

# Linear regression

**Lesson #1:** Linear regression with one variable

**Lesson #2:** Linear regression with multiple variables

**Duration:** 5 hrs

# Lesson #2: Linear regression with multiple variables

- **Duration:** 2 hrs
- **Outline:**
  - Model representation
  - Gradient descend for multiple variables
    - Feature Scaling
    - Learning rate
    - Features and polynomial regression
    - Normal equation

## Multiple features (variables).

Size (feet <sup>2</sup> )	Price (\$1000)
$x$	$y$
2104	460
1416	232
1534	315
852	178
...	...

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

## Multiple features (variables).

Size (feet <sup>2</sup> )	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
...	...	...	...	...

## Multiple features (variables).

Size (feet <sup>2</sup> )	Number of bedrooms	Number of floors	Age of home (years)	Price (\$1000)	
$x_1$	$x_2$	$x_3$	$x_4$	$y$	
2104	5	1	45	460	
1416	3	2	40	232	(m = 47)
1534	3	2	30	315	
852	2	1	36	178	
...	...	...	...	...	

Notation:

$n$  = number of features ( $n = 4$ )

$x^{(i)}$  = input (features) of  $i^{th}$  training example.

$x_j^{(i)}$  = value of feature  $j$  in  $i^{th}$  training example.

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

$$x_3^{(2)} = 2$$

Hypothesis:

Previously:  $h_{\theta}(x) = \theta_0 + \theta_1 x$  (one variable)

Multiple variables:

$$h_{\theta}(x) = \theta_0 + \theta_1 x_1 + \theta_2 x_2 + \theta_3 x_3 + \theta_4 x_4$$

$$E.g., h_{\theta}(x) = 80 + 0.1x_1 + 0.01x_2 + 3x_3 - 2x_4$$

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

For convenience of notation, define  $x_0 = 1$  ( $x_0^{(i)} = 1$ )

$$x = \begin{bmatrix} x_0 \\ x_1 \\ x_2 \\ \vdots \\ x_n \end{bmatrix} \in \mathbb{R}^{n+1}$$

$$\theta = \begin{bmatrix} \theta_0 \\ \theta_1 \\ \theta_2 \\ \vdots \\ \theta_n \end{bmatrix} \in \mathbb{R}^{n+1}$$

$$[\theta_0 \theta_1 \theta_2 \dots \theta_n] \begin{bmatrix} x_0 \\ x_1 \\ x_2 \\ \vdots \\ x_n \end{bmatrix}$$

$$\begin{aligned} h_{\theta}(x) &= \theta_0 x_0 + \theta_1 x_1 + \dots + \theta_n x_n \\ &= \theta^T x \end{aligned}$$

$$\theta^T x$$

1x(n + 1)  
matrix

Multivariate linear regression.

# Lesson #2: Linear regression with multiple variables

- **Duration:** 2 hrs
- **Outline:**
  - Model representation
  - Gradient descent for multiple variables



Hypothesis:  $h_{\theta}(x) = \theta^T x = \overset{x_0 = 1}{\theta_0 x_0} + \theta_1 x_1 + \theta_2 x_2 + \cdots + \theta_n x_n$

Parameters:  $\boxed{\theta_0, \theta_1, \dots, \theta_n}$   $\theta$   $(n + 1)$  dimensional vector

Cost function:

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

Gradient descent:

Repeat {  
     $\theta_j := \theta_j - \alpha \frac{\partial}{\partial \theta_j} \underset{J(\theta)}{J(\theta_0, \dots, \theta_n)}$   
} (simultaneously update for every  $j = 0, \dots, n$ )

## Gradient Descent

Previously (n=1):

Repeat {

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

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

(simultaneously update  $\theta_0, \theta_1$ )

}

New algorithm ( $n \geq 1$ ):

Repeat {  $\frac{\partial}{\partial \theta_j} J(\theta)$  .

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

(simultaneously update  $\theta_j$  for  
 $j = 0, \dots, n$ )

}

$x_0^{(i)} = 1$

---

$$\theta_0 := \theta_0 - \alpha \frac{1}{m} \sum_{i=1}^m (h_{\theta}(x^{(i)}) - y^{(i)}) x_0^{(i)}$$

$$\theta_1 := \theta_1 - \alpha \frac{1}{m} \sum_{i=1}^m (h_{\theta}(x^{(i)}) - y^{(i)}) x_1^{(i)}$$

$$\theta_2 := \theta_2 - \alpha \frac{1}{m} \sum_{i=1}^m (h_{\theta}(x^{(i)}) - y^{(i)}) x_2^{(i)}$$

...

# Lesson #2: Linear regression with multiple variables

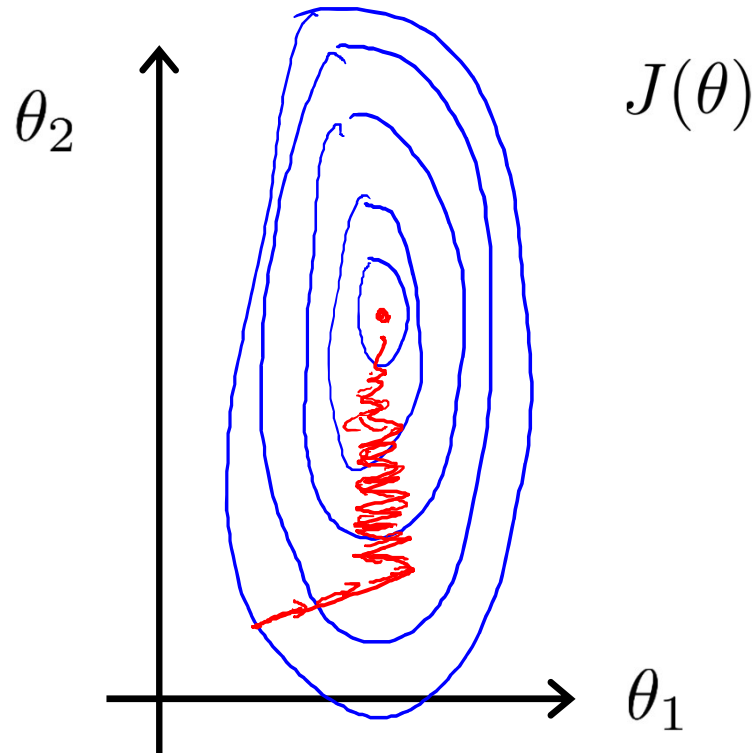
- **Duration:** 2 hrs
- **Outline:**
  - Model representation
  - Gradient descent for multiple variables
    - Feature Scaling
  - Learning rate
  - Features and polynomial regression
  - Normal equation

## Feature Scaling

Idea: Make sure features are on a similar scale.

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

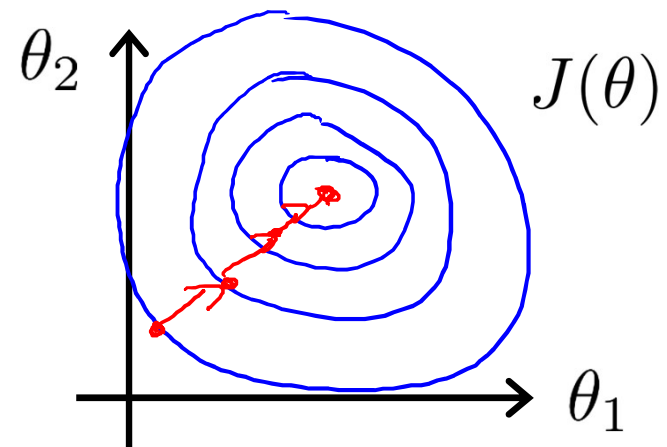
$x_2 = \text{number of bedrooms (1-5)}$



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

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

$$-1 \leq x_1 \leq 1 \quad -1 \leq x_2 \leq 1$$



## Feature Scaling

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

$$x_0 = 1$$

$$0 \leq x_1 \leq 3$$

$$-2 \leq x_2 \leq 0.5$$

$$-100 \leq x_3 \leq 100 \quad \text{Too big}$$

$$-0.0001 \leq x_4 \leq 0.0001 \quad \text{Too small}$$

## Mean normalization

Replace  $x_i$  with  $x_i - \mu_i$  to make features have approximately zero mean  
(Do not apply to  $x_0 = 1$ ).

E.g.  $x_1 = \frac{\text{size} - 1000}{2000}$       Avg. size 1000

$$x_2 = \frac{\#bedrooms - 2}{4} \quad 1-5 \text{ bedrooms}$$

$$-0.5 \leq x_1 \leq 0.5, -0.5 \leq x_2 \leq 0.5$$

$$x_i \leftarrow \frac{x_i - \mu_i}{s_i} \quad \begin{array}{l} \mu_i: \text{average value of } x_i \text{ in training set} \\ s_i: \text{range (max - min) or standard deviation} \end{array}$$

## Lesson #2: Linear regression with multiple variables

- **Duration:** 2 hrs
- **Outline:**
  - Model representation
  - Gradient descent for multiple variables
    - Feature Scaling
    - Learning rate
    - Features and polynomial regression
    - Normal equation

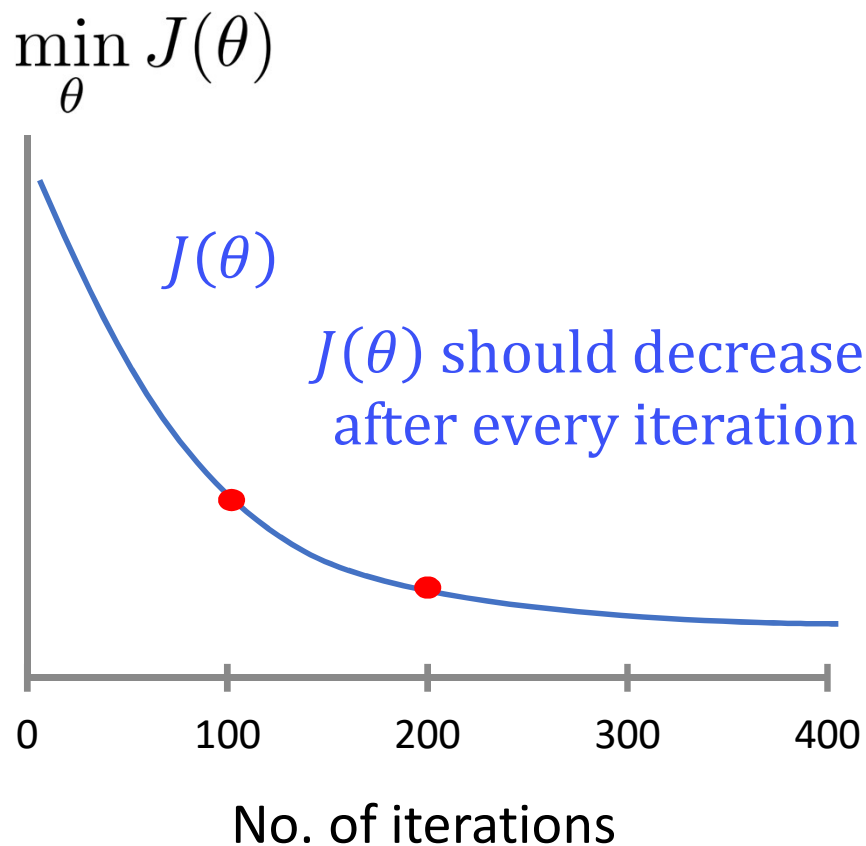
## Gradient descent

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

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



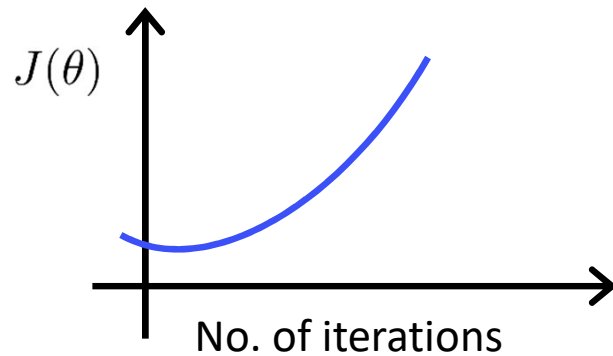
## Making sure gradient descent is working correctly.



Example automatic convergence test:

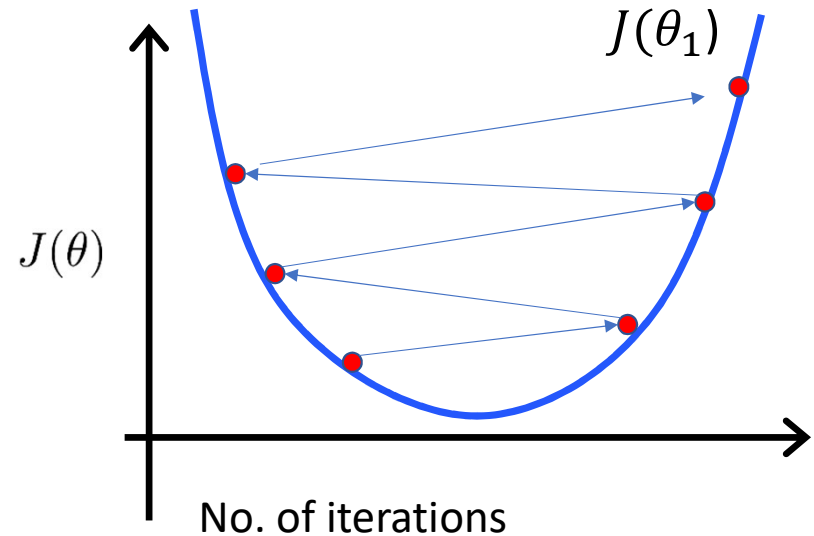
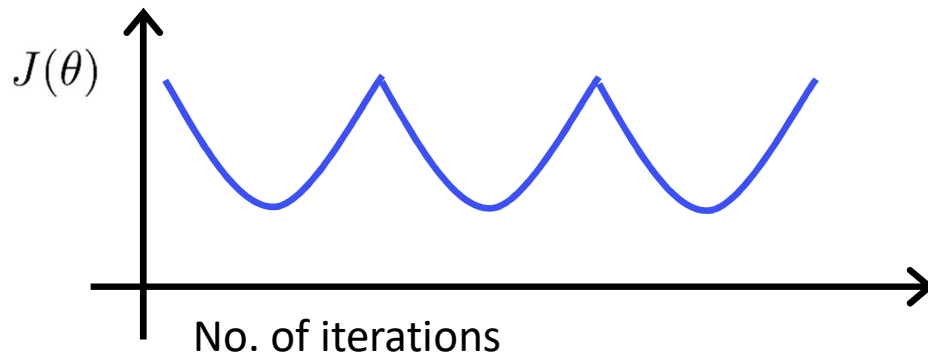
Declare convergence if  $J(\theta)$  decreases by less than  $10^{-3}$  in one iteration.

## Making sure gradient descent is working correctly.



Gradient descent not working.

Use smaller  $\alpha$ .



- For sufficiently small  $\alpha$ ,  $J(\theta)$  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:  $J(\theta)$  may not decrease on every iteration; may not converge.

To choose  $\alpha$ , try

$\dots, 0.001, 0.003, 0.01, 0.03, 0.1, 0.3, 1, \dots$

# Lesson #2: Linear regression with multiple variables

- **Duration:** 2 hrs
- **Outline:**
  - Model representation
  - Gradient descent for multiple variables
    - Feature Scaling
    - Learning rate
    - Features and polynomial regression
    - Normal equation

## Housing prices prediction

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

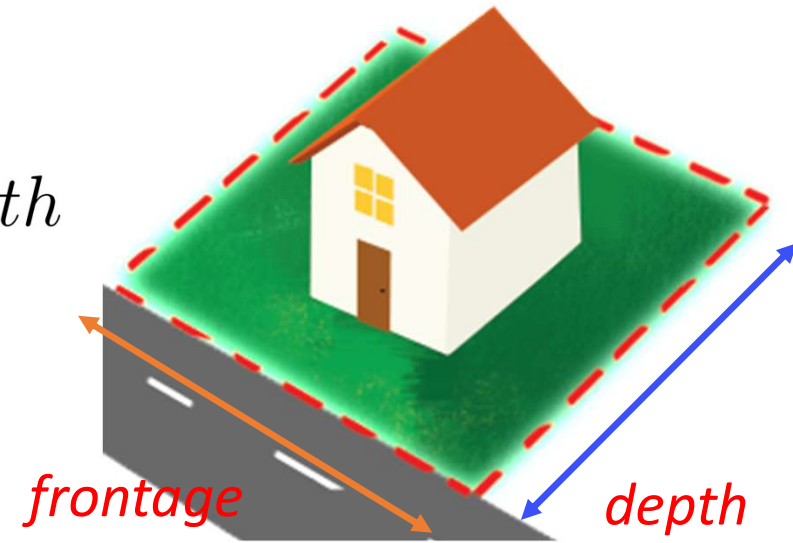
$x_1$

$x_2$

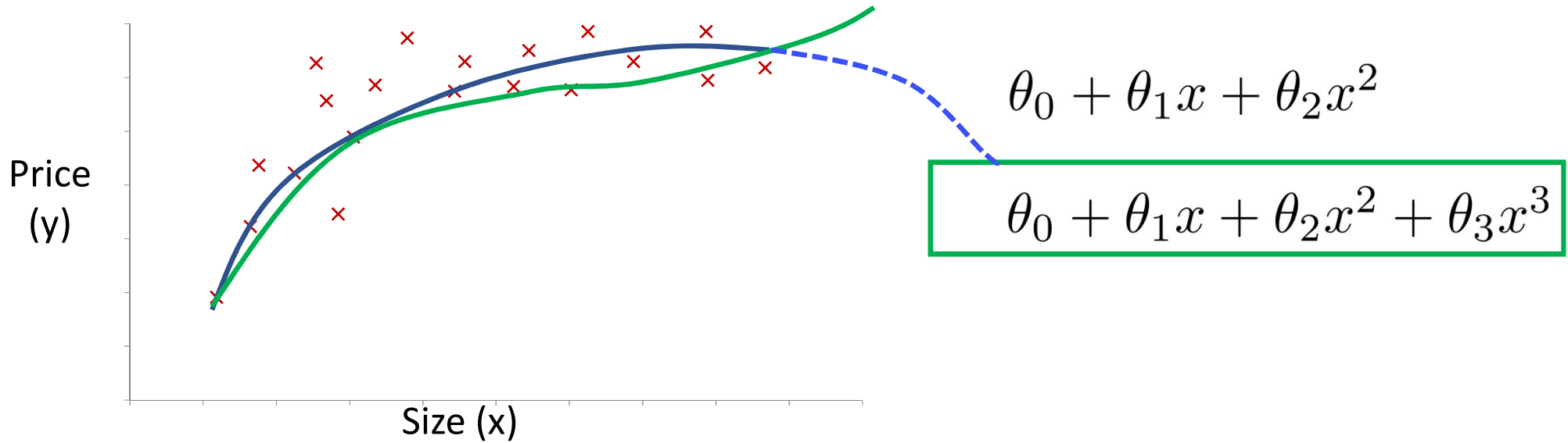
Area:

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

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



## Polynomial regression



$$\theta_0 + \theta_1 x + \theta_2 x^2$$

$$\theta_0 + \theta_1 x + \theta_2 x^2 + \theta_3 x^3$$

$$\begin{aligned} h_{\theta}(x) &= \theta_0 + \theta_1 x_1 + \theta_2 x_2 + \theta_3 x_3 \\ &= \theta_0 + \theta_1(\text{size}) + \theta_2(\text{size})^2 + \theta_3(\text{size})^3 \end{aligned}$$

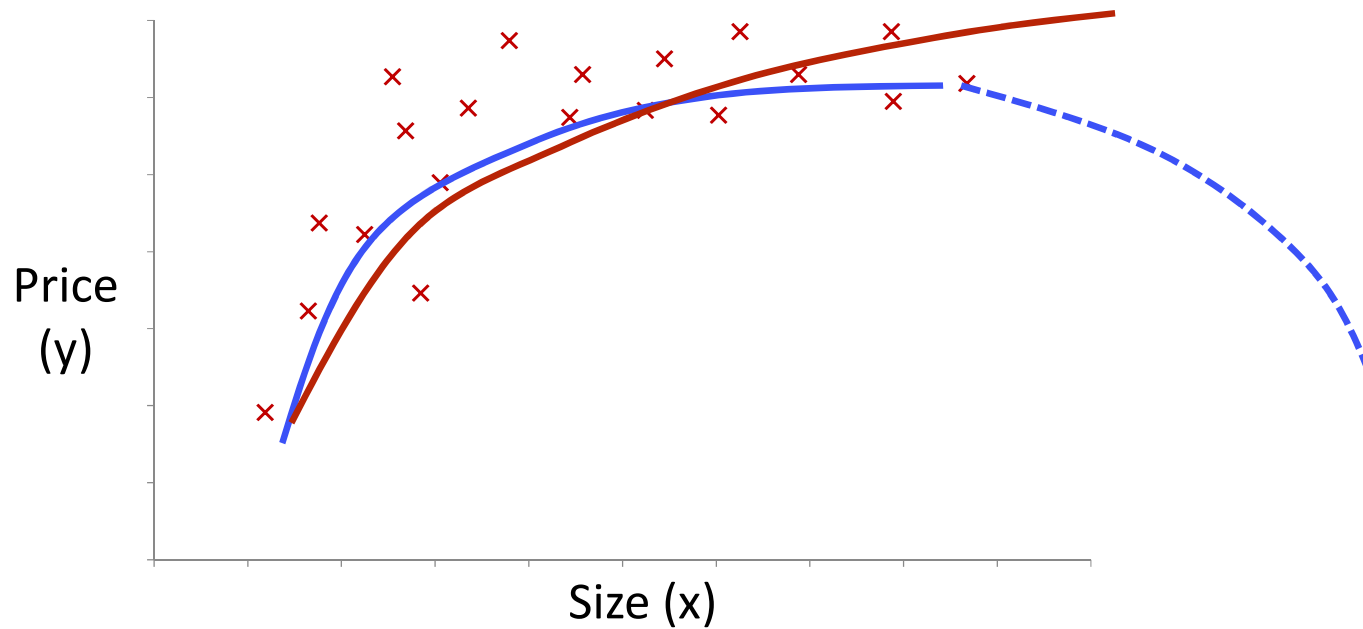
$$\begin{aligned} x_1 &= (\text{size}) \\ x_2 &= (\text{size})^2 \\ x_3 &= (\text{size})^3 \end{aligned}$$

*size* : 1 – 1000

*size*<sup>2</sup>: 1 – 1000,0000

*size*<sup>3</sup>: 1 – 10<sup>9</sup>

## Choice of features



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

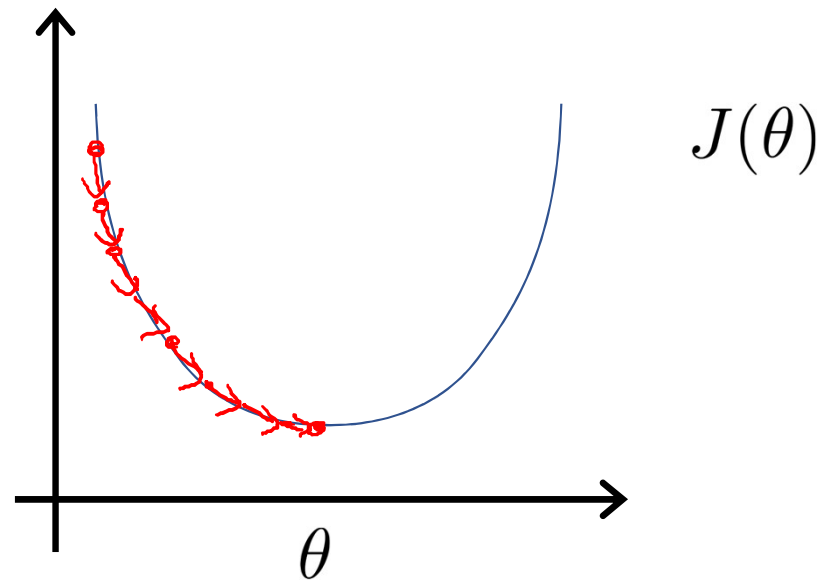
$$h_{\theta}(x) = \theta_0 + \theta_1(\text{size}) + \theta_2\sqrt{(\text{size})}$$

## Lesson #2: Linear regression with multiple variables

- **Duration:** 2 hrs
- **Outline:**
  - Model representation
  - Gradient descent for multiple variables
    - Feature Scaling
    - Learning rate
    - Features and polynomial regression
  - Normal equation



# Gradient Descent



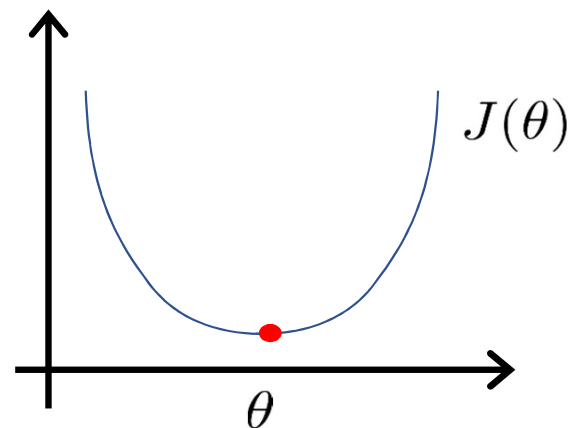
Normal equation: Method to solve for  $\theta$  analytically.

Intuition: If 1D ( $\theta \in \mathbb{R}$ )

$$J(\theta) = a\theta^2 + b\theta + c$$

$$\frac{d}{d\theta} J(\theta) = \dots = 0$$

Solve for  $\theta$



---

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

$$\frac{\partial}{\partial \theta_j} J(\theta) = \dots = 0 \quad (\text{for every } j)$$

Solve for  $\theta_0, \theta_1, \dots, \theta_n$

Examples:  $m = 4$ .

	Size (feet <sup>2</sup> )	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

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

$m \times (n + 1)$

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

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

$m$  – dimensional vector

**Examples:**  $m = 5$ .

	Size (feet <sup>2</sup> )	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

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

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

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

**$m$  examples**  $(x^{(1)}, y^{(1)}), \dots, (x^{(m)}, y^{(m)})$  ;  **$n$  features.**

$$x^{(i)} = \begin{bmatrix} x_0^{(i)} \\ x_1^{(i)} \\ x_2^{(i)} \\ \vdots \\ x_n^{(i)} \end{bmatrix} \in \mathbb{R}^{n+1}$$

$$X = \begin{bmatrix} \dots & (x^{(1)})^T & \dots \\ \dots & (x^{(2)})^T & \dots \\ & \vdots & \\ \dots & (x^{(m)})^T & \dots \end{bmatrix}$$

$m \times (n + 1)$

E.g. If  $x^{(i)} = \begin{bmatrix} 1 \\ x_1^{(i)} \end{bmatrix}$

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

$$X = \begin{bmatrix} 1 & x_1^{(1)} \\ 1 & x_2^{(1)} \\ & \vdots \\ 1 & x_m^{(1)} \end{bmatrix} \quad y = \begin{bmatrix} y^{(1)} \\ y^{(2)} \\ \vdots \\ y^{(m)} \end{bmatrix}$$

$m \times 2$

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

$(X^T X)^{-1}$  is inverse of matrix  $X^T X$ .

Set  $A = X^T X$   
 $(X^T X)^{-1} = A^{-1}$

Matlab `pinv(x' * x) * x' * y`

$$\text{pinv}(X^T * X) * X^T * y$$

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

**NO** Feature Scaling

$$0 \leq x_1 \leq 1$$

$$0 \leq x_2 \leq 1000$$

$$0 \leq x_3 \leq 10^{-5}$$

$m$  training examples,  $n$  features.

### Gradient Descent

- Need to choose  $\alpha$ .
- Needs many iterations.
- Works well even when  $n$  is large.

### Normal Equation

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

## Normal equation

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

- What if  $X^T X$  is non-invertible? (singular/degenerate)
- Matlab: `pinv(X' * X) * X' * y`  
`inv(X' * X) * X' * y`



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

- Redundant features (linearly dependent).

E.g.  $x_1 = \text{size in feet}^2$

~~$x_2 = \text{size in m}^2$~~

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

- Too many features (e.g.  $m \leq n$ ).
  - Delete some features, or use regularization.