

## Regression Analysis

Often in statistics, we use sample data to investigate the relationship among variables. Finding a mathematical model that best describes the relationship between a variable (**response or dependent** variable) and other variables (**regressors, predictor or independent** variables) whose values we believe can be used to predict that response variable is considered a key task in regression analysis.

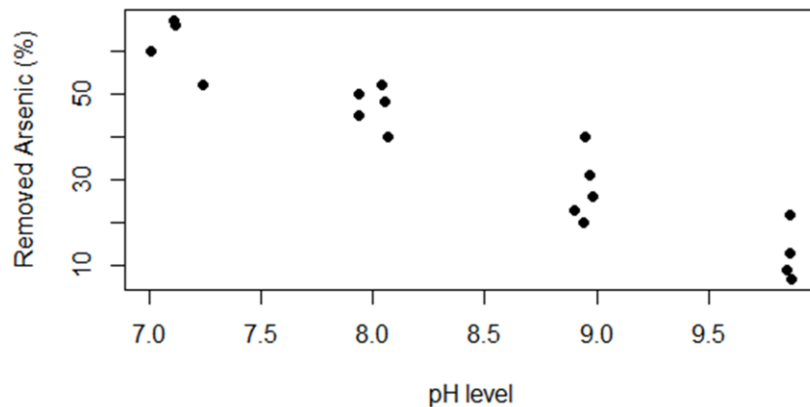
The variable to be predicted is called the **response variable**.

The variable(s) that are used to predict the response variable are called **predictor variable(s)**.

### Example 1:

Arsenic is found in many ground-waters and some surface waters. Recent health effects research has prompted the Environmental Protection Agency to reduce allowable arsenic levels in drinking water so that many water systems are no longer compliant with standards. This has spurred interest in the development of methods to remove arsenic. The accompanying data on  $x = \text{pH}$  and  $y = \text{arsenic removed (\%)}$  by a particular process was read from a scatter plot in the article “Optimizing Arsenic Removal During Iron Removal: Theoretical and Practical Considerations” (J. of Water Supply Res. and Tech., 2005: 545–560).

$y$	$x$
60	7.01
67	7.11
66	7.12
52	7.24
50	7.94
45	7.94
52	8.04
48	8.05
40	8.07
23	8.90
20	8.94
40	8.95
31	8.97
26	8.98
9	9.85
22	9.86
13	9.86
7	9.87



Scatterplot of arsenic percentage removed and pH level of water  
The above scatterplot shows the relationship between the arsenic percentage removed and the pH level of water.

### Main Goals of Regression Analysis

- ❖ Describe the relationship between predictor(s) and response variable
- ❖ Based on the identified relationship, predict the y-value for given x-value

Linear regression model is

$$y = \beta_0 + \beta_1 x_1 + \beta_2 x_2 + \beta_3 x_3 + \cdots + \beta_k x_k + \varepsilon$$

$y$  is response variable

$x_1, x_2, \dots, x_k$  are the predictor variables

$\beta_0, \beta_1, \dots, \beta_k$  are regression parameters (coefficients)

$\varepsilon$  is the error term

This is a **multiple linear regression** model because it has more than one predictor.

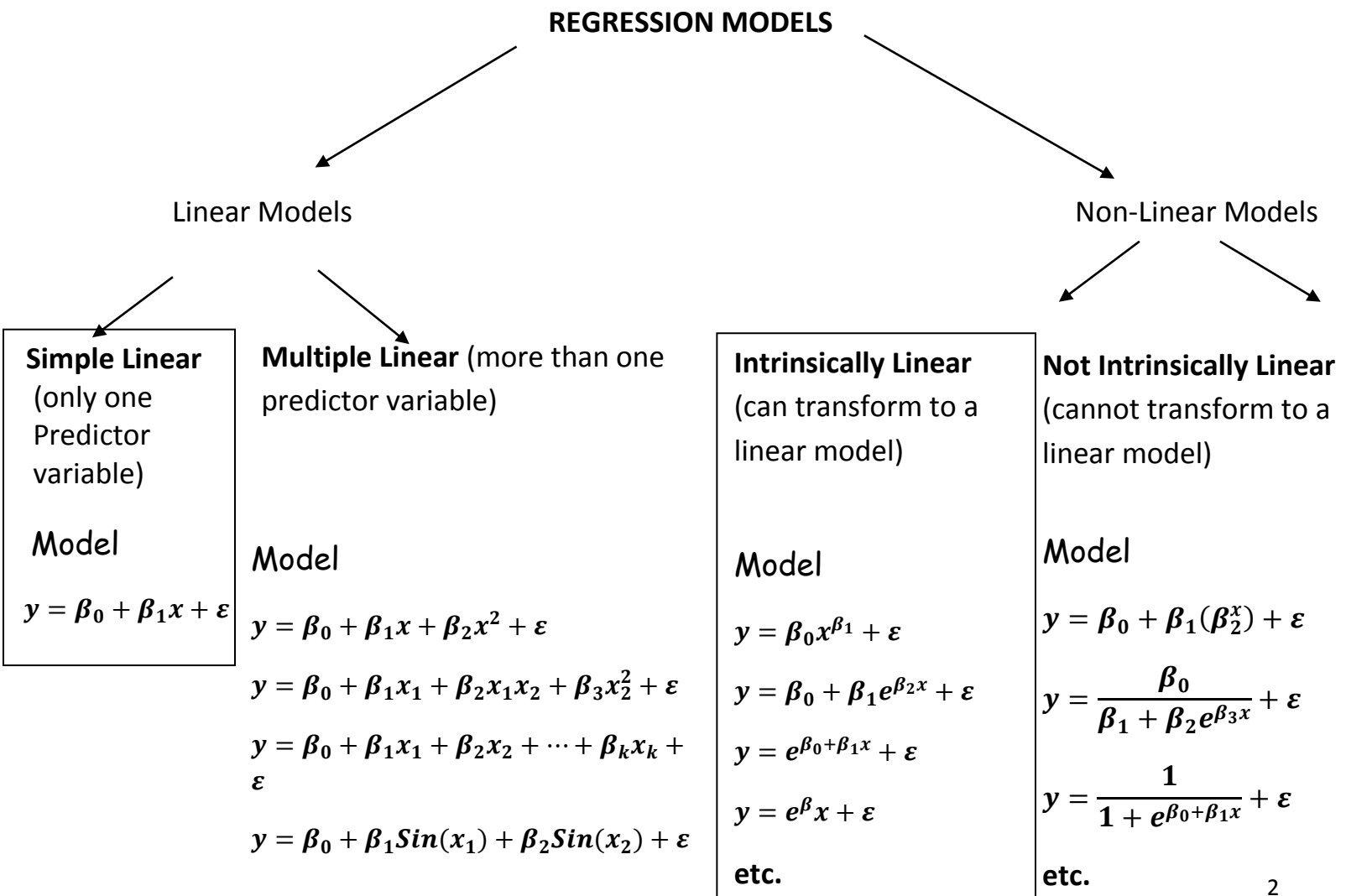
The **linear** regression models have linearity in the parameters.

For example,

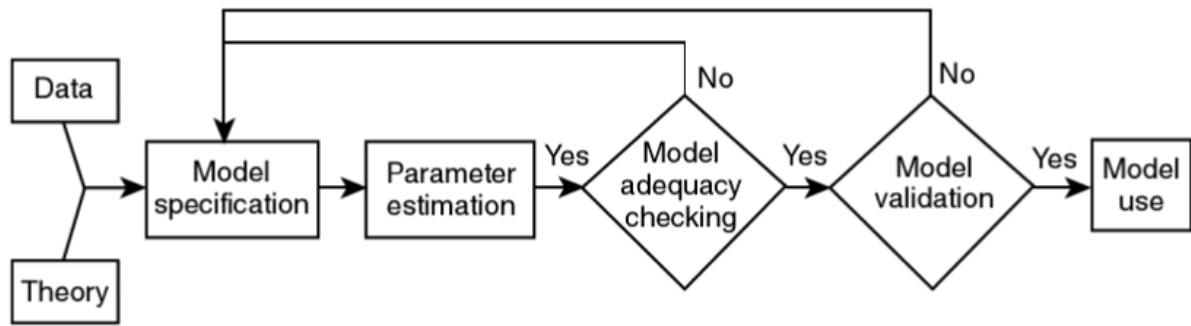
is a **linear** regression model.

However,

is **NOT a linear** regression model. This is non-linear in  $\beta$ .



## Regression model building process



**Figure:** Regression model building process.

### Simple Linear Regression

Only one predictor variable

Relationship is linear (straight line)

The **simple linear** regression model is a regression model that has a **straight line** relationship between the response variable and the predictor variable, and **only one predictor** variable.

#### **Model:**

$$y = \beta_0 + \beta_1 x + \varepsilon$$

$y$  is response variable

$x$  is the predictor variable

$\beta_0, \beta_1$  are regression parameters (coefficients)

$\varepsilon$  is the random error

#### **Assumptions:**

1. The relationship between the response and predictor is linear.
2. The error term has mean zero. That is,  $E(\varepsilon) = 0$
3. The error term has a constant variance for each given  $x$  value. That is,  $Var(\varepsilon) = \sigma^2$ .  
This is called homoscedasticity (homogeneity of variance)
4. The error terms are uncorrelated (independent). That is,  $Cov(\varepsilon_i, \varepsilon_j) = 0$  for  $i \neq j$
5. The error terms have a normal distribution.

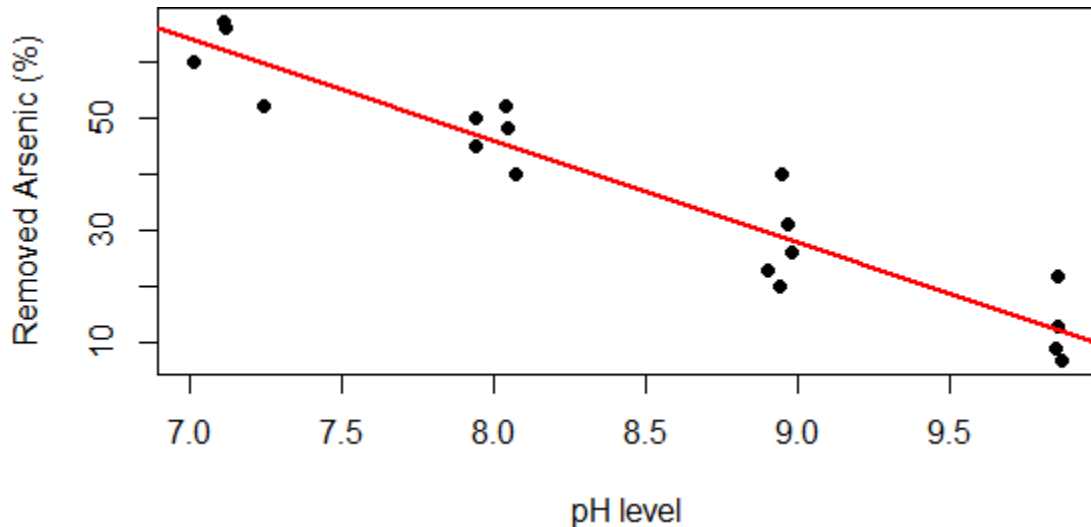
With these assumptions about error, the simple linear regression model for  $n$  number of observations can be written as

### Estimate Regression Parameters (or Coefficients)

There is more than one straight line that can be used to explain a linear relationship between response and predictors. We need a measure of how close an estimated regression line is to the data points.

The **method of least squares** chooses the estimates,  $\widehat{\beta}_0$  and  $\widehat{\beta}_1$  to minimize the sum of squared vertical deviations from the line. The regression line  $\hat{y} = \widehat{\beta}_0 + \widehat{\beta}_1 x$  is the **fitted least squares regression line**.

$\hat{y}_i = \widehat{\beta}_0 + \widehat{\beta}_1 x_i$  is called the **fitted value** for the  $i^{th}$  observation.



**Figure:** Fitted regression line overlaid on the scatterplot of percentage of arsenic removed and pH level of water.

Vertical deviation =

This is called residual (or observed error).

In least squares, we minimize