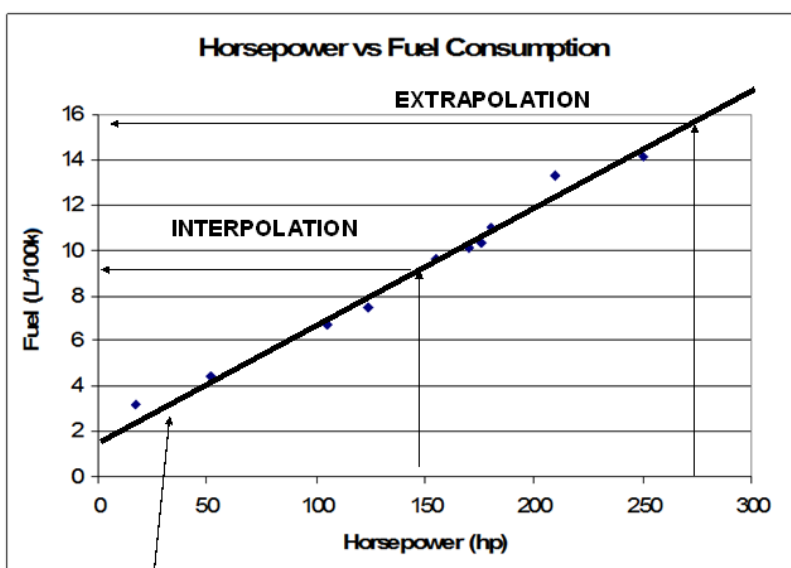


3.2 LINEAR REGRESSION

- Regression is an analytical technique for modeling the relationship between two variables.
- When the relationship between the two variables is **linear**, a simple mathematical model can be developed by finding **the line of best fit**.
- The equation of this line of best fit can be used to make predictions by **Interpolation** (estimates within the given data range) and **Extrapolation** (estimates outside of the given data range)

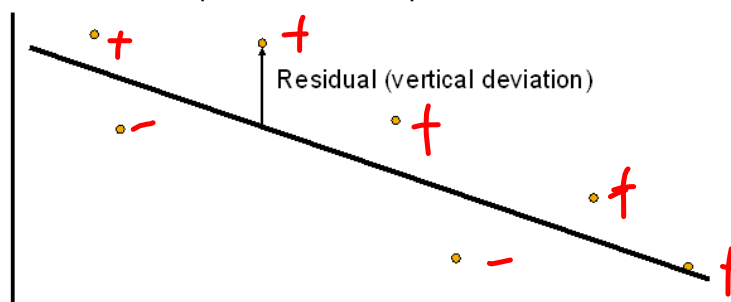
3.2 LINEAR REGRESSION



LINE OF BEST FIT: $y = 20x + 1.75$

3.2 LINEAR REGRESSION

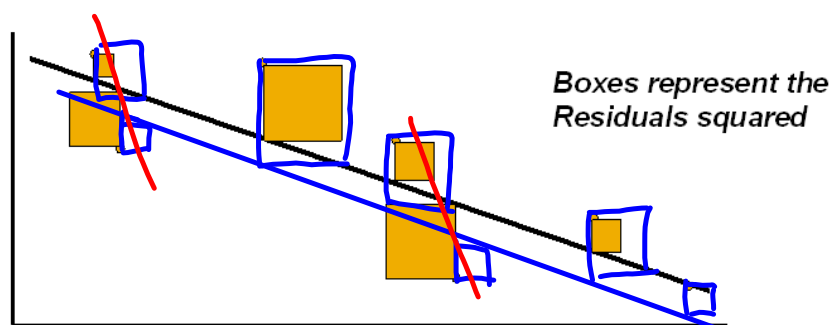
- An analytic method for determining the line of best fit can be determined by the *method of least squares*
- In order to understand the least squares method, we must first define the term *residual*
- In regression analysis, a *residual* is defined to be the vertical distance from a particular data point to the line of best fit



3.2 LINEAR REGRESSION

For the line of best fit in the *method of least squares*

- The sum of the residuals is zero (sum of the distance above the line is equal to the sum of the distance below the line)
- The sum of the squares of the residuals has the least possible value. (Boxes shown below are the smallest possible)



3.2 LINEAR REGRESSION

Statisticians have developed the following formula to determine the equation of the line of best fit using the least squares method

The equation of a line is given by: $y = ax + b$ $y = mx + b$

WHERE... $m a = \frac{n(\sum xy) - (\sum x)(\sum y)}{n(\sum x^2) - (\sum x)^2}$ (slope)

and

$$b = \bar{y} - a\bar{x} \quad \text{or} \quad b = \frac{\sum y - a\sum x}{n}$$

$$\bar{y} = \frac{\sum y}{n}$$

$$\bar{x} = \frac{\sum x}{n}$$

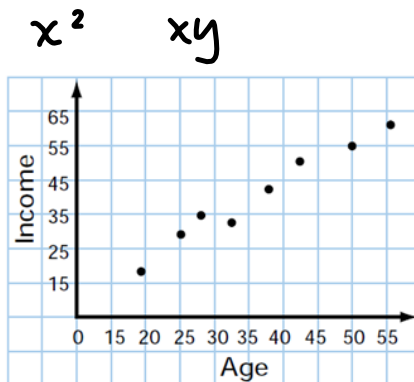
$$b = \frac{\sum y}{n} - m \frac{\sum x}{n}$$

$$= \frac{\sum y - m\sum x}{n}$$

3.2 LINEAR REGRESSION

The table and scatter plot show data for the full-time employees of a company:

| x | y |
|-------------|-----------------------|
| Age (years) | Annual Income (\$000) |
| 33 | 33 |
| 25 | 31 |
| 19 | 18 |
| 44 | 52 |
| 50 | 56 |
| 54 | 60 |
| 38 | 44 |
| 29 | 35 |



3.2 LINEAR REGRESSION

In order to calculate the line of best fit using the least squares method, the following table and calculations are set up.

| Age, x | Income, y | x^2 | xy |
|----------------|----------------|----------------------|---------------------|
| 33 | 33 | 1089 | 1089 |
| 25 | 31 | 625 | 775 |
| 19 | 18 | 361 | 342 |
| 44 | 52 | 1936 | 2288 |
| 50 | 56 | 2500 | 2800 |
| 54 | 60 | 2916 | 3240 |
| 38 | 44 | 1444 | 1672 |
| 29 | 35 | 841 | 1015 |
| $\sum x = 292$ | $\sum y = 329$ | $\sum x^2 = 11\,712$ | $\sum xy = 13\,221$ |

$$(\sum x)^2$$

$$\bar{x} = \frac{\sum x}{n}$$

$$\bar{y} = \frac{\sum y}{n}$$

$$\begin{aligned}a &= \frac{n(\sum xy) - (\sum x)(\sum y)}{n(\sum x^2) - (\sum x)^2} \\&= \frac{8(13\,221) - (292)(329)}{8(11\,712) - (292)^2} \\&= \frac{9700}{8432} \\&\doteq 1.15\end{aligned}$$

$$\begin{aligned}\bar{x} &= \frac{\sum x}{n} & \bar{y} &= \frac{\sum y}{n} & b &= \bar{y} - a\bar{x} \\&= \frac{292}{8} & &= \frac{329}{8} & &= 41.125 - 1.15(36.5) \\&= 36.5 & &= 41.125 & &= -0.85\end{aligned}$$

$$\begin{aligned}y &= ax + b \\&= 1.15x - 0.85\end{aligned}$$

Therefore, the equation of the line of best fit is $y = 1.15x - 0.85$.

3.2 LINEAR REGRESSION

- The slope a indicates only how y varies with x on the line of best fit
- The slope a does NOT tell anything about the strength of the correlation between the two variables (the correlation coefficient r does)
- It is possible to have a weak correlation with a large slope or a strong correlation with a small slope

3.2 LINEAR REGRESSION

Outliers

•an observation that is numerically distant from the rest of the data

Why Do We Have Outliers?

- Measurement Error
- Miscoding
- Misinterpretation
- Entered incorrectly
- Relationship is non-linear
- etc.

3.2 LINEAR REGRESSION

How to deal with outliers

- Most common method is to ignore outliers or not even consider them when doing regression analysis
- May make sense to leave outliers in dataset if they make a meaningful cluster and/or there are many of them
- Apply a different regression model (non-linear)

Good rule of thumb:

When in doubt, present results with and without outliers.