

# Artificial Intelligence

## Lecture 10. Linear Regression

Spring 2022

Prof. Jonghoon Chun, Ph.D.

E-mail: [jchun@mju.ac.kr](mailto:jchun@mju.ac.kr)

Lecture Note: <http://lms.mju.ac.kr>

# Agenda

---

- Classification concept review
- Linear Regression
- Logistic Regression

# Classification Concept Review

---

- Supervised learning (classification)

- Supervision: The training data (observations, measurements, etc.) are accompanied by **labels** indicating the class of the observations
- New data is classified based on the training set

- Unsupervised learning (clustering)

- The class labels of training data is unknown
- Given a set of measurements, observations, etc. with the aim of establishing the existence of classes or clusters in the data

# Regression vs. Classification

---

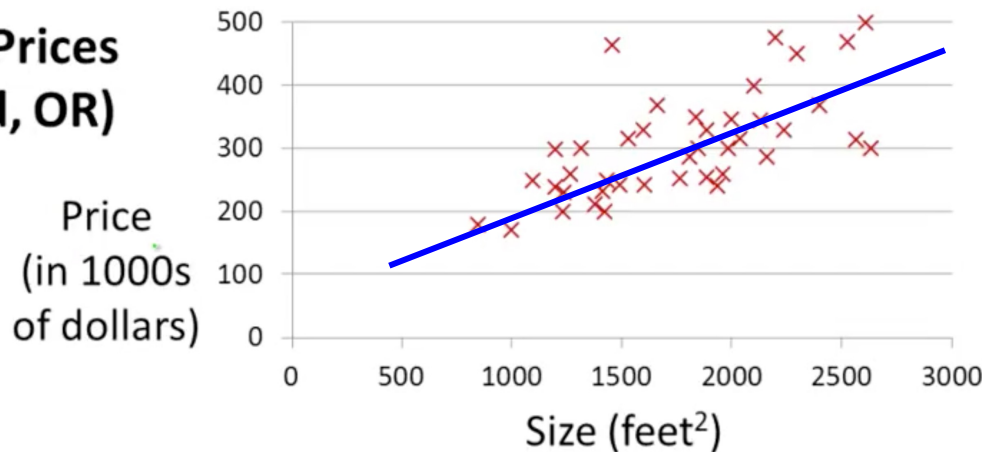
- Regression problem
  - predict results within a continuous output
  - map input variables to some continuous function
  - E.g., Given data about the size of houses on the real estate market, try to predict their price. Price as a function of size is a continuous output.
- Classification problem
  - predict results in a discrete output
  - map input variables into discrete categories (e.g., positive, negative reviews)
  - E.g., Given a patient with a tumor, predict whether the tumor is malignant or benign.

# LINEAR REGRESSION

# Regression Problem – an example

- Predict Housing Prices

## Housing Prices (Portland, OR)



### Supervised Learning

Given the “right answer” for each example in the data.

### Regression Problem

Predict real-valued output

Classification: Discrete- Valued

# Training Set

- Training set of housing prices

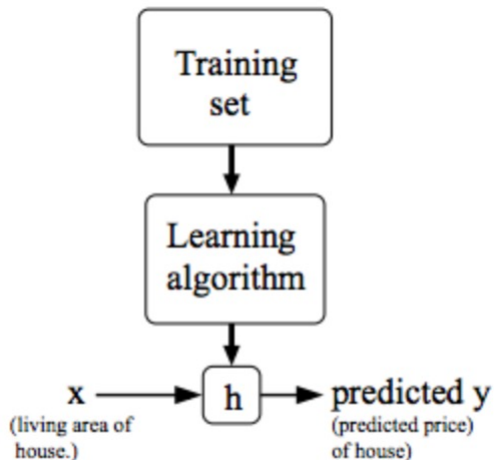
Size in feet <sup>2</sup> (x)	Price (\$) in 1000's (y)
2104	460
1416	232
1534	315
852	178
...	...

}  $m = 47$

- Notation:
  - $m$  = Number of training examples
  - $x$ 's = "input" variable / features
  - $y$ 's = "output" variable / "target" variable
  - $(x, y)$  : one training example
  - $(x^{(i)}, y^{(i)})$  :  $i^{\text{th}}$  training example

# Supervised Learning Problem

- Our goal is,
  - Given a training set, to learn a function  $h : X \rightarrow Y$  so that  $h(x)$  is a “good” predictor for the corresponding value of  $y$ .
- For historical reasons, this function  $h$  is called a **hypothesis**.



How do we represent  $h$ ?

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

Shorthand:  $h(x)$

Linear regression with one variable

Univariate linear regression

When the target variable is continuous,  
we call the learning problem a **regression** problem!



# Cost Function

- Problem to solve
  - how to fit the best possible straight line to our data

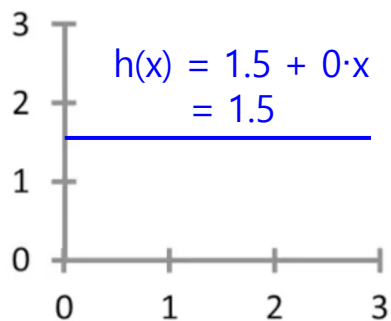
Training Set	Size in feet <sup>2</sup> (x)	Price (\$) in 1000's (y)
	2104	460
	1416	232
	1534	315
	852	178
	...	...

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

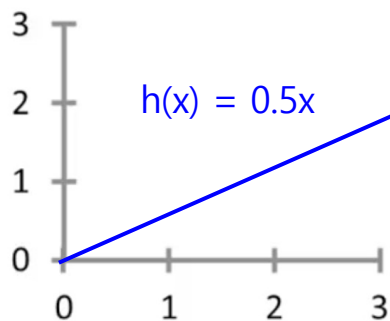
- $\theta_i$ 's : Parameters
- How to choose  $\theta_i$ 's?

# Cost Function

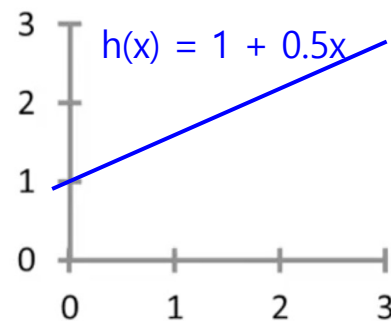
$$h_{\theta}(x) = \theta_0 + \theta_1 x$$



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

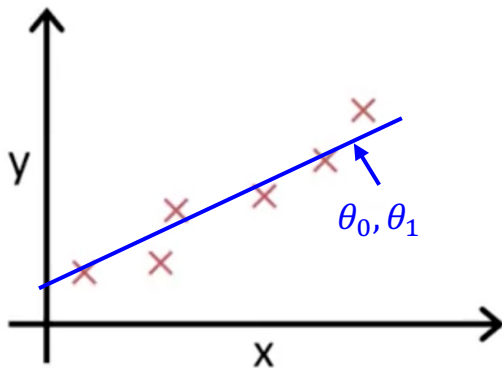


$$\theta_0 = 0$$
$$\theta_1 = 0.5$$



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

# Cost Function



$$\text{Minimize}_{\theta_0, \theta_1} \frac{1}{2m} \sum_{i=1}^m \underbrace{(h_{\theta}(x^{(i)}) - y^{(i)})^2}_{h_{\theta}(x^{(i)}) = \theta_0 + \theta_1 x^{(i)}}$$

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

$$\text{Minimize}_{\theta_0, \theta_1} \underbrace{J(\theta_0, \theta_1)}_{\substack{\text{Cost Function} \quad \text{Squared Error Function}}}$$

- Choose  $\theta_0, \theta_1$  so that  $h_{\theta}(x)$  is close to  $y$  for our training samples  $(x, y)$

# Our Goal

---

- Have some function  $J(\theta_0, \theta_1)$

We want  $\min_{\theta_0, \theta_1} J(\theta_0, \theta_1)$

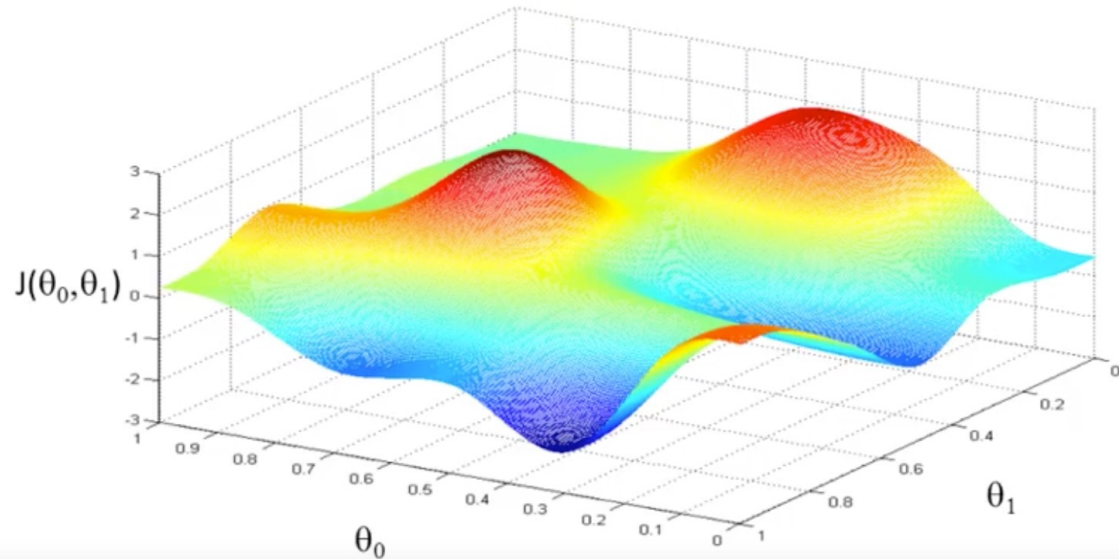
$$J(\theta_0, \theta_1, \theta_2, \dots, \theta_n)$$
$$\min_{\theta_0, \dots, \theta_n} J(\theta_0, \dots, \theta_n)$$

- Outline:

- Start with some  $\theta_0, \theta_1$  (e.g.,  $\theta_0 = 0, \theta_1 = 0$ )
- Keep changing  $\theta_0, \theta_1$  to reduce  $J(\theta_0, \theta_1)$   
until we hopefully end up at a minimum

An algorithm called "**Gradient Descent**" for minimizing the cost function  $J$  for the linear regression!

# Gradient Descent



# Gradient Descent Algorithm

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

Simultaneously update  
 $\theta_0$  and  $\theta_1$

Correct: Simultaneous update

```
temp0 :=  $\theta_0 - \alpha \frac{\partial}{\partial \theta_0} J(\theta_0, \theta_1)$   
temp1 :=  $\theta_1 - \alpha \frac{\partial}{\partial \theta_1} J(\theta_0, \theta_1)$   
 $\theta_0 :=$  temp0  
 $\theta_1 :=$  temp1
```

Incorrect:

```
temp0 :=  $\theta_0 - \alpha \frac{\partial}{\partial \theta_0} J(\theta_0, \theta_1)$   
 $\theta_0 :=$  temp0  
temp1 :=  $\theta_1 - \alpha \frac{\partial}{\partial \theta_1} J(\theta_0, \theta_1)$   
 $\theta_1 :=$  temp1
```

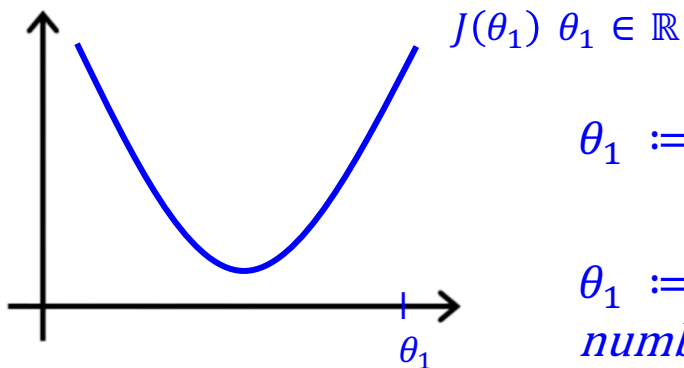
# Gradient Descent Algorithm Intuition

$$\begin{aligned} &\text{repeat until convergence } \{ \\ &\quad \theta_j := \theta_j - \alpha \frac{\partial}{\partial \theta_j} J(\theta_0, \theta_1) \quad (\text{simultaneously update} \\ &\quad \quad \quad j = 0 \text{ and } j = 1) \\ &\} \end{aligned}$$

- Assume a simpler case, we want to minimize the cost function with just one parameter

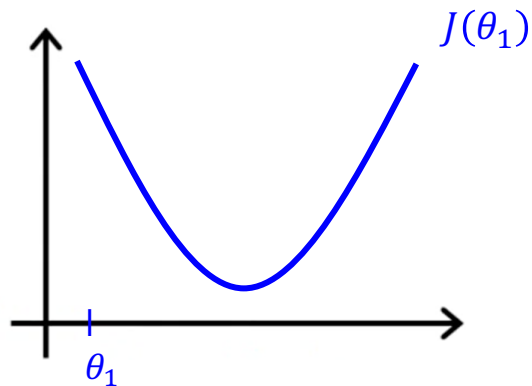
$$\min_{\theta_1} J(\theta_1), \quad \theta_1 \in \mathbb{R}$$

# Gradient Descent Algorithm Intuition



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

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



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

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

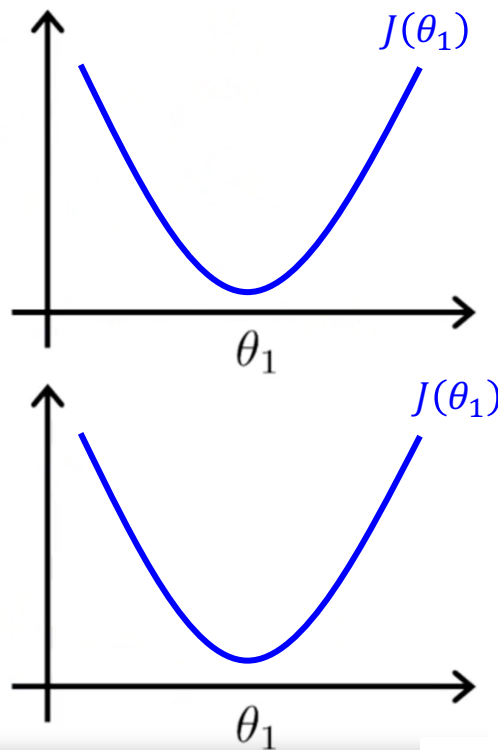


# Learning Rate

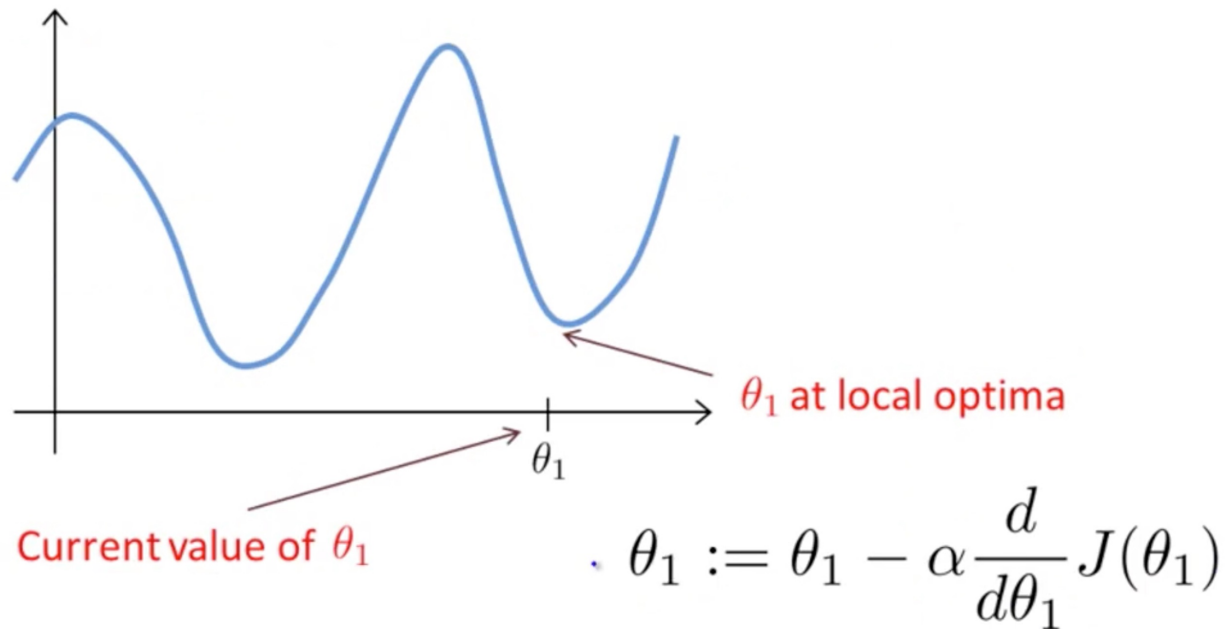
$$\theta_1 := \theta_1 - \alpha \frac{\partial}{\partial \theta_1} J(\theta_1)$$

If  $\alpha$  is too small, gradient descent can be slow.

If  $\alpha$  is too large, gradient descent can overshoot the minimum. It may fail to converge, or even diverge.



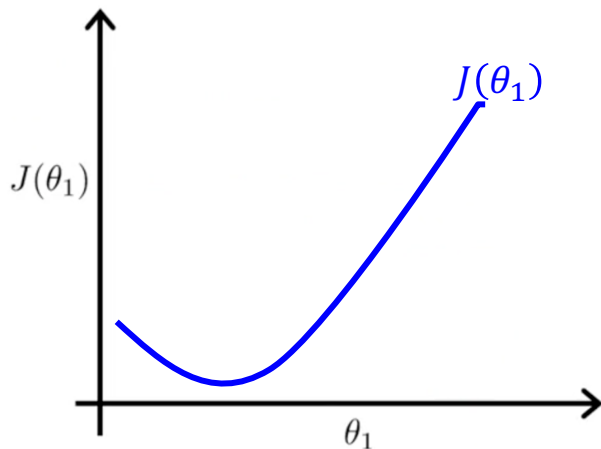
# Local Minimum Problem



# Local Minimum Problem

- Gradient descent can converge to a local minimum, even with the learning rate  $\alpha$  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  $\alpha$  over time.

# Gradient Descent for Linear Regression

Gradient descent algorithm

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

Linear Regression Model

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

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

$$\min_{\theta_0, \theta_1} J(\theta_0, \theta_1)$$

# Gradient Descent for Linear Regression

---

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

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

# Gradient Descent for Linear Regression

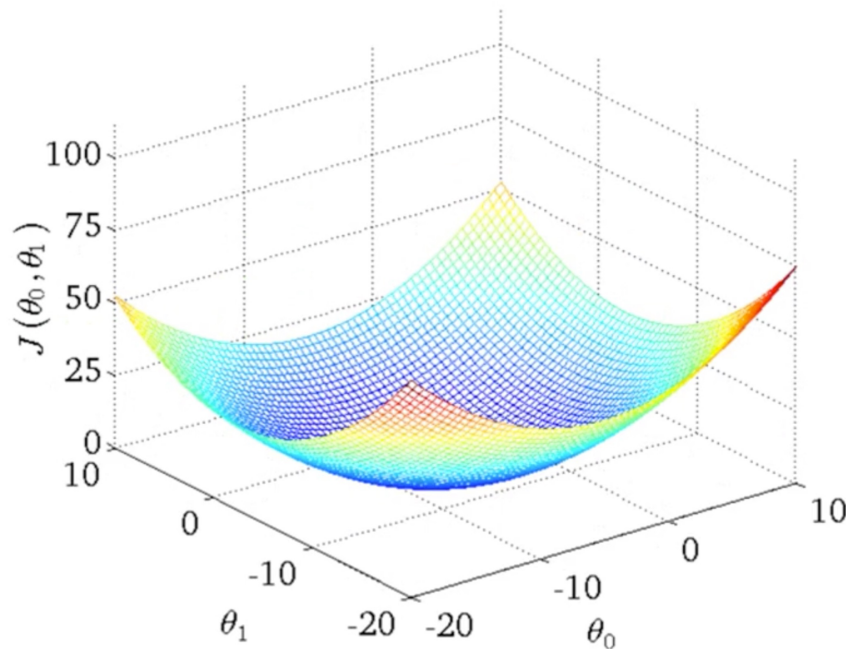
$$\begin{aligned} &\text{repeat until convergence } \{ \\ &\quad \theta_0 := \theta_0 - \alpha \frac{1}{m} \sum_{i=1}^m (h_{\theta}(x^{(i)}) - y^{(i)}) \\ &\quad \theta_1 := \theta_1 - \alpha \frac{1}{m} \sum_{i=1}^m (h_{\theta}(x^{(i)}) - y^{(i)}) \cdot x^{(i)} \\ &\quad \} \end{aligned}$$

update  
 $\theta_0$  and  $\theta_1$   
Simultaneously

## "Batch" Gradient Descent

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

# Gradient Descent for Linear Regression



Cost function for linear regression is always going to be a  
**"Convex Function"!**

# Remarks

---

- Solving for the minimum of the cost function
  - Gradient descent: an iterative algorithm
  - **Normal equation method**: numerically solving for the minimum of the cost function without needing the multiple steps of gradient descent
- Gradient descent will scale better to larger data sets than that normal equation method.
- Generalization of gradient descent algorithm
  - Will make it much more powerful



# MULTIVARIATE LINEAR REGRESSION

# 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
...	...	...	...	...

## ■ 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

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

# Multivariate Linear Regression

---

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

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

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

- General form:

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

For convenience of notation, define  $x_0 = 1$

$$\begin{aligned} h_{\theta}(x) &= \theta_0 x_0 + \theta_1 x_1 + \theta_2 x_2 + \dots + \theta_n x_n \\ &= \theta^T x \end{aligned}$$

# Multivariate Linear Regression

$$h_{\theta}(x) = \theta^T x$$

$$x = \begin{bmatrix} x_0 \\ x_1 \\ x_2 \\ \vdots \\ x_n \end{bmatrix} \in \mathbb{R}^{n+1} \quad \theta = \begin{bmatrix} \theta_0 \\ \theta_1 \\ \theta_2 \\ \vdots \\ \theta_n \end{bmatrix} \in \mathbb{R}^{n+1}$$

$\underbrace{[\theta_0 \ \theta_1 \ \dots \ \theta_n]}_{\theta^T}$

$\begin{bmatrix} x_0 \\ x_1 \\ x_2 \\ \vdots \\ x_n \end{bmatrix}$

(n + 1) X 1 matrix

# 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 {

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

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

}

---


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

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

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

...

## ACKNOWLEDGMENTS

- Many of the slides and examples are borrowed from the open course by Coursera "**Machine Learning**" by A. Ng which were modified into their current form.

END