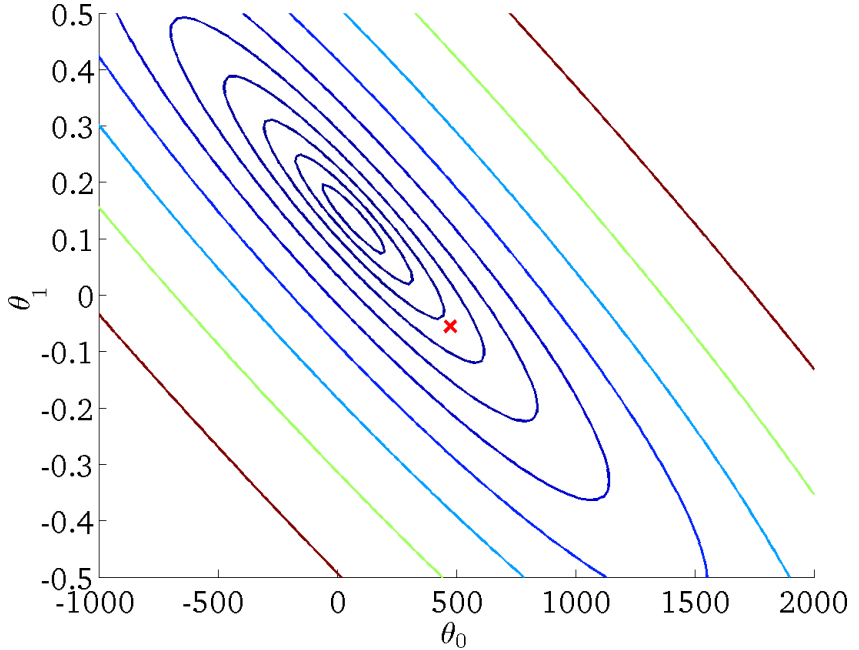


Cost function: contour curve

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

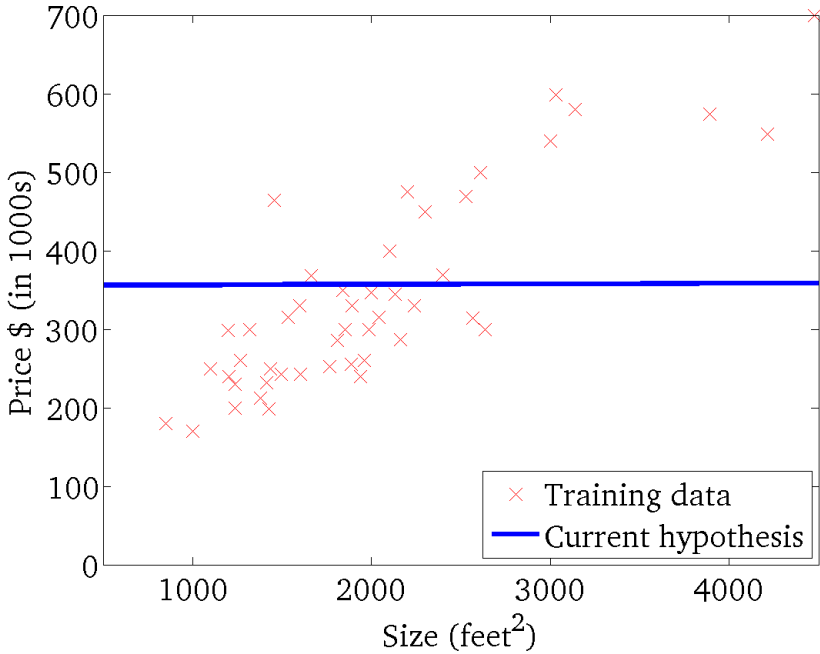


$$J(\theta_0, \theta_1)$$

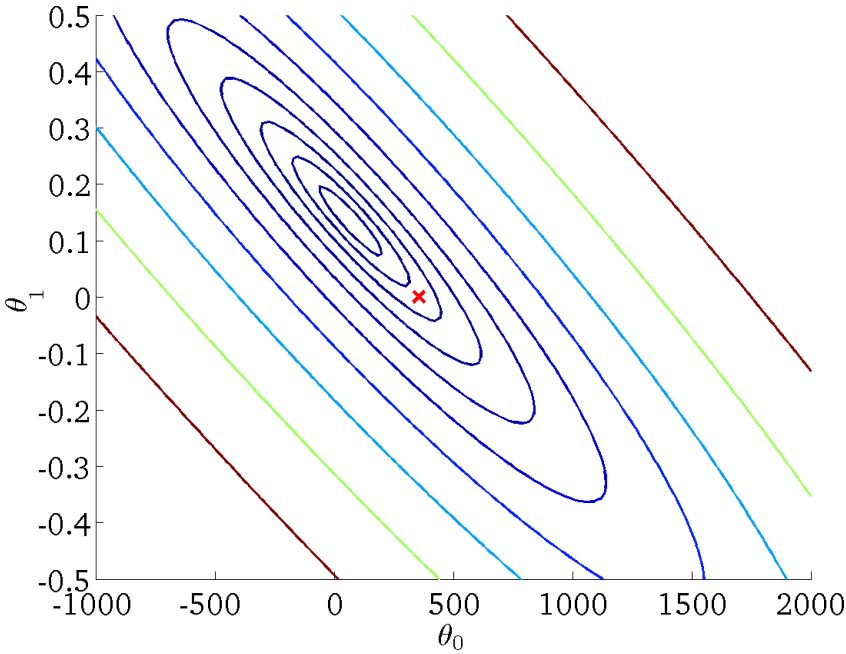


Cost function: contour curve

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

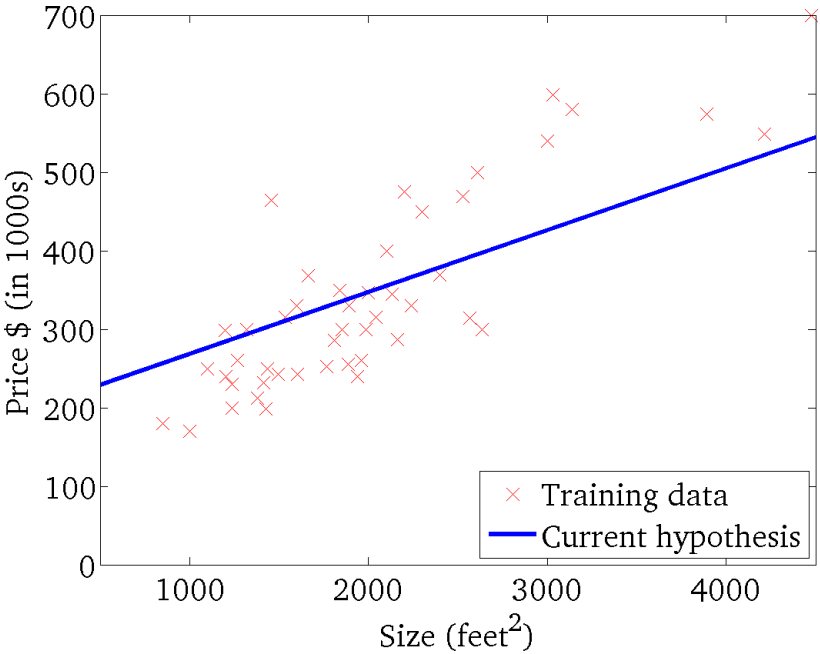


$$J(\theta_0, \theta_1)$$

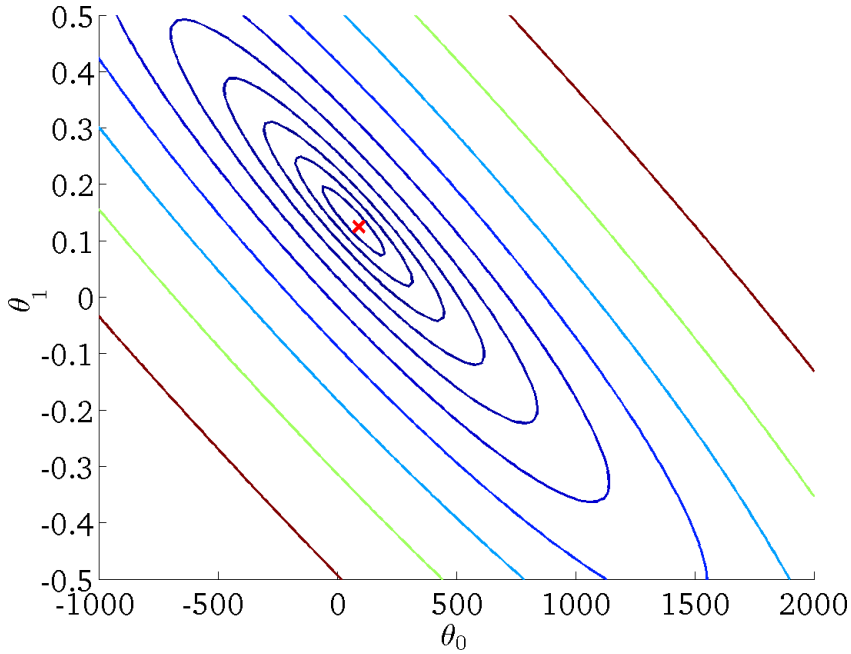


Cost function: contour curve

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



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



A horizontal bar at the top of the slide, divided into a red section on the left and a teal section on the right. The text "Gradient Descent" is written in white serif font on the teal section.

Gradient Descent

Question

25

□ hypothesis

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

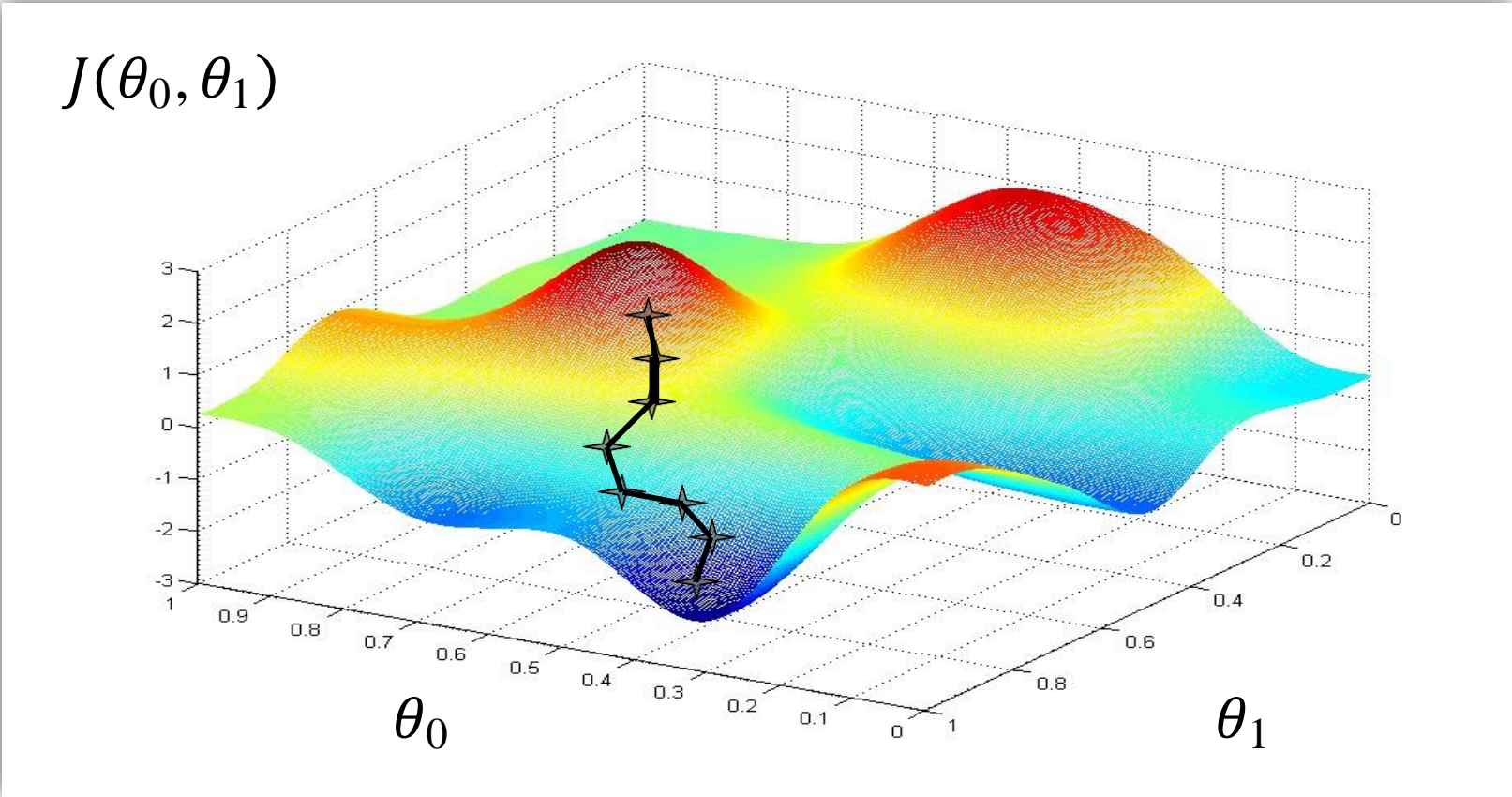
□ purpose

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

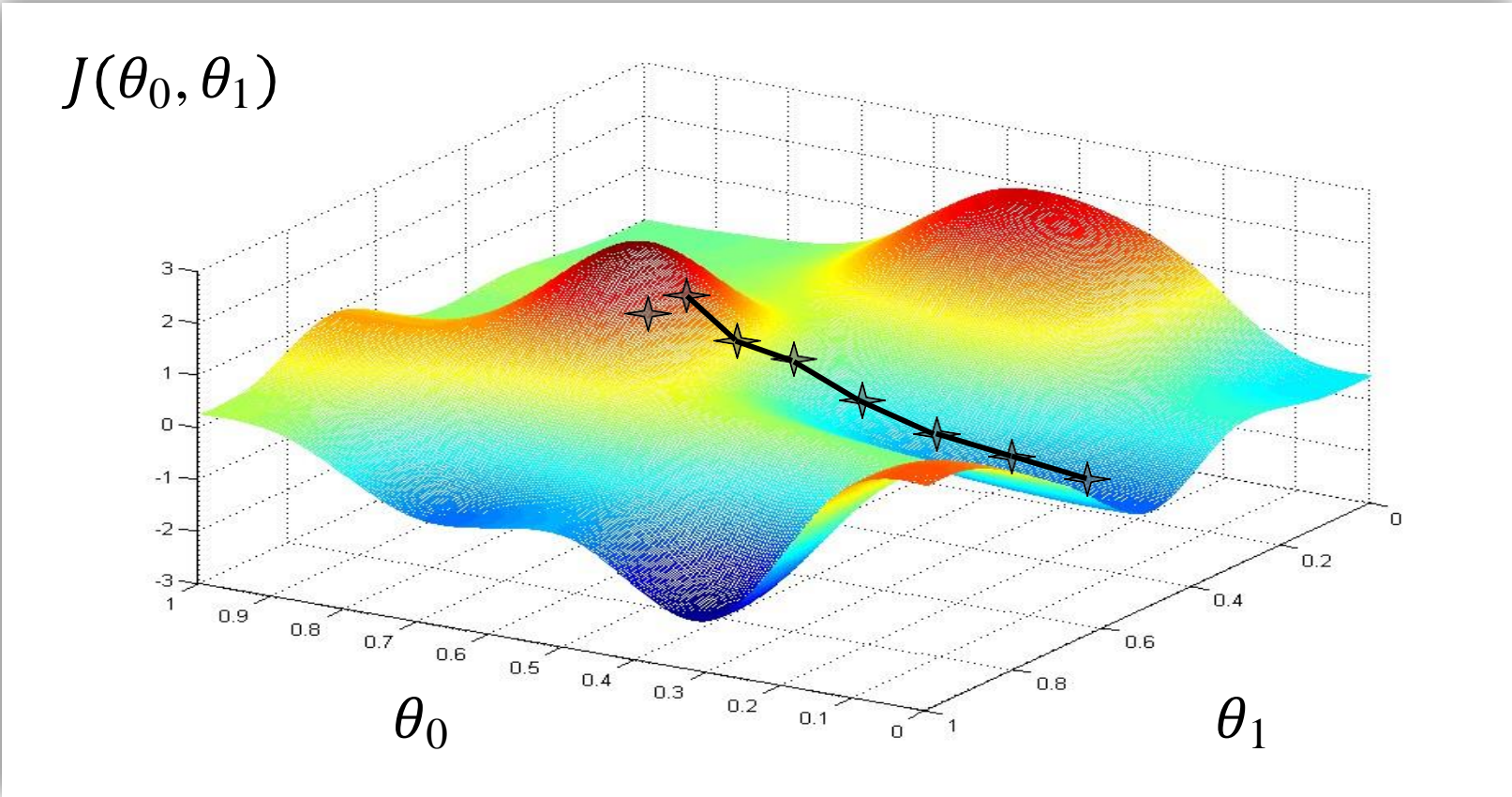
□ General method

- ❖ Start with a random initial value for the parameters θ_0 and θ_1 . [eg zero value]
- ❖ Change the value of the parameters in such a way that the value of the cost function $J(\theta_0, \theta_1)$ decreases.
- ❖ Repeat the above operation until we reach a minimum value for the cost function. [convergence]

Gradient Descent: Global Optimum



Gradient Descent: Local Optimum



Gradient descent algorithm

28

repeat until convergence {
 $\theta_j := \theta_j - \alpha \frac{\partial}{\partial \theta_j} J(\theta_0, \theta_1)$ (for $j = 0$ and $j = 1$)
Learning rate

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

Correct implementation. Update parameters value
simultaneously

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

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

$$\theta_0 := \theta_0 + \Delta\theta_0$$

$$\theta_1 := \theta_1 + \Delta\theta_1$$

Gradient descent algorithm

29

repeat until convergence {
 $\theta_j := \theta_j - \alpha \frac{\partial}{\partial \theta_j} J(\theta_0, \theta_1)$ (for $j = 0$ and $j = 1$)
Learning rate

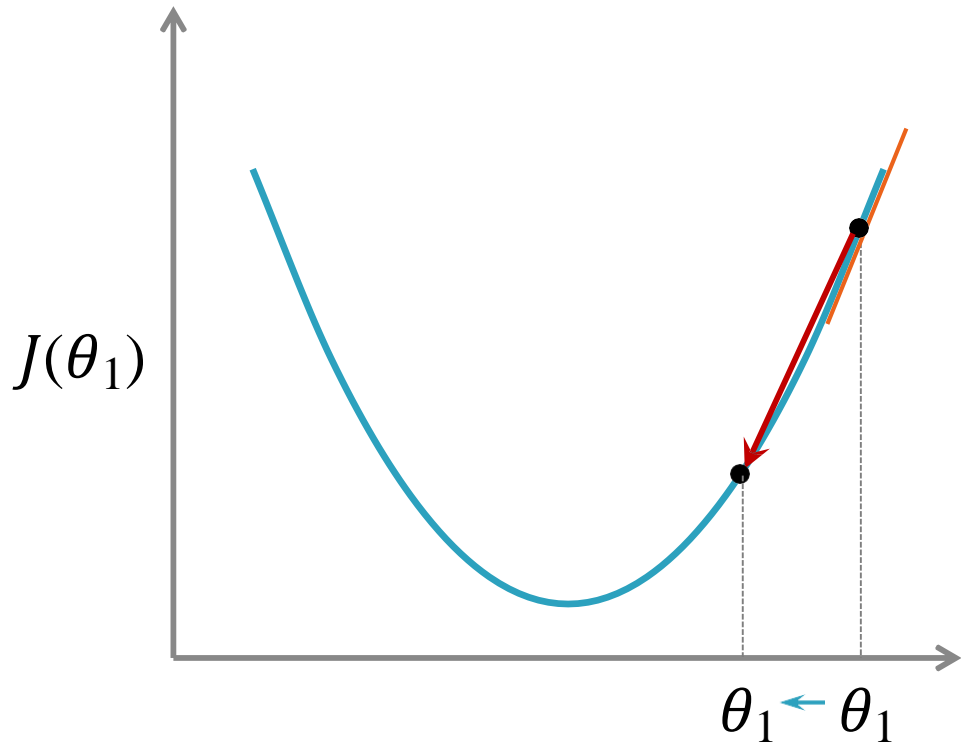
$$h_{\theta}(x) = \theta_0 + \theta_1 x$$

Incorrect implementation. Update parameters value sequentially

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

Gradient descent algorithm

30



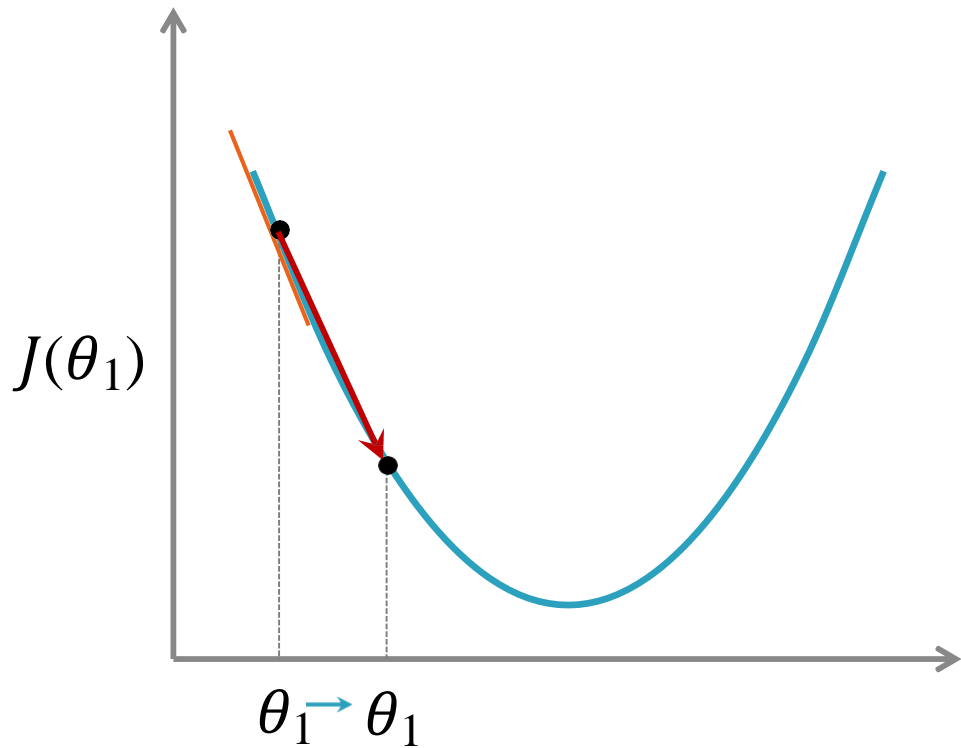
Positive slope

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

$$\theta_1 := \theta_1 - \underbrace{\alpha}_{\geq 0} (\geq 0)$$

Gradient descent algorithm

31



negative slope

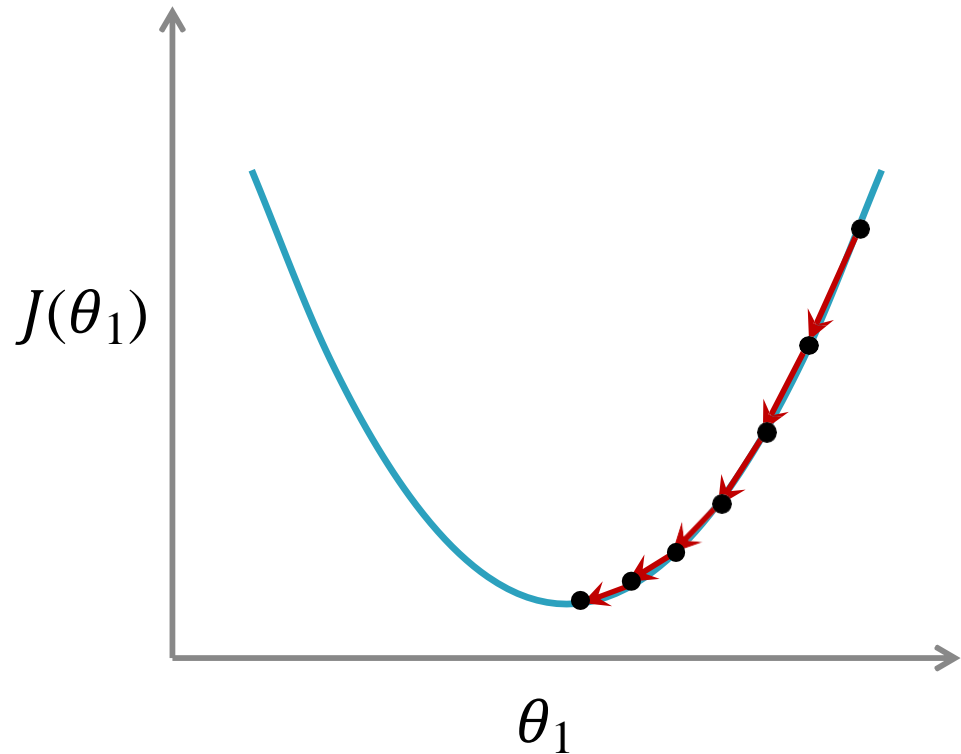
$$\theta_1 := \theta_1 - \alpha \frac{\partial}{\partial \theta_1} J(\theta)$$

$$\theta_1 := \theta_1 - \underbrace{\alpha}_{\leq 0} (\leq 0)$$

Gradient descent algorithm: learning rate

32

If the **learning rate** is too **small**, the gradient descent will converge slowly.

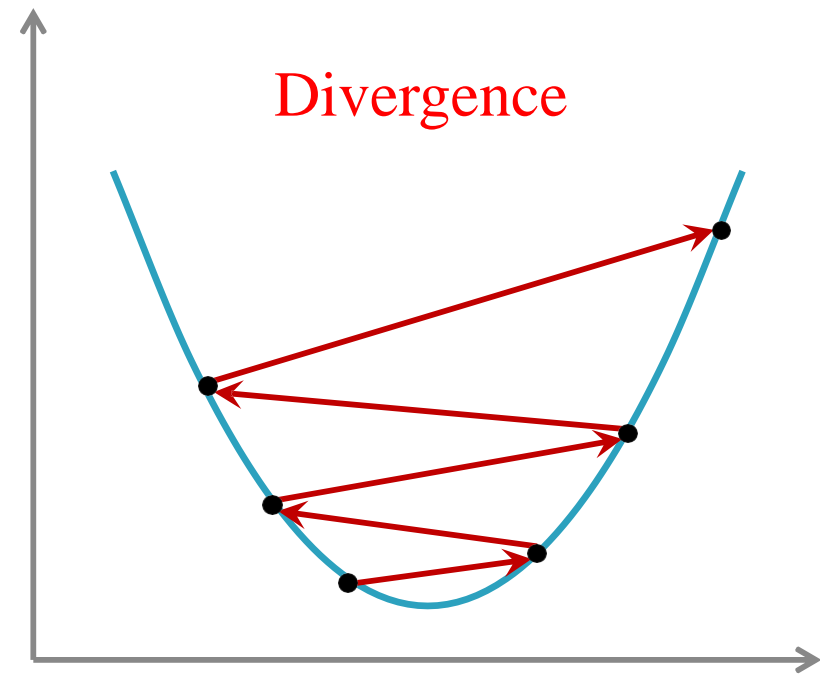
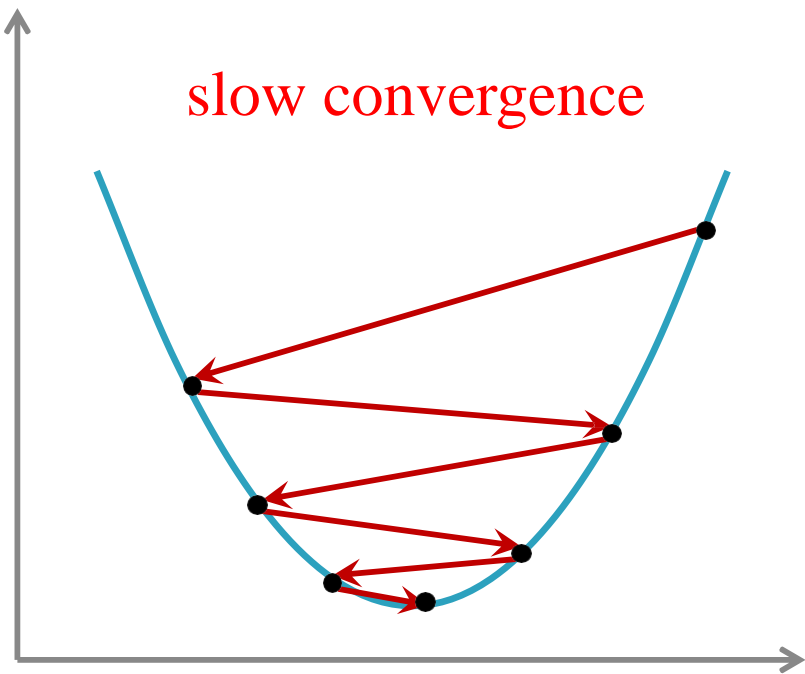


$$\theta_1 := \theta_1 - \alpha \frac{\partial}{\partial \theta_1} J(\theta_1)$$

Gradient descent algorithm: learning rate

33

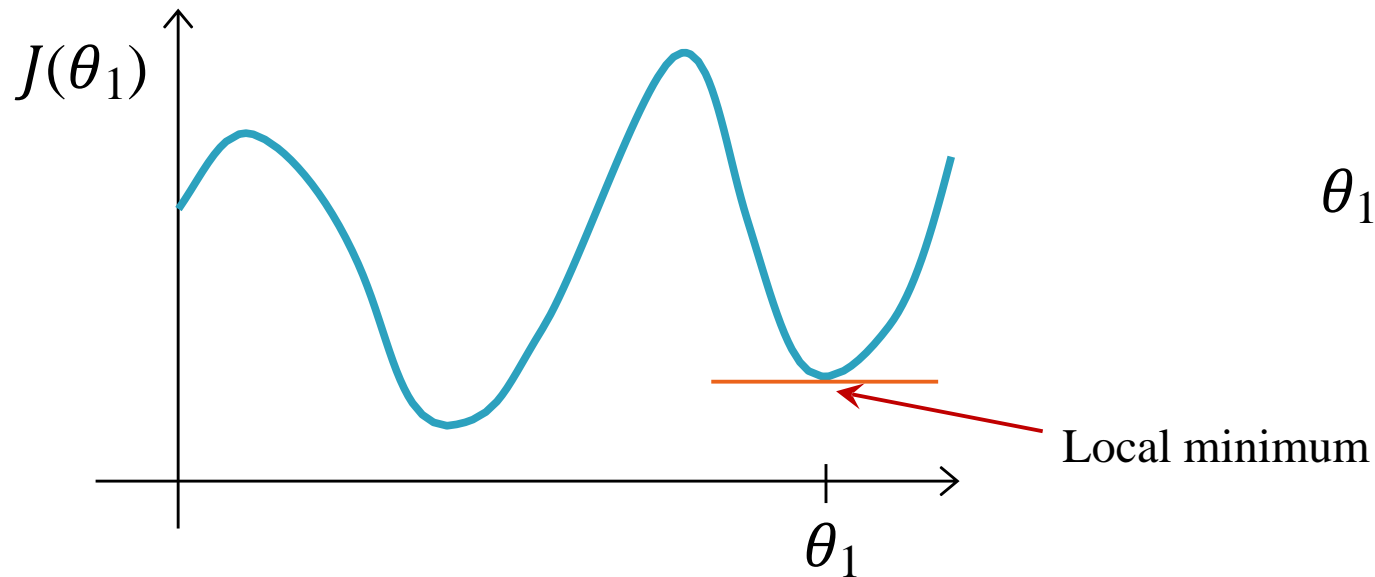
If the **learning rate** is too **large**, the gradient descent may converge slowly or even diverge.



Gradient Descent Algorithm: Convergence

34

Convergence When the value of the parameter θ_1 is in a **local minimum**.



$$\theta_1 := \theta_1 - \alpha \underbrace{\frac{\partial}{\partial \theta_1} J(\theta_1)}_0$$

Application of Gradient Descent in linear

regression

Gradient Descent and linear regression

36

repeat until convergence {

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

}

linear regression

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

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

The gradient of the cost function

37

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

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

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

Decreasing gradient and linear regression

38

repeat until convergence {

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

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

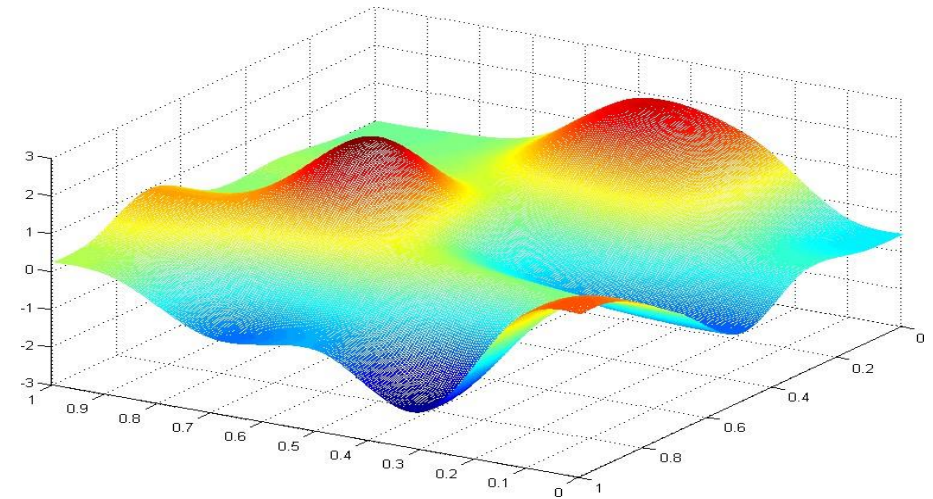
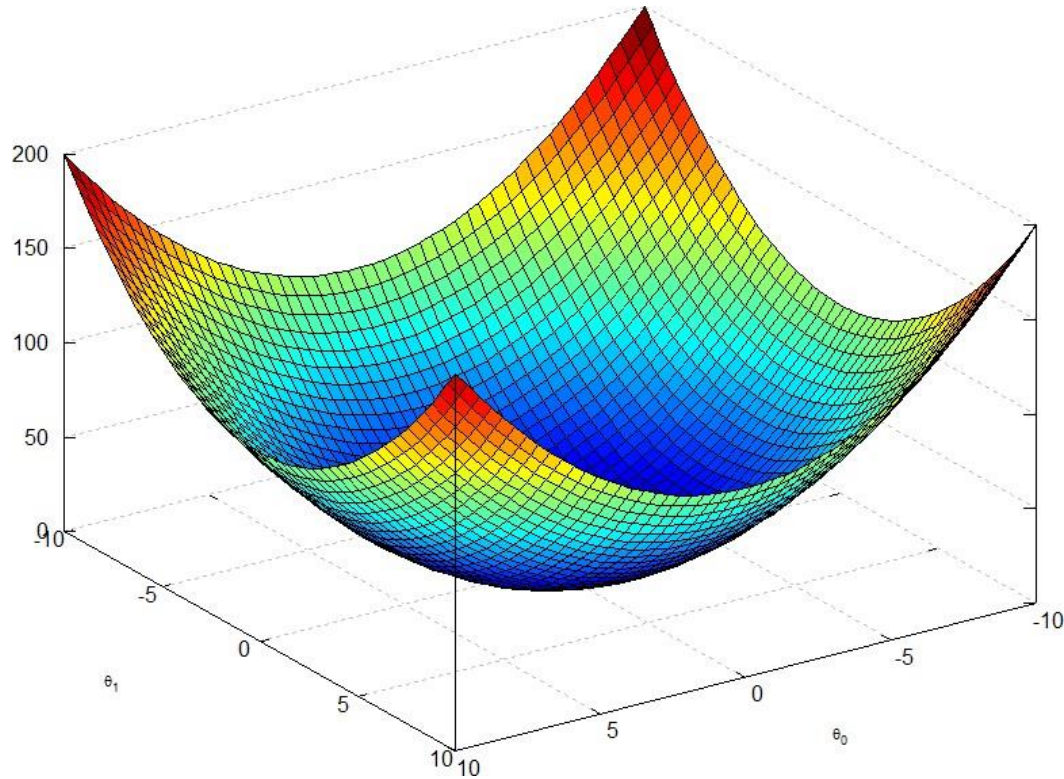
}

} Simultaneous update

Gradient Descent and linear regression

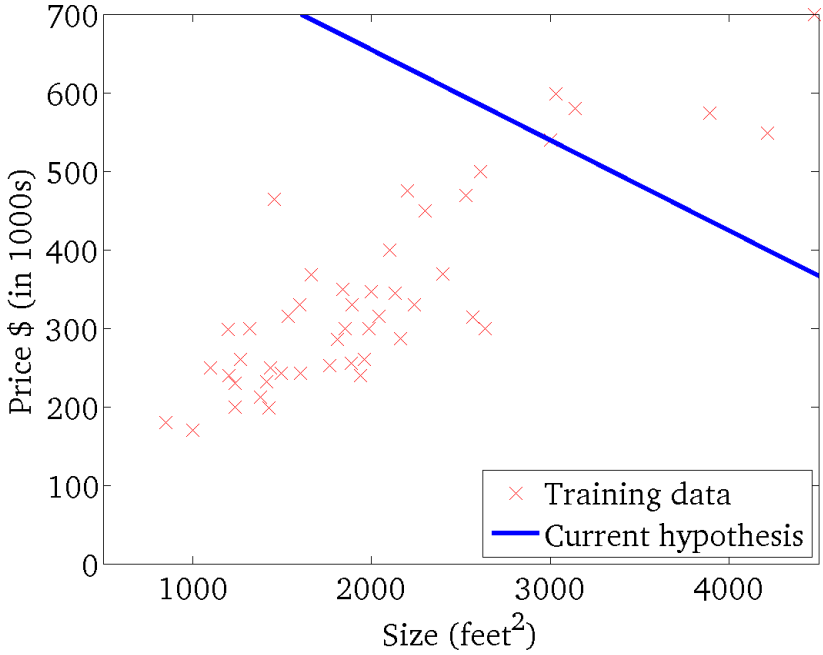
39

□ **Attention.** In linear regression, the cost function is a **cognate function**, and as a result, the decreasing gradient necessarily converges to the **global optimum** in case of convergence.

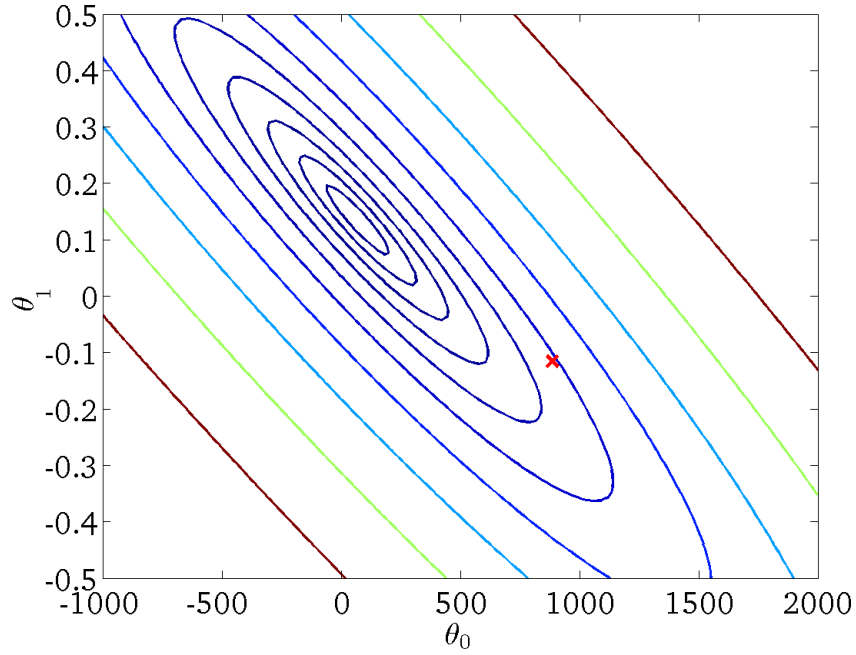


Gradient Descent

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

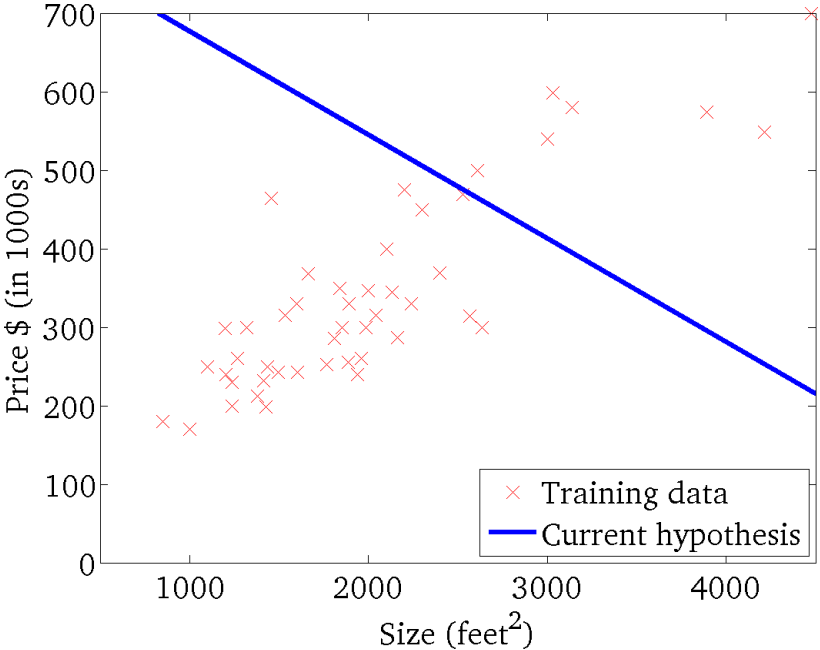


$$J(\theta_0, \theta_1)$$

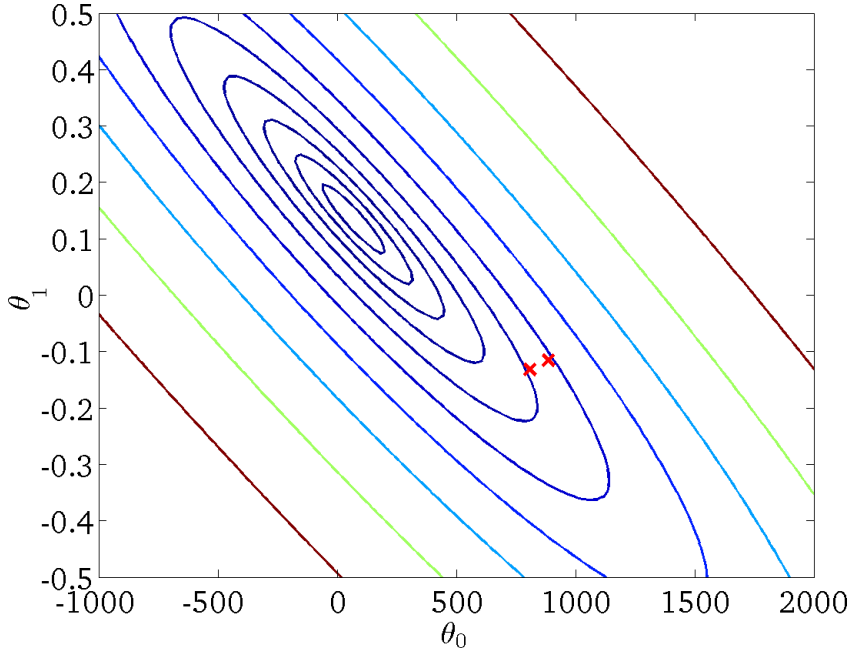


Gradient Descent

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

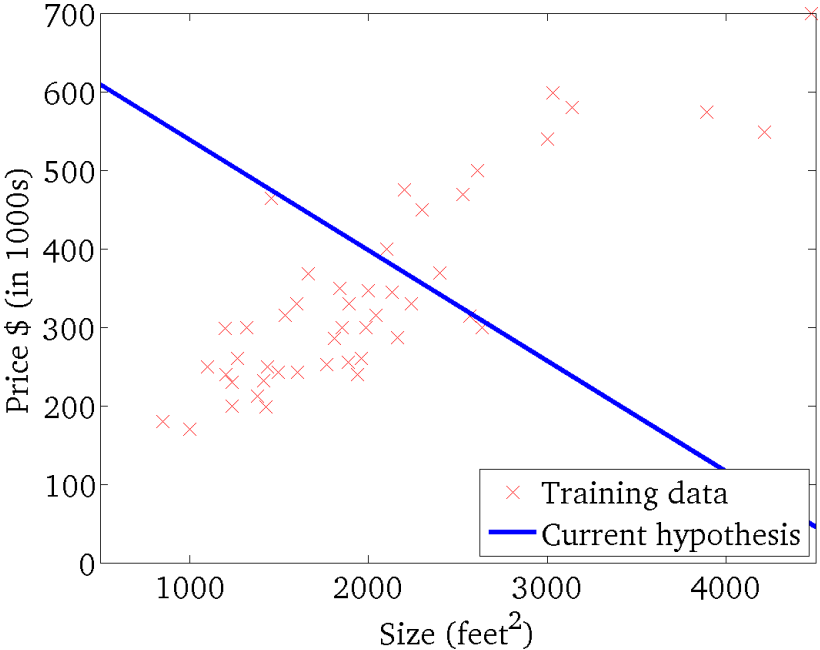


$$J(\theta_0, \theta_1)$$

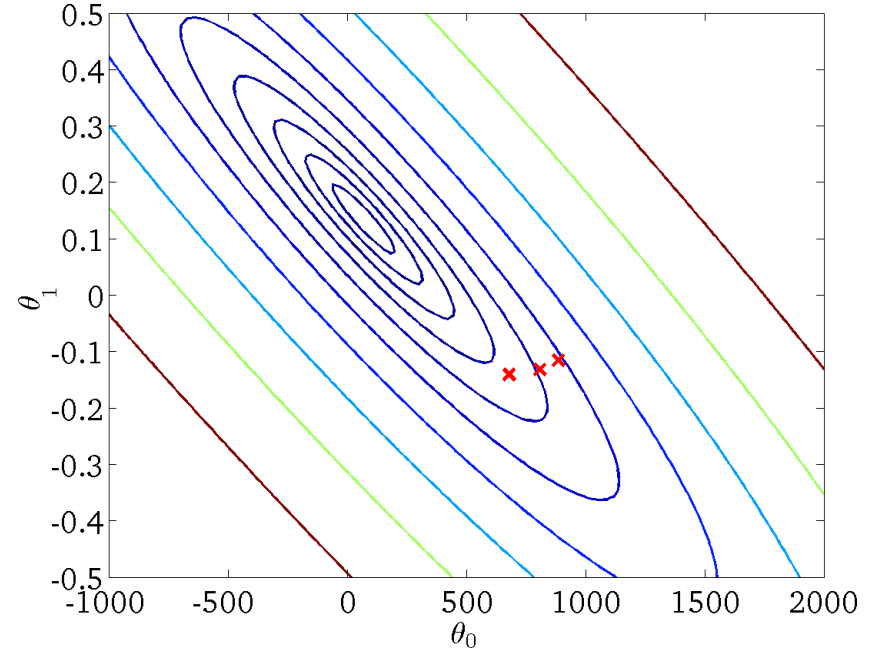


Gradient Descent

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

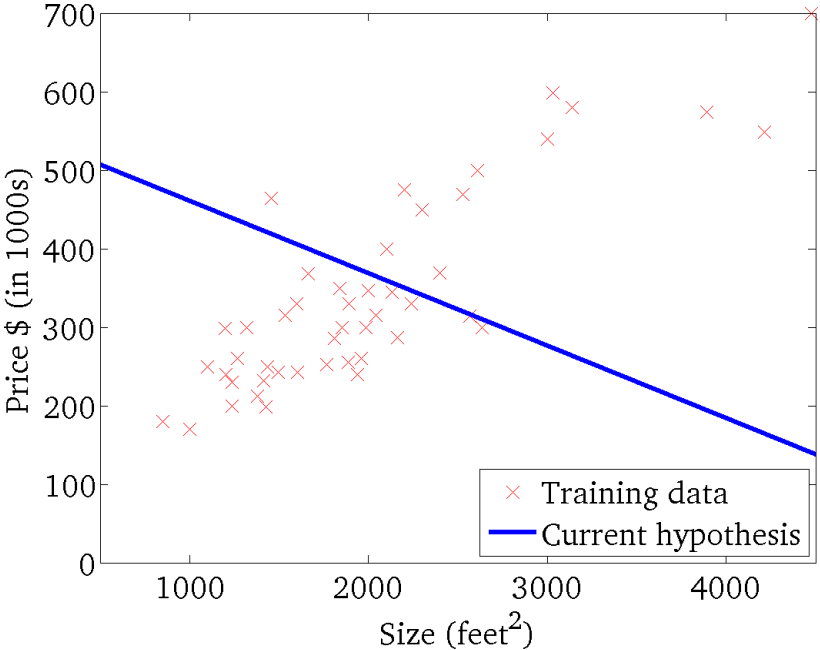


$$J(\theta_0, \theta_1)$$

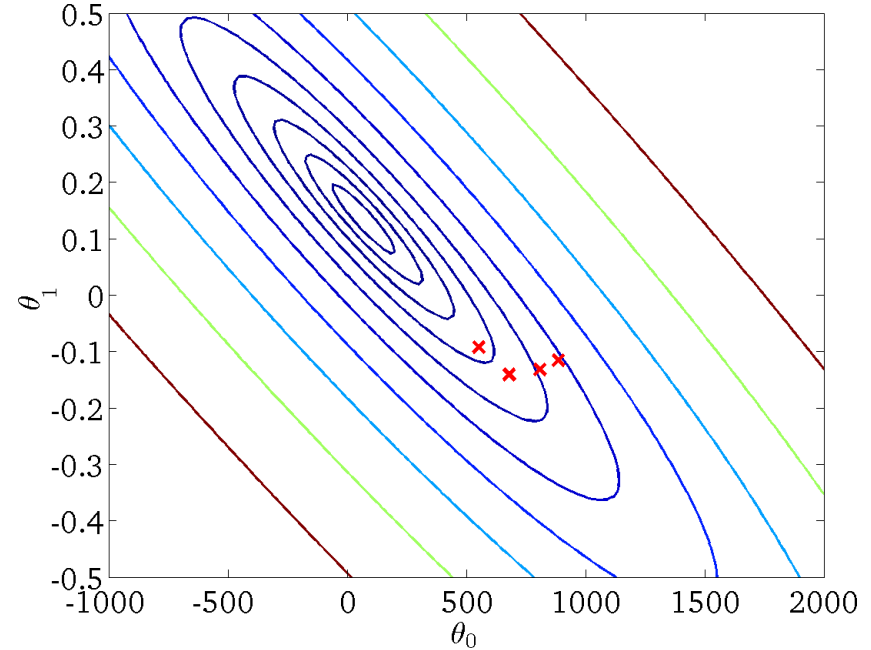


Gradient Descent

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

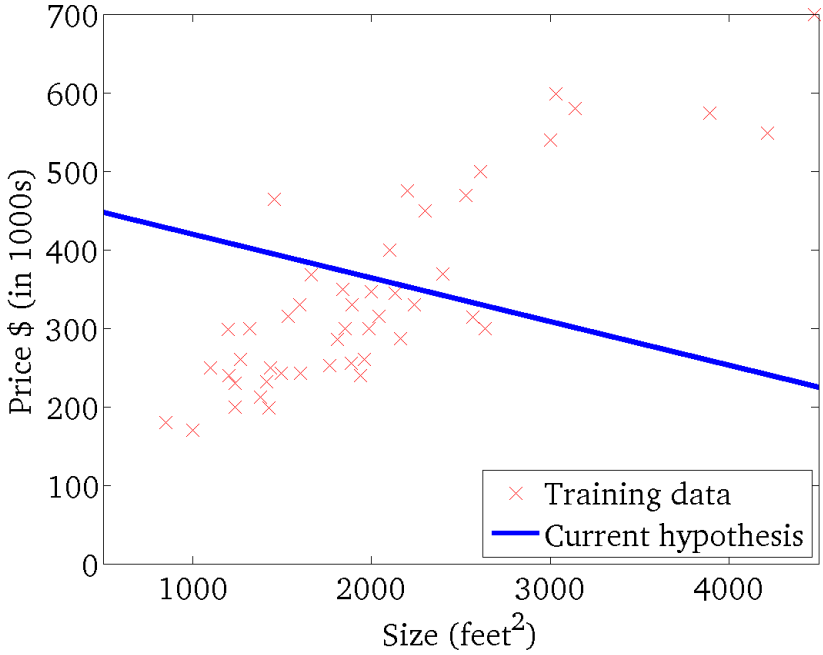


$$J(\theta_0, \theta_1)$$

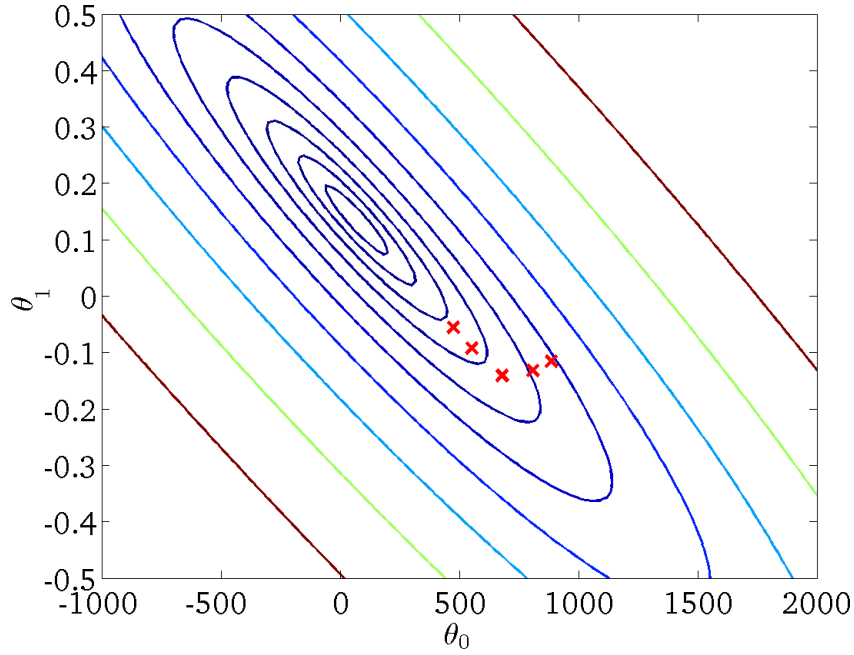


Gradient Descent

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

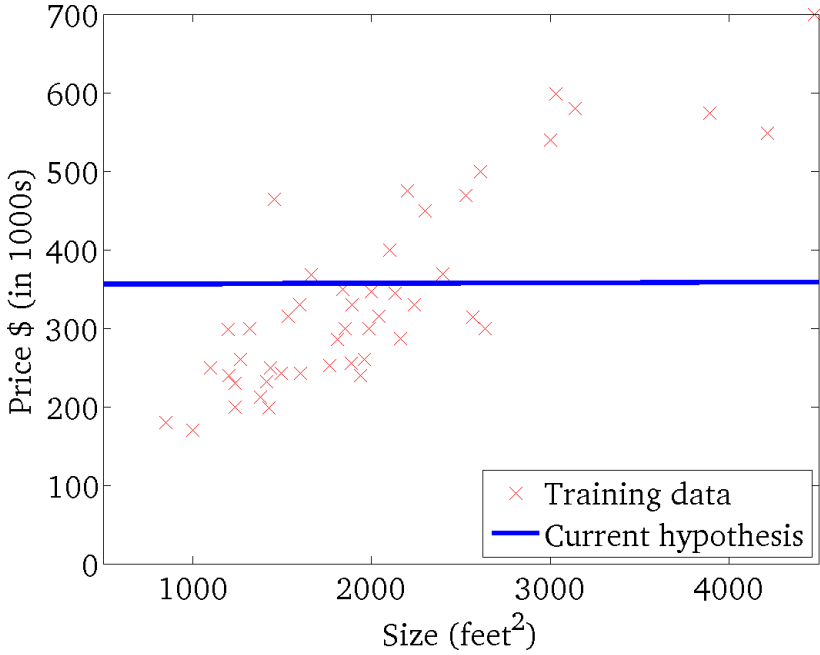


$$J(\theta_0, \theta_1)$$

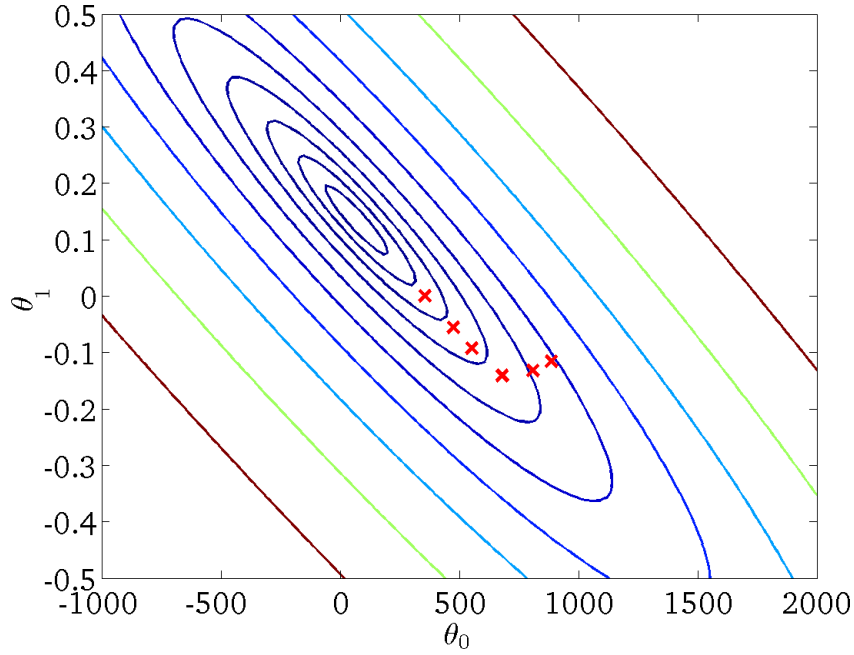


Gradient Descent

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

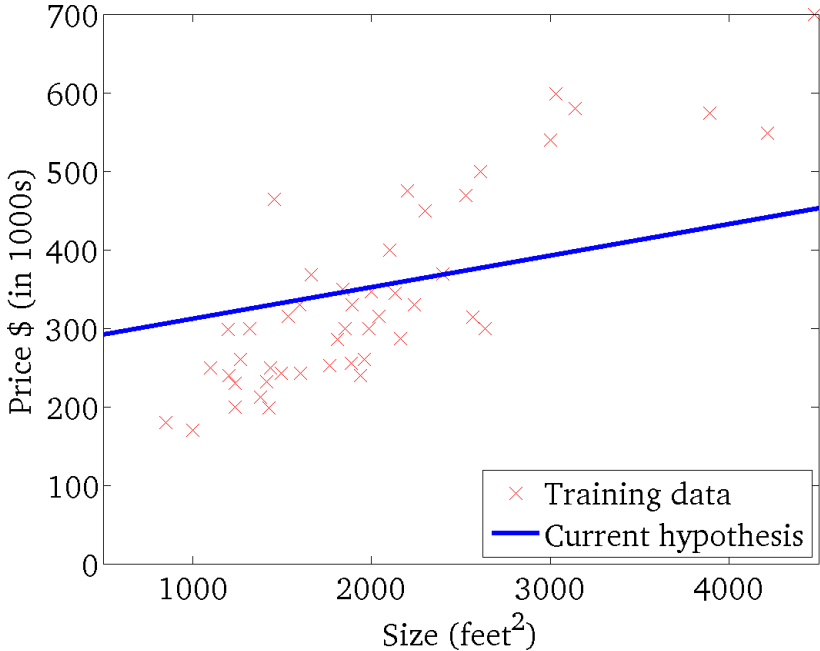


$$J(\theta_0, \theta_1)$$

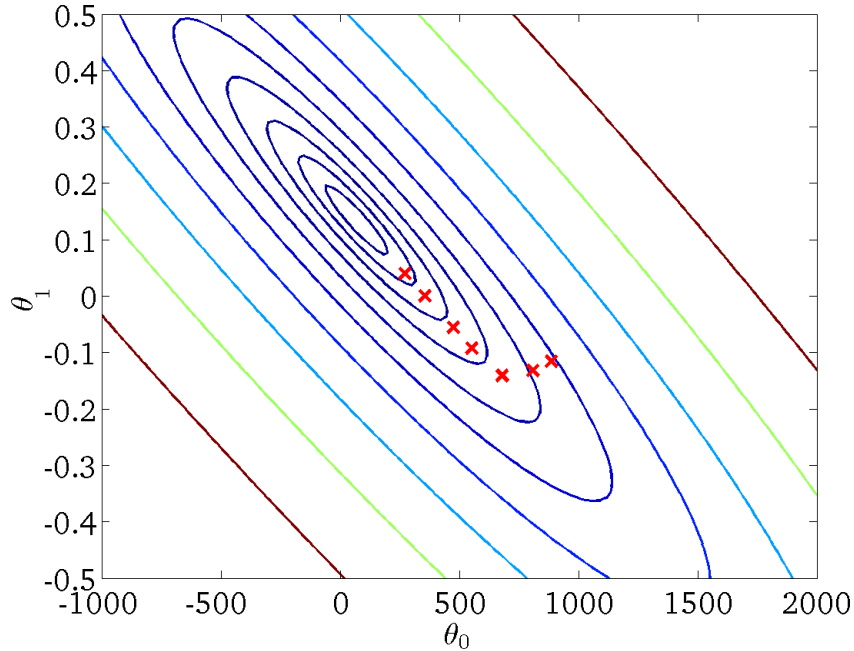


Gradient Descent

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

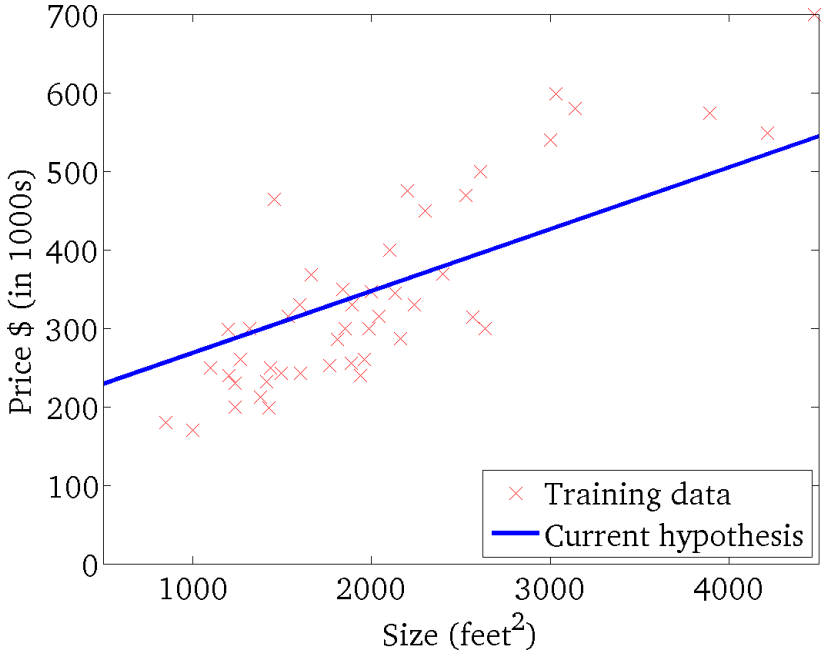


$$J(\theta_0, \theta_1)$$

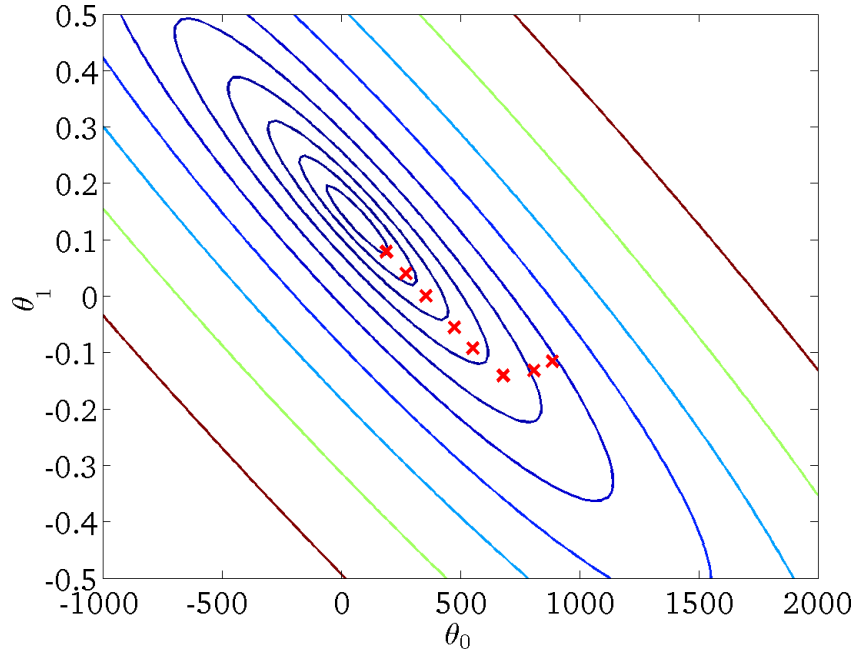


Gradient Descent

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

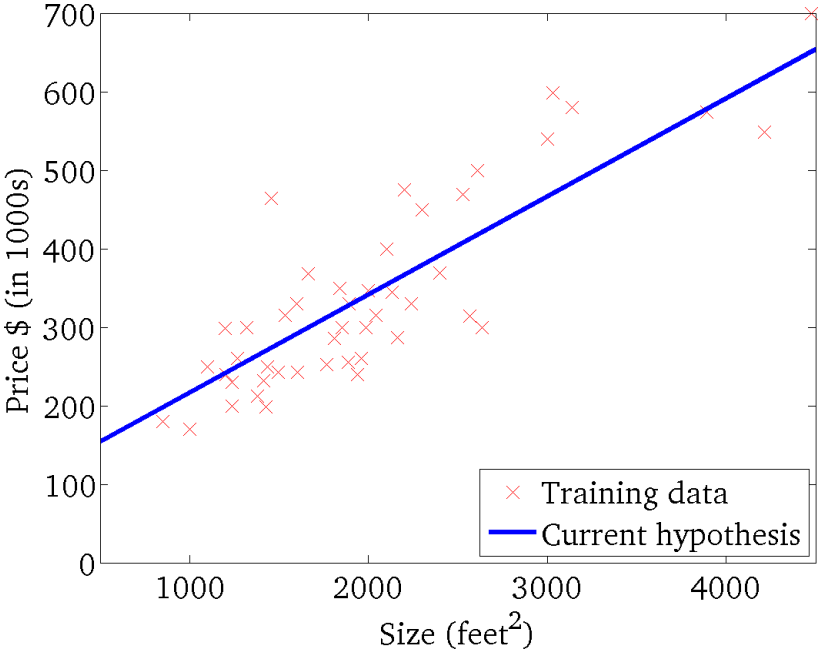


$$J(\theta_0, \theta_1)$$

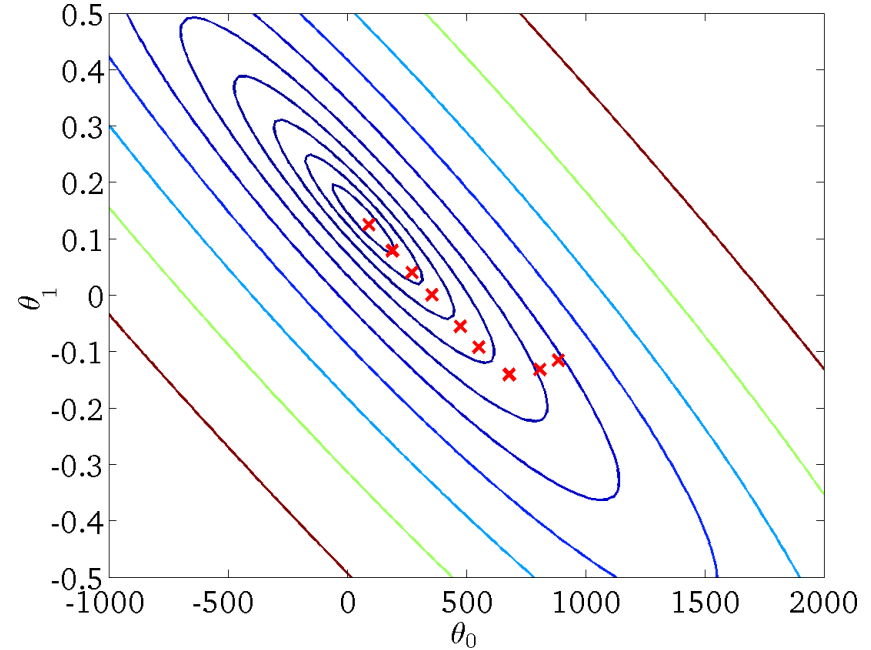


Gradient Descent

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



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



Batch Gradient Descent

49

□ **Batch Gradient Descent.** In each iteration of the algorithm, all training samples are used to update the parameter values.

$$\theta_j := \theta_j - \alpha \frac{\partial}{\partial \theta_j} J(\theta) \quad (\text{for } j = 0 \text{ and } j = 1)$$

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

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

A horizontal decorative bar at the top of the slide, consisting of a red rectangular section on the left and a teal rectangular section on the right.

Multivariate linear regression

Univariate regression (one feature)

51

- Linear regression with one feature.

Price(1000 dollars) y	Meterage (Square foot) x
460	2104
232	1416
315	1534
178	852
...	...

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

Multivariate regression (multiple features)

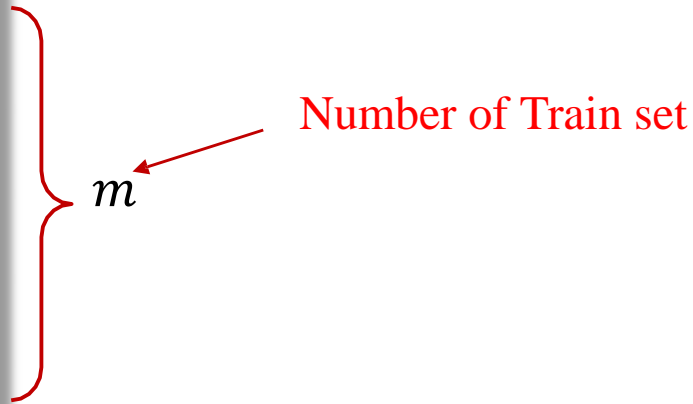
52

□ Linear regression with multiple features.

Price(1000 dollars) y	age of the house(Year) x_4	Number of floors x_3	Number of bedrooms x_2	Meterage (Square foot) x_1
460	45	1	5	2104
232	40	2	3	1416
315	30	2	3	1534
178	36	2	2	852
...

Multivariate regression (multiple features)

Price(1000 dollars) y	age of the house(Year) x_4	Number of floors x_3	Number of bedrooms x_2	Meterage (Square foot) x_1
460	45	1	5	2104
232	40	2	3	1416
315	30	2	3	1534
178	36	1	2	852
...



□ symbols.

- n Number of features
- $x^{(i)}$ Inputs in the i -th training sample
- $x_j^{(i)}$ The j th feature value in the i -th training sample

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

hypothesis

54

- Univariate linear regression.

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

- Multivariate linear regression.

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

- For simplicity, we define $x_0 = 1$

$$x = \begin{bmatrix} x_0 \\ x_1 \\ \vdots \\ x_n \end{bmatrix} \quad \theta = \begin{bmatrix} \theta_0 \\ \theta_1 \\ \vdots \\ \theta_n \end{bmatrix} \quad \Rightarrow \quad h_{\theta}(x) = \theta^T x$$

Gradient Descent in multivariate linear

regression

Gradient Descent

56

□hypothesis

$$h_{\theta}(x) = \theta^T x = \theta_0 + \theta_1 x_1 + \theta_2 x_2 + \dots + \theta_n x_n$$

□parameters

$$\theta = (\theta_0, \theta_1, \dots, \theta_n)$$

□Cost Function

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

repeat until convergence {

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

}

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

Practical tricks in multivariate regression

Scaling features

Determine the learning rate

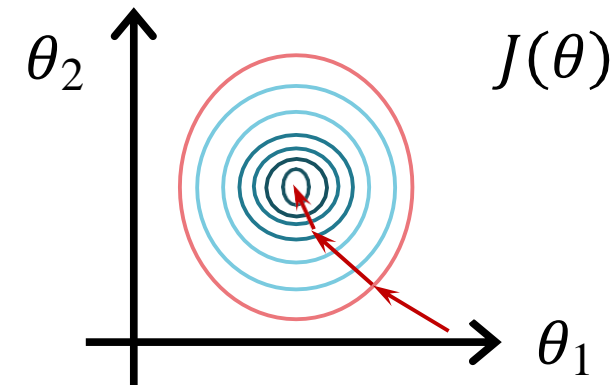
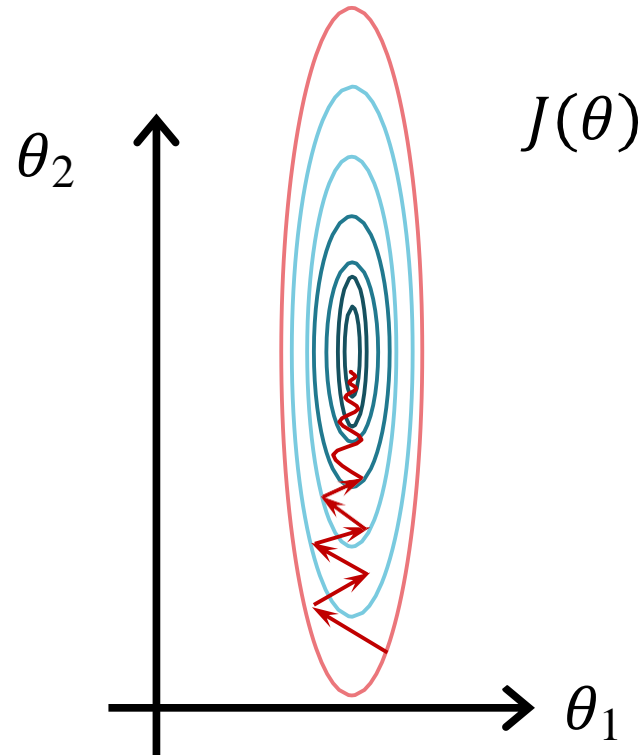
Scaling of features (normalization)

58

- **Ideas.** Ensuring that feature values are on the same scale.
- **Target.** Increasing the speed of convergence in decreasing gradient.
- **Example.** The value of each feature is in a small coefficient from -1 to +1

x_1 : house size (0 to 2000)

x_2 : Number of bedrooms (1 to 5).



Scaling features

59

□ **Scaling**. The value of each feature is in a small coefficient from -1 to +1

□ **Example**

$$x_1 = \frac{\text{size} - 1000}{2000} \quad - 0.5 \leq x_1 \leq 0.5$$

$$x_2 = \frac{\# \text{ bedrooms} - 2}{5} \quad - 0.5 \leq x_2 \leq 0.5$$

□ **Average normalization**

$$x_i = \frac{x_i - \mu_i}{s_i}$$

→ Average

→ Deviation from the norm

Gradient descent

60

□ gradient descent

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

□ Question. How can you make sure that gradient descent works correctly?

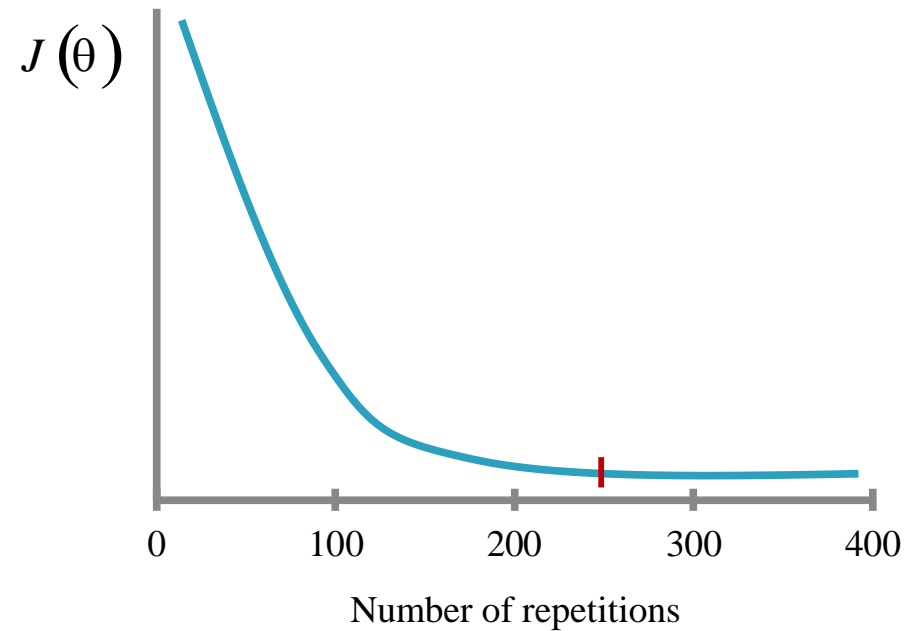
□ Question. What is the right value for the learning rate?

Correct function for gradient descent

61

□ Convergence test

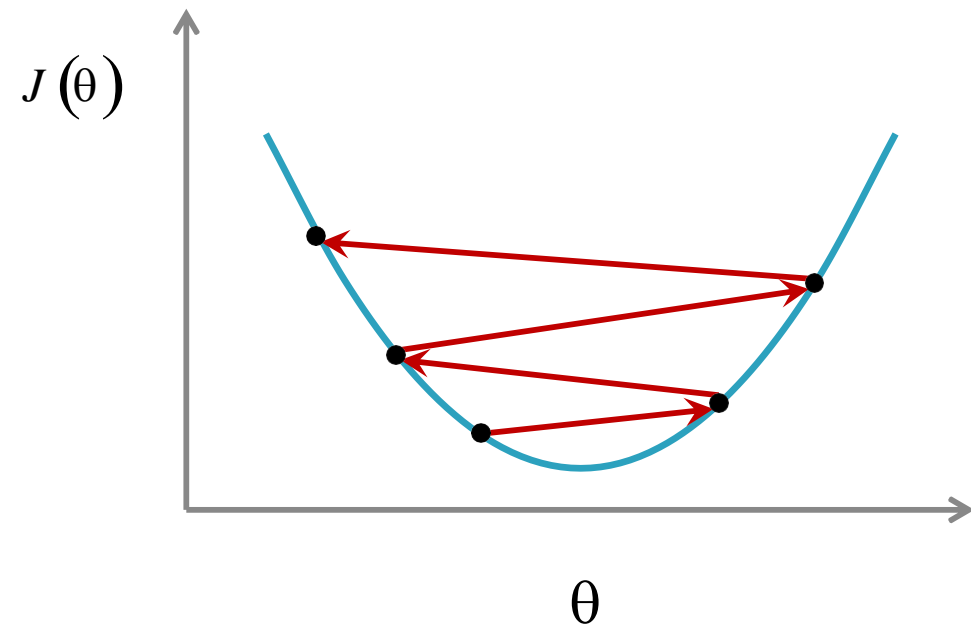
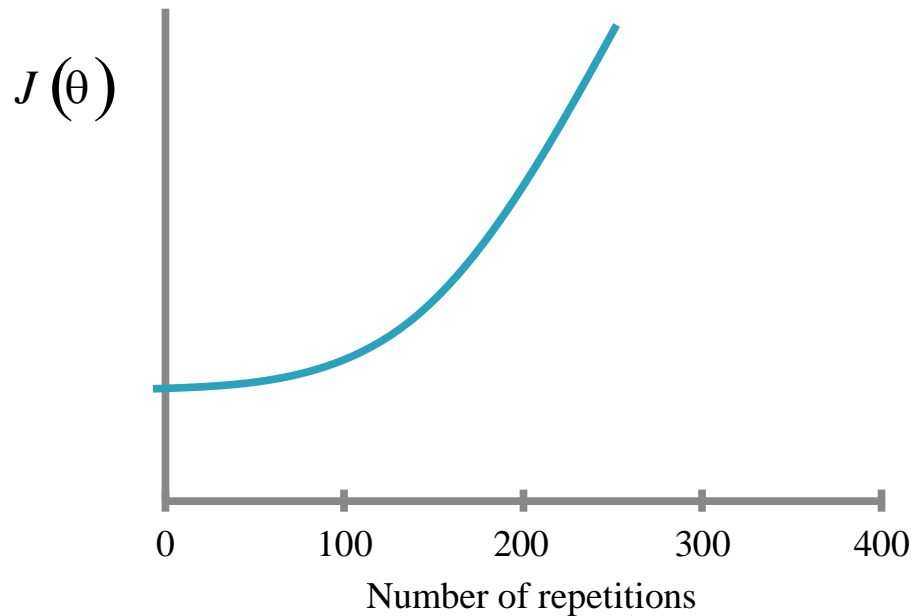
If the value of $J(\theta)$ changes by an amount less than 10^{-3} in one iteration, convergence has occurred.



Indications of gradient descent malfunction

62

- lack of convergence
- solution. Use smaller values for the learning rate, but if the learning rate is too small the convergence will be very slow.



Summary

63

- Learning rate.
 - Very small: Very slow convergence
 - Very large: slow convergence or no convergence

- Choice of learning rate.
 - In order to choose a suitable value for the learning rate, try the following values:

..., 0.001, 0.003, 0.01, 0.03, .01, 0.3, 1.0, ...

Multinomial regression

Pricing a Home: Selecting Features

65

$$h_{\theta}(x) = \theta_0 + \theta_1 \times \underbrace{(\text{frontage})}_{x_1} \times \theta_2 + \underbrace{(\text{depth})}_{x_2}$$

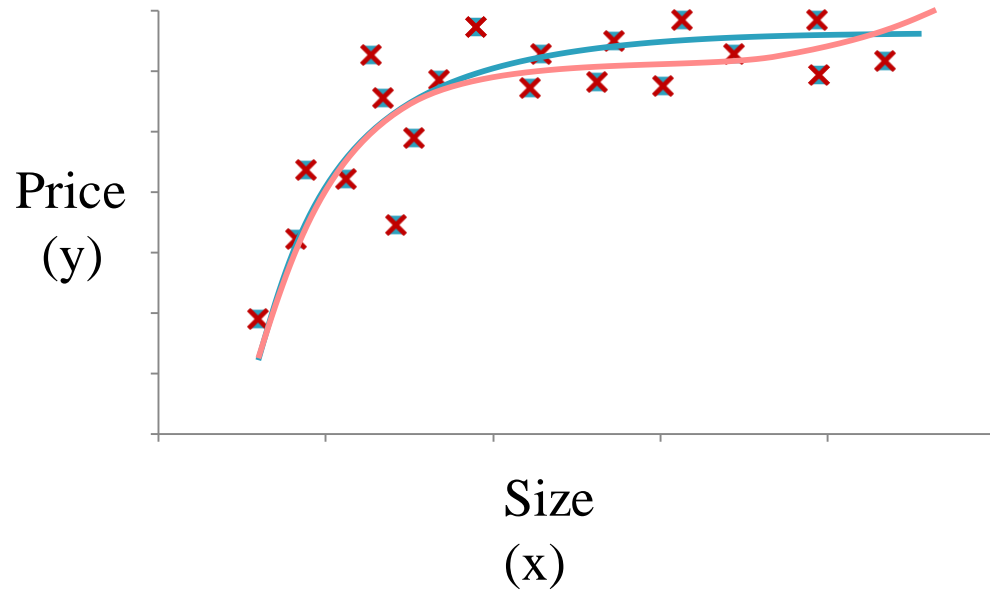
$$\text{Area} = \text{frontage} \times \text{depth}$$

$$h_{\theta}(x) = \theta_0 + \theta_1 \times (\text{Area})$$

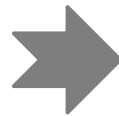


Polynomial regression

66



$$h_{\theta}(x) = \theta_0 + \theta_1x + \theta_2x^2 + \theta_3x^3$$



Polynomial regression

$$h_{\theta}(x) = \theta_0 + \theta_1x + \theta_2x^2$$

$$h_{\theta}(x) = \theta_0 + \theta_1x + \theta_2x^2 + \theta_3x^3$$

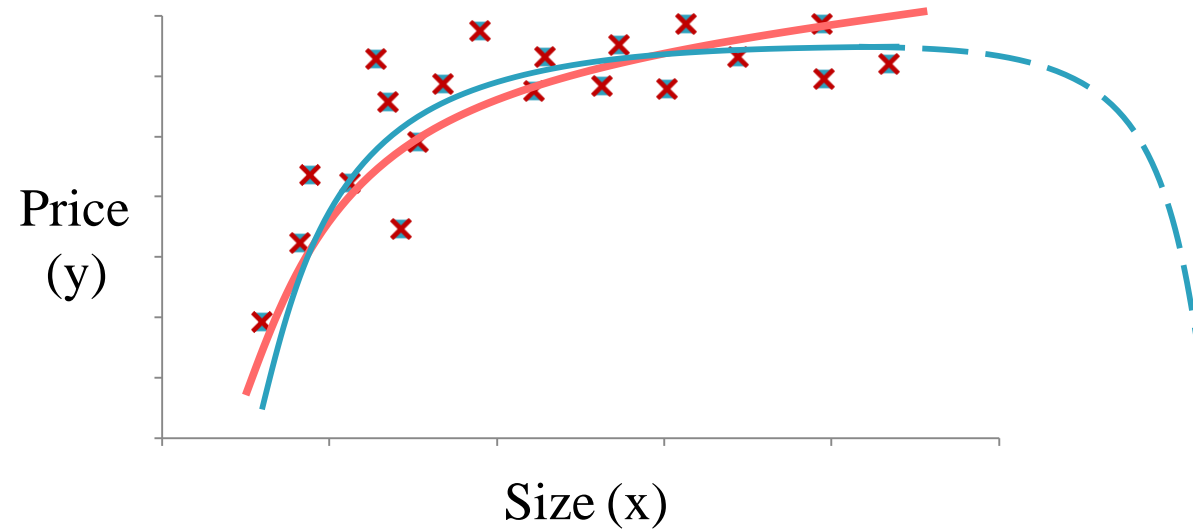
$$x_1 = (\text{Size})^1 \quad 1 \leq x_1 \leq 10^3$$

$$x_2 = (\text{Size})^2 \quad 1 \leq x_2 \leq 10^6$$

$$x_3 = (\text{Size})^3 \quad 1 \leq x_3 \leq 10^9$$

features Selection

67



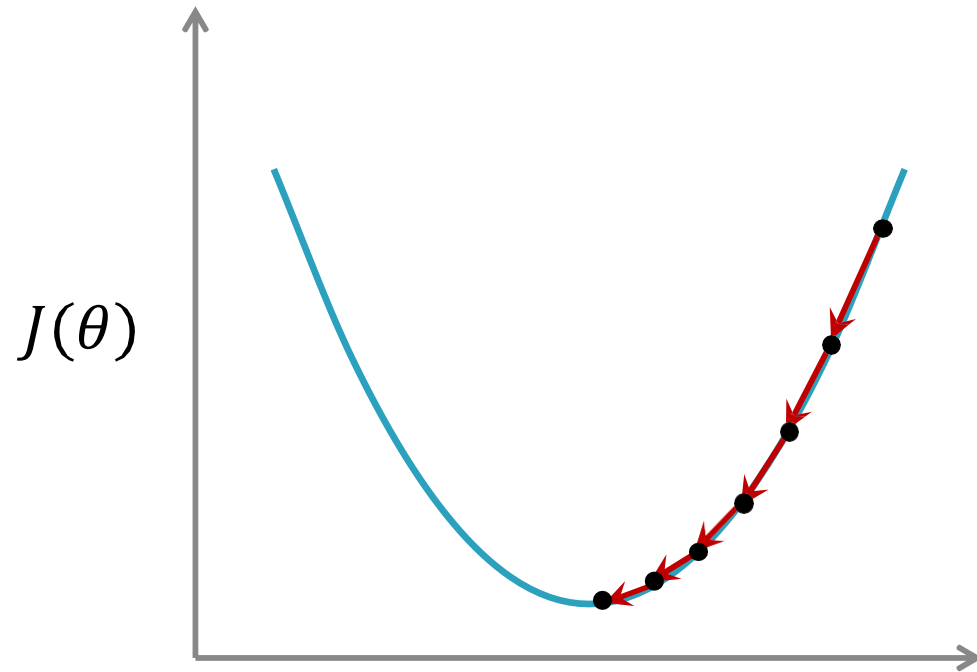
$$h_{\theta}(x) = \theta_0 + \theta_1 (\text{size}) + \theta_2 (\text{size})^2$$

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

Linear regression: normal equation

Decreasing gradient and normal equation

69



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

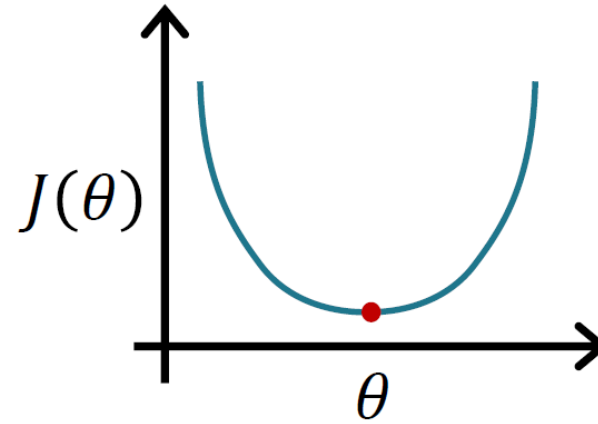
□ **Normal equation.** An analytical method to determine the value of parameters

Normal equation

70

$$\theta \in \mathbb{R}: J(\theta) = a\theta^2 + b\theta + c$$

$$\frac{d}{d\theta}J(\theta) \stackrel{\text{def}}{=} 0$$



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

$$\frac{\partial}{\partial \theta_j} J(\theta) \stackrel{\text{def}}{=} 0$$

$$(j = 0, 1, 2, \dots, n)$$

Normal equation: Example (m = 4)

71

	Size (feet ²)	Number of bedrooms	Number of floors	Age of home (years)	Price (\$1000)
x_0	x_1	x_2	x_3	x_4	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

$$X = \begin{bmatrix} 1 & 2104 & 5 & 1 & 45 \\ 1 & 1416 & 3 & 2 & 40 \\ 1 & 1534 & 3 & 2 & 30 \\ 1 & 852 & 2 & 1 & 36 \end{bmatrix}$$

$$\theta = \begin{bmatrix} \theta_0 \\ \theta_1 \\ \theta_2 \\ \theta_3 \\ \theta_4 \end{bmatrix}$$

$$y = \begin{bmatrix} 460 \\ 232 \\ 315 \\ 178 \end{bmatrix}$$

$$X\theta = y$$

Normal equation: Example (m = 5)

72

	Size (feet ²)	Number of bedrooms	Number of floors	Age of home (years)	Price (\$1000)
x_0	x_1	x_2	x_3	x_4	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	3000	4	1	38	540

$$X = \begin{bmatrix} 1 & 2104 & 5 & 1 & 45 \\ 1 & 1416 & 3 & 2 & 40 \\ 1 & 1534 & 3 & 2 & 30 \\ 1 & 852 & 2 & 1 & 36 \\ 1 & 3000 & 4 & 1 & 38 \end{bmatrix}$$

$$\theta = \begin{bmatrix} \theta_0 \\ \theta_1 \\ \theta_2 \\ \theta_3 \\ \theta_4 \end{bmatrix}$$

$$y = \begin{bmatrix} 460 \\ 232 \\ 315 \\ 178 \\ 540 \end{bmatrix}$$

$$X\theta = y$$

Normal equation: general state

73

□ m educational sample; n feature

$$X\theta = y$$

$$\mathbf{x}^{(i)} = \begin{bmatrix} x_0^{(i)} \\ x_1^{(i)} \\ x_2^{(i)} \\ \vdots \\ x_n^{(i)} \end{bmatrix} \in \mathbb{R}^{n+1} \quad X = \begin{bmatrix} \text{---} (x^{(1)})^T \text{---} \\ \text{---} (x^{(2)})^T \text{---} \\ \text{---} (x^{(3)})^T \text{---} \\ \vdots \\ \text{---} (x^{(m)})^T \text{---} \end{bmatrix} \in \mathbb{R}^{m \times (n+1)} \quad \mathbf{y} = \begin{bmatrix} y^{(1)} \\ y^{(2)} \\ y^{(3)} \\ \vdots \\ y^{(m)} \end{bmatrix} \in \mathbb{R}^m$$

Design matrix

Normal equation

74

□ Solving linear equations

$$X\theta = y$$

$$X^T X\theta = X^T y \quad \leftarrow \text{Normal equation}$$

$$\theta = \underbrace{(X^T X)^{-1} X^T}_{X^+} y \quad \leftarrow \text{Pseudo-inverse matrix}$$

Octave:

```
theta = pinv(X'*X)*X'*y
```

Gradient Descent and normal equation

75

normal equation

- No need to choose α
- No need to repeat

- It is slow for large n due to the need to calculate the inverse of the $X^T X$ matrix.

$$n < 10000$$

Gradient Descent

- need to choose α
- need to repeat

- It works well even for very large values of n .

$$n \geq 10000$$

Normal equation and irreversibility of $X^T X$

76

□ Normal equation

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

□ Question. But what if $X^T X$ is not invertible?

Octave: **`theta = pinv(X'*X)*X'*y`**


pseudo reverse

Causes of $X^T X$ irreversibility

77

- Redundancy of features (linear dependence).

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

$$x_2 = \text{size (m}^2\text{)}$$

$$x_1 = (3.28)^2 x_2$$

- A very large number of features (eg $n \geq m$).

❖ solution. Remove some features by setting [Next]

Regression with local weighting

Parametric and non-parametric learning methods

79

□ Parametric learning methods.

- There is a fixed set of parameters.
- We do not need a training set to predict new data.
- Example: regression and logistic regression [continue].

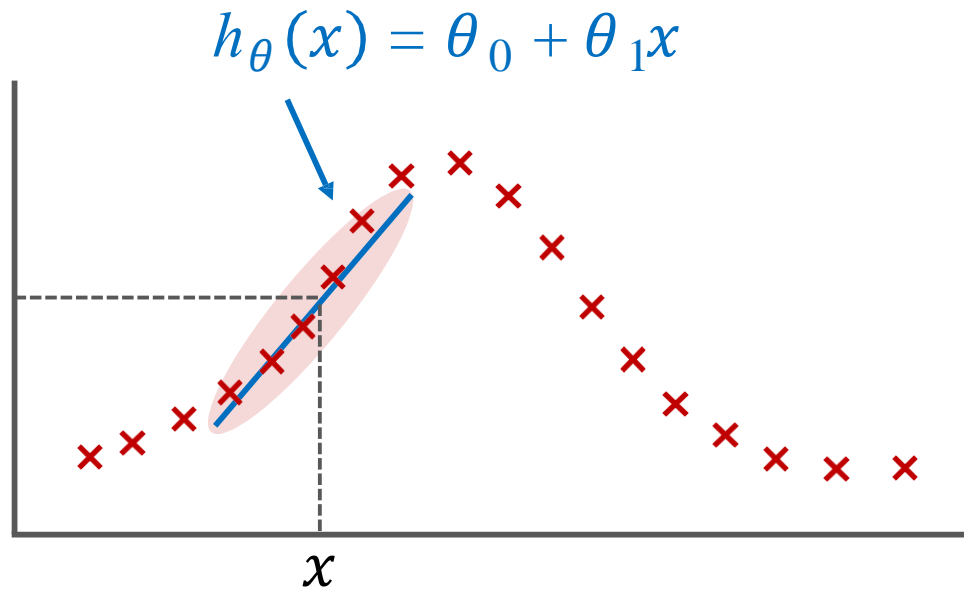
□ Non-parametric learning methods.

- The number of parameters increases (linearly) with the increase in the size of the training set.
- In order to predict new data, we need the entire training set.
- Example: regression and logistic regression [continue].

Regression with local volatility

80

□ Ideas. Giving more importance to closer data.



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

$$w^{(i)} = \exp \left(- \frac{\left(x^{(i)} - x \right)^2}{2r^2} \right)$$

Bandwidth

Regression: probabilistic interpretation

Regression: probabilistic interpretation

82

□ model

$$y^{(i)} = \theta^T x^{(i)} + \epsilon^{(i)}$$

error

□ Error

- ❖ Considering the influence of other factors
 - ❖ such as unintended features
- ❖ Considering the effect of noise.

$$\epsilon^{(i)} \sim N(0, \sigma^2) = \frac{1}{\sigma\sqrt{2\pi}} \exp\left(-\frac{\epsilon^2}{2\sigma^2}\right)$$

□ The errors are independent and follow a uniform Gaussian distribution. [*iid*]

$$y^{(i)} \sim N(\theta^T x^{(i)}, \sigma^2) = \frac{1}{\sigma\sqrt{2\pi}} \exp\left(-\frac{(y^{(i)} - \theta^T x^{(i)})^2}{2\sigma^2}\right)$$

Maximum likelihood estimation

83

□ Exponentiation function. Probability of seeing training data as a function of parameters θ

$$L(\theta) = p(Y|X; \theta) = \prod_{i=1}^m \frac{1}{\sigma\sqrt{2\pi}} \exp\left(-\frac{(y^{(i)} - \theta^T x^{(i)})^2}{2\sigma^2}\right)$$

□ The logarithm of the exponential function.

$$\begin{aligned} l(\theta) = \ln L(\theta) &= \ln \prod_{i=1}^m \frac{1}{\sigma\sqrt{2\pi}} \exp\left(-\frac{(y^{(i)} - \theta^T x^{(i)})^2}{2\sigma^2}\right) \\ &= \sum_{i=1}^m \ln \left[\frac{1}{\sigma\sqrt{2\pi}} \exp\left(-\frac{(y^{(i)} - \theta^T x^{(i)})^2}{2\sigma^2}\right) \right] \\ &= m \ln \frac{1}{\sigma\sqrt{2\pi}} - \frac{1}{2\sigma^2} \sum_{i=1}^m (y^{(i)} - \theta^T x^{(i)})^2 \end{aligned}$$

Maximum likelihood estimation

84

Maximum likelihood estimation.

Choosing a value for the parameter θ so that $l(\theta)$ is maximized.

$$\begin{aligned}\theta_{MLE} &= \arg \max_{\theta} l(\theta) \\ &= \arg \max_{\theta} m \ln \frac{1}{\sigma \sqrt{2\pi}} - \frac{1}{2\sigma^2} \sum_{i=1}^m (y^{(i)} - \theta^T x^{(i)})^2 \\ &= \arg \max_{\theta} -\frac{1}{2} \sum_{i=1}^m (y^{(i)} - \theta^T x^{(i)})^2 \\ &= \arg \min_{\theta} \frac{1}{2} \sum_{i=1}^m (y^{(i)} - \theta^T x^{(i)})^2 \longleftarrow \text{Cost function}\end{aligned}$$