

# Machine Learning and Data Analytics

## ME 5013- Fall 2019

### Lecture 06

- Gradient Descent cont.
- Polynomial Regression



The University of Texas at San Antonio™

Adel Alaeddini, PhD

Associate Professor of Mechanical Engineering

Advanced Data Engineering Lab

[adel.alaeddini@utsa.edu](mailto:adel.alaeddini@utsa.edu)

Hypothesis:  $h_{\beta}(x) = \beta_0 + \beta_1 x$

Cost Function:  $\frac{1}{2n} \sum_{i=1}^n (y_i - \beta_0 - \beta_1 x^{(i)})^2$

Parameters:  $\beta_0, \beta_1$



$x_i$

Gradient descent:

Repeat until convergence{

$$\beta_0 := \beta_0 + \alpha \frac{1}{n} \sum_{i=1}^n (y_i - \beta_0 - \beta_1 x^{(i)})$$

$$\beta_1 := \beta_1 + \alpha \frac{1}{n} \sum_{i=1}^n (y_i - \beta_0 - \beta_1 x^{(i)}) x^{(i)}$$

}

Hypothesis:  $h_{\beta}(x) = \beta_0 + \beta_1 x_1 + \dots + \beta_p x_p$

Cost Function:  $\frac{1}{2n} \sum_{i=1}^n (y_i - \beta_0 - \beta_1 x_1^{(i)} - \dots - \beta_p x_p^{(i)})^2$

Parameters:  $\beta_0, \beta_1, \dots, \beta_p$

Gradient descent:

Repeat until convergence{

$$\beta_0 := \beta_0 + \alpha \frac{1}{n} \sum_{i=1}^n (y_i - \beta_0 - \beta_1 x_1^{(i)} - \dots - \beta_p x_p^{(i)})$$

$$\beta_1 := \beta_1 + \alpha \frac{1}{n} \sum_{i=1}^n (y_i - \beta_0 - \beta_1 x_1^{(i)} - \dots - \beta_p x_p^{(i)}) x_1^{(i)}$$

...

$$\beta_p := \beta_p + \alpha \frac{1}{n} \sum_{i=1}^n (y_i - \beta_0 - \beta_1 x_1^{(i)} - \dots - \beta_p x_p^{(i)}) x_p^{(i)}$$

}

## Feature Scaling

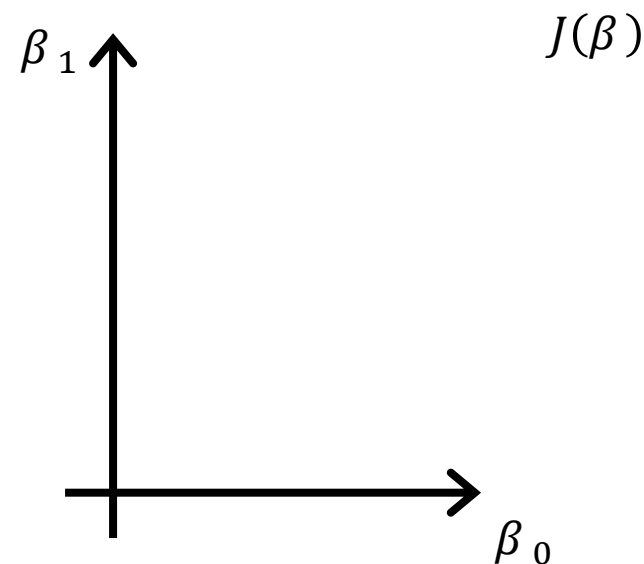
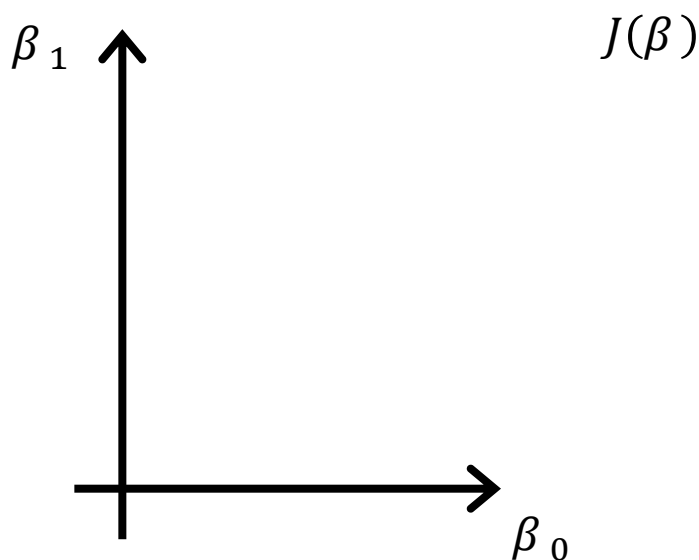
Idea: Make sure features are on a similar scale.

E.g.  $x_1 = \text{size (0-2000 feet}^2\text{)}$

$x_2 = \text{number of bedrooms (1-5)}$

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

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



Get every feature into approximately a  $-1 \leq x_i \leq 1$  range.

## Mean normalization

- Replace  $x_i$  with  $x_i - \bar{x}_i$ , where  $\bar{x}_i$  is the average value of  $x_i$  in the training set, to make features have approximately zero mean (Do not apply to  $x_0$ ).
- Divide  $x_i$  by  $R_i$ , where  $R_i$  is the range of  $x_i$  in the training set,

OR

Divide  $x_i$  by  $S_i$ , where  $S_i$  is the standard deviation of  $x_i$  in the training set

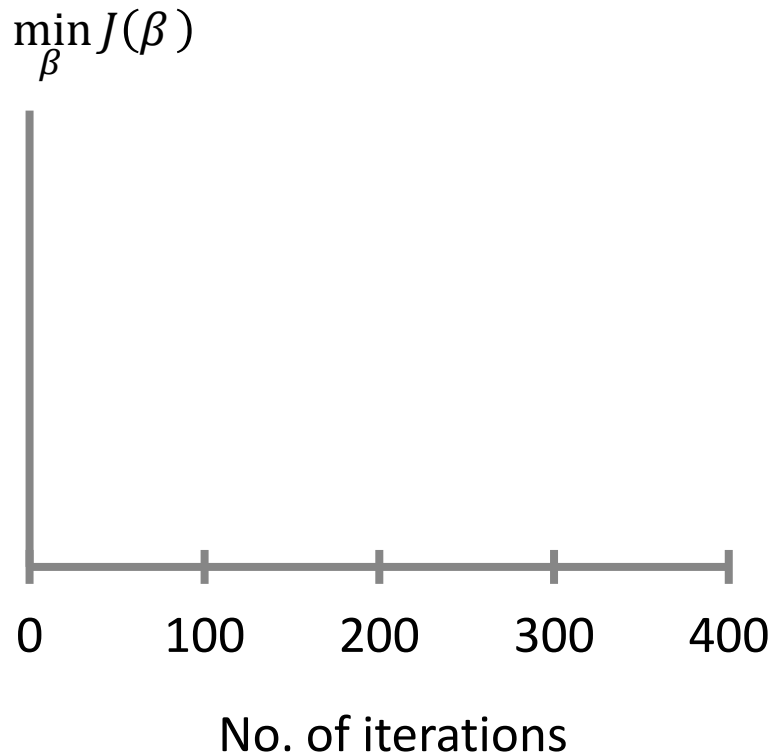
E.g. 
$$x_2 = \frac{\#bedrooms - 2}{5}$$

## Gradient descent

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

- “Debugging”: How to make sure gradient descent is working correctly.
- How to choose learning rate  $\alpha$ .

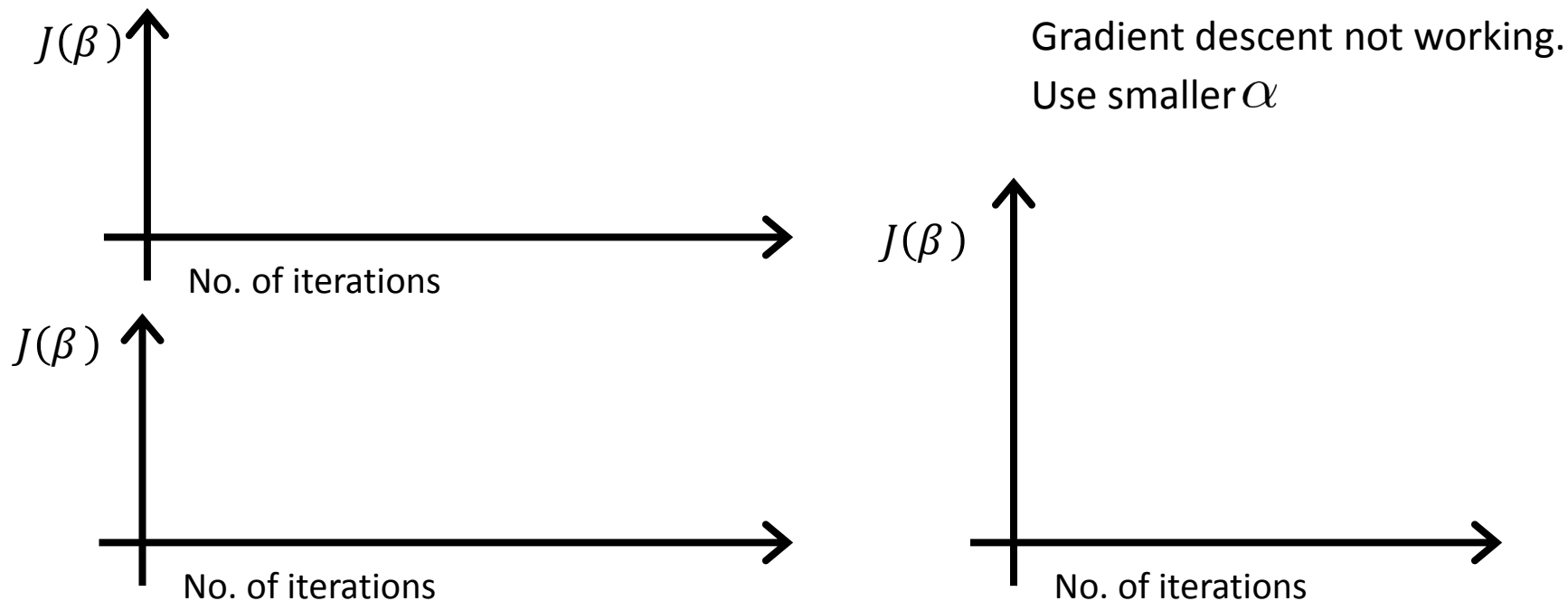
## Making sure gradient descent is working correctly.



Example automatic convergence test:

Declare convergence if  $J(\beta)$  decreases by less than  $10^{-3}$  in one iteration.

## Making sure gradient descent is working correctly.



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

## Summary:

- If  $\alpha$  is too small: slow convergence.
- If  $\alpha$  is too large:  $J(\beta)$  may not decrease on every iteration; may not converge.

To choose  $\alpha$ , try

$\dots, 0.001, \quad , 0.01, \quad , 0.1, \quad , 1, \dots$



$n$  training examples,  $p$  features.

### Gradient Descent

- Need to choose  $\alpha$ .
- Needs many iterations.
- Works well even when  $p$  is large.
- More robust to linearly dependencies

### Normal Equation

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

What if  $X^T X$  is non-invertible?

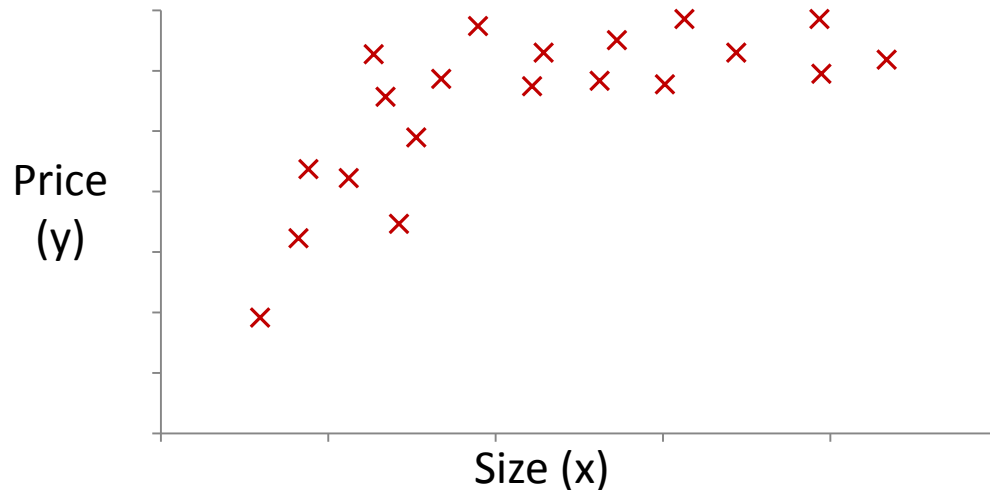
- Redundant features (linearly dependent).  
E.g.  $x_1 = \text{size in feet}^2$   
 $x_2 = \text{size in m}^2$
- Too many features (e.g.  $m \leq n$ ).
  - Delete some features, or use regularization.

## Housing prices prediction

$$h_{\beta}(x) = \beta_0 + \beta_1 \text{frontage} + \cdots + \beta_p \text{depth}$$



## Polynomial regression



$$\beta_0 + \beta_1 x + \beta_2 x^2$$

$$\beta_0 + \beta_1 x + \beta_2 x^2 + \beta_3 x^3$$

$$h_{\beta}(x) = \beta_0 + \beta_1 x + \beta_2 x^2 + \beta_3 x^3$$

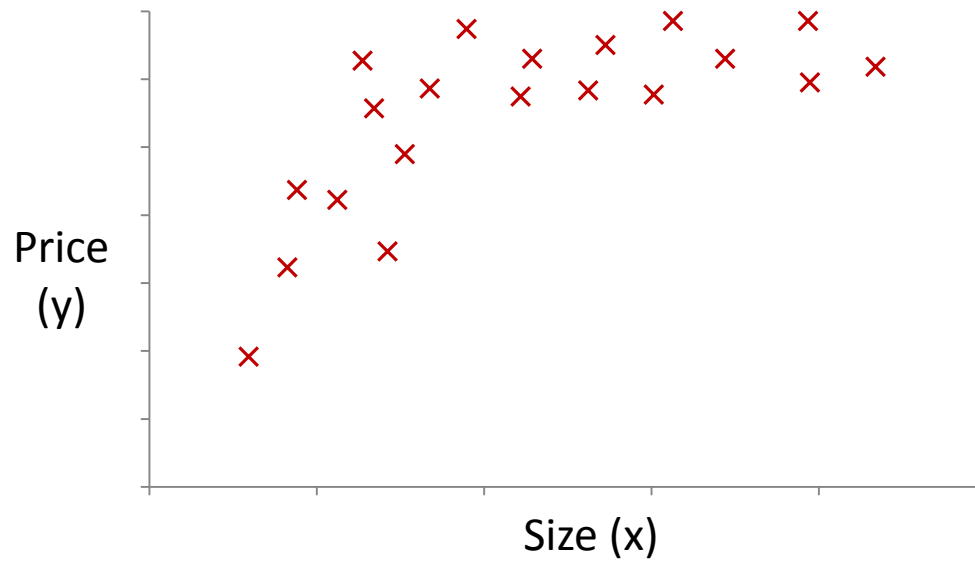
$$\beta_0 + \beta_1(\text{size}) + \beta_2(\text{size})^2 + \beta_3(\text{size})^3$$

$$x_1 = (\text{size})$$

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

$$x_3 = (\text{size})^3$$

## Choice of features



$$h_{\beta}(x) = \beta_0 + \beta_1(\text{size}) + \beta_2(\text{size})^2$$

$$h_{\beta}(x) = \beta_0 + \beta_1(\text{size}) + \beta_2\sqrt{(\text{size})}$$