# STATISTICS FOR DATA SCIENCE

**Satya Vani NL**

Department of Science & Humanities

# STATISTICS FOR DATA SCIENCE

**Uncertainties in Least Squares Coefficients**
**The More Spread in the x Values , the Better**

**Satya Vani NL**

Department of Science & Humanities

Consider Bivariate data $(x_i, y_i)$ for $i=1,2,3,\ldots\ldots\ldots\ldots n$

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

The line $y_i = \beta_0 + \beta_1 x_i + \varepsilon_i$, $\varepsilon_i$ is the error, that best fits the data in the sense of minimizing the sum of the squared errors. It is called the least squares regression line
$\widehat{\beta_0}$, $\widehat{\beta_1}$ are estimates of $\beta_0$, $\beta_1$.

$\widehat{y}_i = \widehat{\beta_0} + \widehat{\beta_1} x_i$ = Fitted line

If $\varepsilon_i$ tend to be large then $(x_i, y_i)$ are widely scattered around the line.
If $\varepsilon_i$ tend to be small then $(x_i, y_i)$ are tightly clustered around the line.

**STATISTICS FOR DATA SCIENCE**

$\widehat{\beta_0}$ , $\widehat{\beta_1}$ are called Least Squares Coefficients and defined as

$$\widehat{\beta_1} = \sum_{i=1}^{n} \left[ \frac{x_i - \bar{x}}{\sum_{i-1}^{n}(x_i - \bar{x})^2} \right] y_i$$

$$\widehat{\beta_0} = \sum_{i=1}^{n} \left[ \frac{1}{n} - \frac{\bar{x}(x_i - \bar{x})}{\sum_{i-1}^{n}(x_i - \bar{x})^2} \right] y_i$$

This indicates that $\widehat{\beta_0}$ , $\widehat{\beta_1}$ are linear combination of $y_i$ .
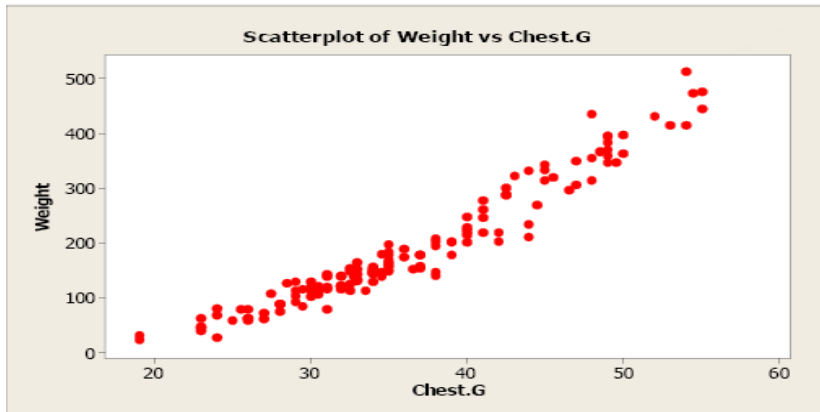
$$S_{\widehat{\beta_0}} = S \sqrt{\left[\frac{1}{n} + \frac{\overline{x}^2}{\sum_{i-1}^{n}(x_i - \overline{x})^2}\right]} \qquad S_{\widehat{\beta_1}} = S \sqrt{\left[\frac{1}{\sum_{i-1}^{n}(x_i - \overline{x})^2}\right]}$$

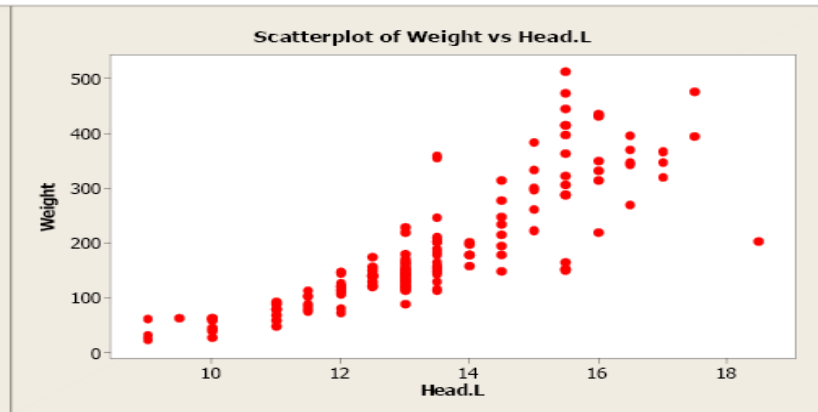$$S_{\widehat{\beta_1}} \ \alpha \ \frac{1}{\sum_{i-1}^{n}(x_i - \overline{x})^2}$$

If $x-$ values are more spread then the uncertainty of estimates $\widehat{\beta_0}, \widehat{\beta_1}$ $are$ Smaller.
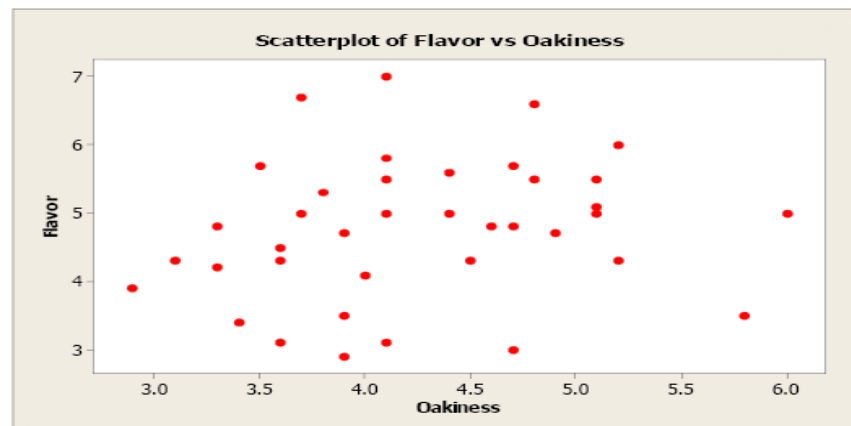
The standard deviation of x is more.

Scatterplot of Weight vs Chest.G

**Strong positive relationship**
**r = 0.96**



Scatterplot of Weight vs Head.L

**Moderate positive relationship**
**r = 0.67**



Scatterplot of Flavor vs Oakiness

**Very weak positive relationship**
**r = 0.07**

Problem: Two engineers are conducting independent experiments to estimate spring constant for a particular spring. The first engineer suggests measuring the length of the spring with no load, then applying loads of 0,1,2,3,& 4 lb. The second engineer suggests using loads of 0, 2, 4, 6 & 8 lb. Which will be more precise?

Sol: X ----- 0, 1, 2, 3, 4                    Y-------- 0, 2, 4, 6, 8

$\sigma_y$ is twice as great as $\sigma_x$ .

Uncertainty of X is twice as large as the uncertainty of Y. Hence, the Engineer, Y 's estimate is twice as precise.

# THANK YOU

**Satya Vani NL**

Department of Science & Humanities
sathyavaninl@pes.edu

+91 80 66186410