

## MTH 445/545

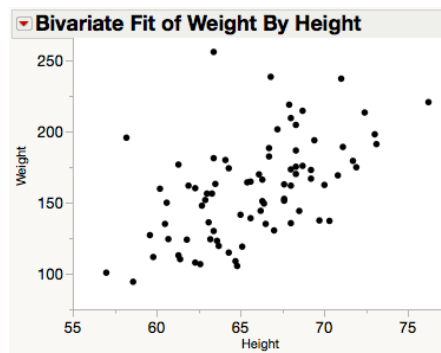
### Linear Regression and Correlation

- Regression analysis is a statistical tool that utilizes the relation between two or more quantitative variables so that one variable can be predicted from the other, or others.
- Some Examples:
  - Height and weight of people
  - Income and expenses of people
  - Production size and production time
  - Soil pH and the rate of growth of plants

1

## Correlation

- An easy way to determine if two quantitative variables are linearly related is by looking at their scatterplot.
- Another way is to calculate the correlation coefficient, denoted usually by  $r$ .
- The **Linear Correlation** measures the strength of the linear relationship between explanatory variable ( $x$ ) and the response variable ( $y$ ). An estimate of this correlation parameter is provided by the Pearson sample correlation coefficient,  $r$ .

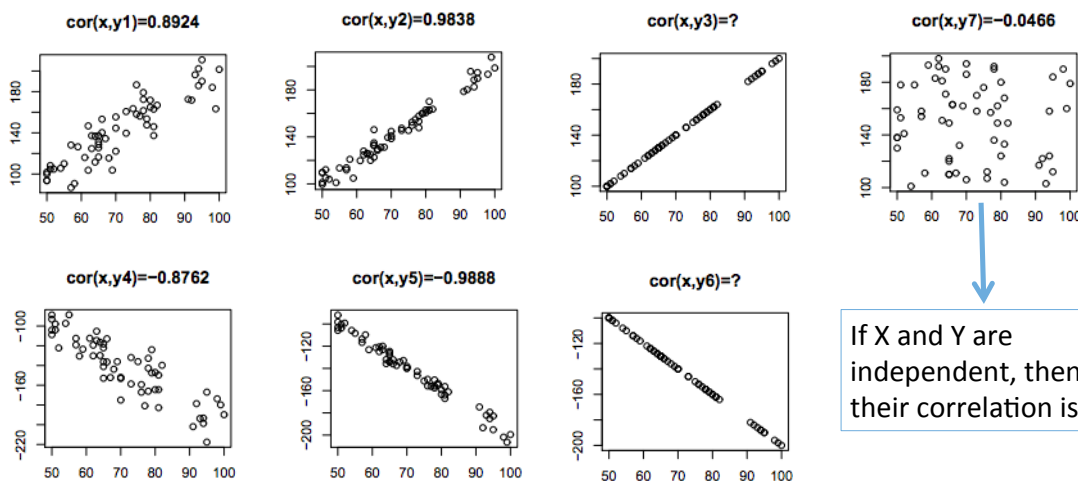


$$r = \frac{\sum (x_i - \bar{x})(y_i - \bar{y})}{\sqrt{\sum (x_i - \bar{x})^2 \sum (y_i - \bar{y})^2}}$$

Note:  $-1 \leq r \leq 1$ .

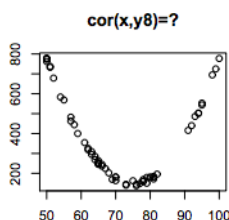
2

## Example Scatterplots with Correlations



3

## Correlation



If the correlation between X and Y is 0, it doesn't mean they are independent. It only means that they are not linearly related.

- Some Guidelines in Interpreting r.

One complain about the correlation is that it can be subjective when interpreting its value. Some people are very happy with  $r \approx 0.6$ , while others are not.

**Note: Correlation does not necessarily imply Causation!**

Value of $ r $	Strength of linear relationship
If $ r  \geq 0.95$	Very Strong
If $0.85 \leq  r  < 0.95$	Strong
If $0.65 \leq  r  < 0.85$	Moderate to Strong
If $0.45 \leq  r  < 0.65$	Moderate
If $0.25 \leq  r  < 0.45$	Weak
If $ r  < 0.25$	Very weak/Close to none

4

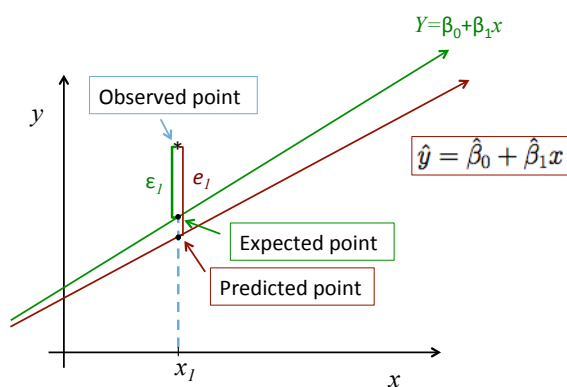


## Model Assumptions

- Model:  $Y_i = (\beta_0 + \beta_1 x_i) + \varepsilon_i$

where,

- $\varepsilon_i$ 's are uncorrelated with a mean of 0 and constant variance  $\sigma^2_{\varepsilon}$ .
- $\varepsilon_i$ 's are normally distributed. (This is needed in the test for the slope.)



Since the underlying (green) line is unknown to us, we can't calculate the values of the error terms ( $\varepsilon_i$ ). The best that we can do is study the residuals ( $e_i$ ).