

# Regression Analysis

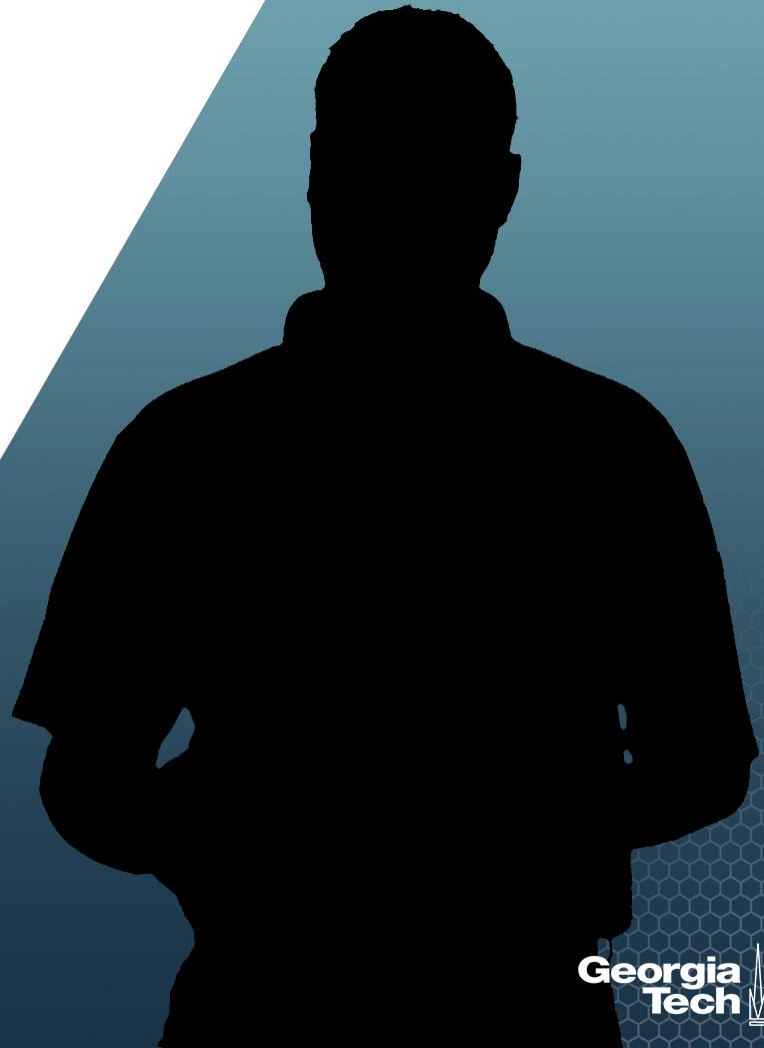
## Multiple Linear Regression

**Nicoleta Serban, Ph.D.**

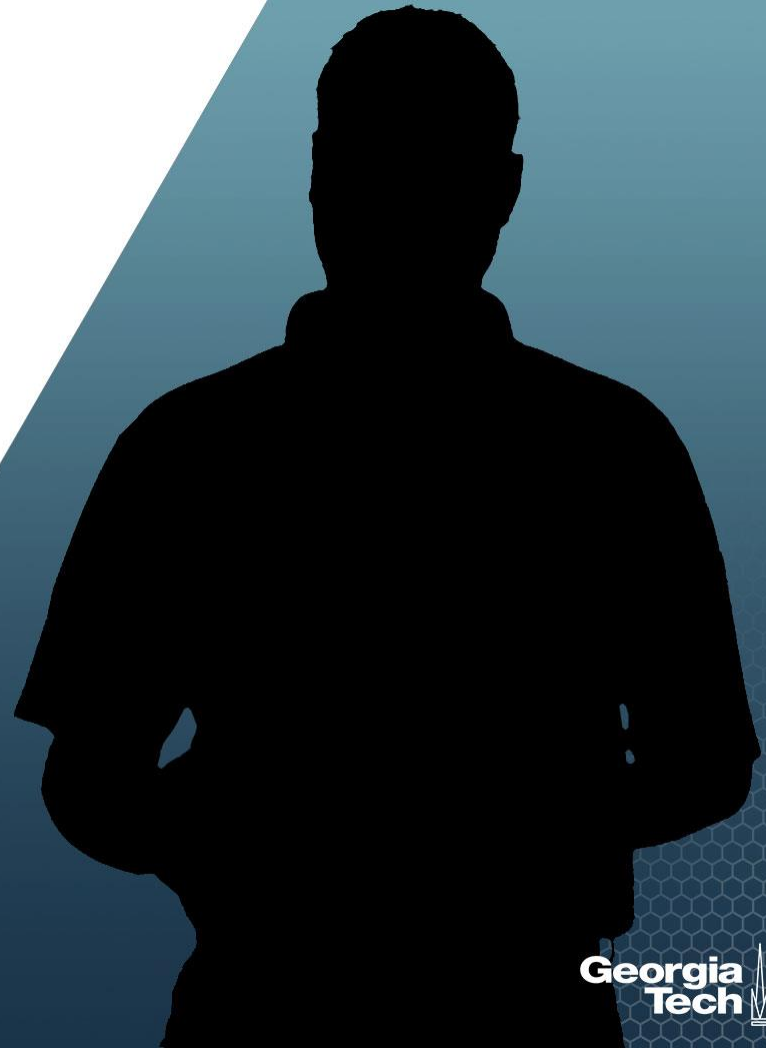
*Professor*

School of Industrial and Systems Engineering

Basics of Multiple Regression



# About This Lesson



# Multiple Linear Regression: Model

**Data:**  $\{(x_{1,1}, \dots, x_{1,p}), y_1\}, \dots, \{(x_{n,1}, \dots, x_{n,p}), y_n\}$

**Model:**  $Y_i = \beta_0 + \beta_1 x_{i,1} + \beta_2 x_{i,2} + \dots + \beta_p x_{i,p} + \varepsilon_i, \quad i = 1, \dots, n$

## Assumptions:

- *Linearity/Mean Zero Assumption:*  $E(\varepsilon_i) = 0$
- *Constant Variance Assumption:*  $\text{Var}(\varepsilon_i) = \sigma^2$
- *Independence Assumption:*  $\{\varepsilon_1, \dots, \varepsilon_n\}$  are independent random variables
- $\varepsilon_i \sim$  Normally distributed *for confidence/prediction intervals, hypothesis testing*

# Multiple Linear Regression: Model

**Data:**  $\{(x_{1,1}, \dots, x_{1,p}), y_1\}, \dots, \{(x_{n,1}, \dots, x_{n,p}), y_n\}$

**Model:**  $Y_i = \beta_0 + \beta_1 x_{i,1} + \beta_2 x_{i,2} + \dots + \beta_p x_{i,p} + \varepsilon_i, \quad i = 1, \dots, n$

## Assumptions:

- *Linearity/Mean Zero Assumption:*  $E(\varepsilon_i) = 0$
- *Constant Variance Assumption:*  $\text{Var}(\varepsilon_i) = \sigma^2$
- *Independence Assumption:*  $\{\varepsilon_1, \dots, \varepsilon_n\}$  are independent random variables
- $\varepsilon_i \sim$  Normally distributed *for confidence/prediction intervals, hypothesis testing*

**The model parameters are:**  $\beta_0, \beta_1, \dots, \beta_p, \sigma^2$

- **Unknown** regardless how much data are observed
- **Estimated** given the model assumptions
- **Estimated** based on data

# Multiple Linear Regression: Model

**Data:**  $\{(x_{1,1}, \dots, x_{1,p}), y_1\}, \dots, \{(x_{n,1}, \dots, x_{n,p}), y_n\}$

**Model:**  $Y_i = \beta_0 + \beta_1 x_{i,1} + \beta_2 x_{i,2} + \dots + \beta_p x_{i,p} + \varepsilon_i, \quad i = 1, \dots, n$

**Model in Matrix Form:**  $Y = X\beta + \varepsilon$

**Response**

$$Y = \begin{pmatrix} Y_1 \\ Y_2 \\ \vdots \\ Y_n \end{pmatrix}$$

**Design Matrix**

$$X = \begin{bmatrix} 1 & x_{1,1} & x_{1,2} & \cdots & x_{1,p} \\ 1 & x_{2,1} & x_{2,2} & \cdots & x_{2,p} \\ \vdots & \vdots & \vdots & \ddots & \vdots \\ 1 & x_{n,1} & x_{n,2} & \cdots & x_{n,p} \end{bmatrix}$$

**Coefficients**

$$\beta = \begin{pmatrix} \beta_0 \\ \beta_1 \\ \vdots \\ \beta_p \end{pmatrix}$$

**Error**

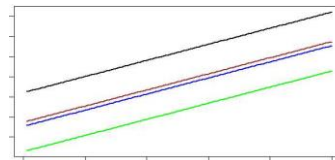
$$\varepsilon = \begin{pmatrix} \varepsilon_1 \\ \varepsilon_2 \\ \vdots \\ \varepsilon_n \end{pmatrix}$$

# Model Flexibility: Main Effects & Interactions

For  $k = 2$  predicting variables, four useful regressions:

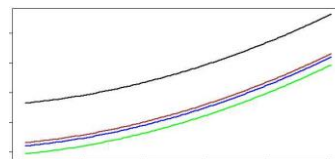
- **1<sup>st</sup> Order Model:**

$$Y = \beta_0 + \beta_1 x_1 + \beta_2 x_2 + \varepsilon$$



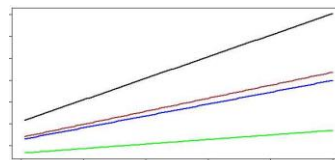
- **2<sup>nd</sup> Order Model:**

$$Y = \beta_0 + \beta_1 x_1 + \beta_2 x_2 + \beta_3 x_1^2 + \beta_4 x_2^2 + \varepsilon$$



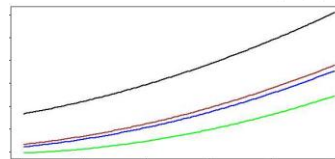
- **1<sup>st</sup> Order Interaction Model:**

$$Y = \beta_0 + \beta_1 x_1 + \beta_2 x_2 + \beta_3 x_1 x_2 + \varepsilon$$



- **2<sup>nd</sup> Order Interaction Model:**

$$Y = \beta_0 + \beta_1 x_1 + \beta_2 x_2 + \beta_3 x_1 x_2 + \beta_4 x_1^2 + \beta_5 x_2^2 + \varepsilon$$



# Quantitative and Qualitative Variables

**Simple Linear Regression:** Linear regression with one quantitative predicting variable

**ANOVA:** Linear regression with one or more qualitative predicting variables

**Multiple Linear Regression:** Multiple quantitative and qualitative predicting variables

# Quantitative and Qualitative Variables

**Multiple Linear Regression:** Multiple quantitative/qualitative predicting variables

$x_1$  quantitative

$x_2$  qualitative with three levels:  $D_1$ ,  $D_2$ , and  $D_3$  dummy variables

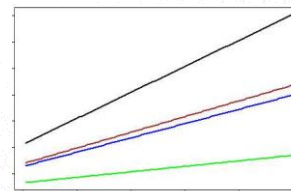
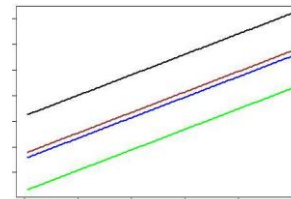
Model:  $y = \beta_0 + \beta_1 x_1 + \beta_2 d_1 + \beta_3 d_2 + \varepsilon$  ➡ **Intercept varies**

If  $d_1=0, d_2=0$ :  $\beta_0 + \beta_1 x_1$

If  $d_1=1, d_2=0$ :  $(\beta_0 + \beta_2) + \beta_1 x_1$

If  $d_1=0, d_2=1$ :  $(\beta_0 + \beta_3) + \beta_1 x_1$

**Parallel regression lines**



**If  $x_1$   $x_2$  interaction: Nonparallel regression lines**



# Linear Regression: Example 1



# Linear Regression: Example 1

## Quantitative Predicting Variables:

$X_1$  = The amount (in hundreds of dollars) spent on advertising in 1999

$X_2$  = The total amount of bonuses paid in 1999

$X_3$  = The market share in each territory

$X_4$  = The largest competitor's sales

## Qualitative Predicting Variable:

$X_5$  = Indicates the region of the office (1 = south, 2 = west, 3 = midwest)

# Linear Regression: Example 3



**Bike sharing systems are of great interest due to their important role in traffic management.**

**Dataset:** Historical data for years 2011-2012 for the bike sharing system in Washington D.C.

# Linear Regression: Example 3

## Qualitative predicting variables:

$X_1$  = Day of the week

$X_2$  = Month of the year

$X_3$  = Hour of the day (ranging 0-23)

$X_4$  = Year (2011, 2012)

$X_5$  = Holiday Indicator

$X_6$  = Weather condition (with four levels  
from good weather for level 1 to  
severe condition for level 4)

## Quantitative predicting variables:

$X_7$  = Normalized temperature

$X_8$  = Normalized humidity

$X_9$  = Wind speed

# Linear Regression: Example 3

## Qualitative predicting variables:

$X_1$  = Day of the week

$X_2$  = Month of the year

$X_3$  = Hour of the day (ranging 0-23)

$X_4$  = Year (2011, 2012)

$X_5$  = Holiday Indicator

$X_6$  = Weather condition (with four levels from good weather for level 1 to severe condition for level 4)

## Quantitative predicting variables:

$X_7$  = Normalized temperature

$X_8$  = Normalized humidity

$X_9$  = Wind speed

**Year:** A quantitative or a qualitative predicting variable?

- If observations are made over many years, then consider it to be *quantitative*
- If observations are made over only a few years, then consider it to be *qualitative*

# Summary

