

$$(8) \quad SD = \sqrt{\frac{\sum (x - \bar{x})^2}{n}}$$

$$\text{Cov}(x, y) = \frac{\sum (x - \bar{x})(y - \bar{y})}{n}$$

$$R = \frac{\text{Cov}(x, y)}{\sigma_x \sigma_y} \quad (\text{Correlation Coefficient})$$

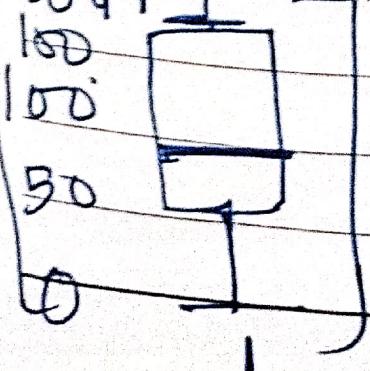
$$Z\text{-Score} = \frac{x - \bar{x}}{\sigma}$$

$$\text{Median} = l + \left(\frac{\frac{N}{2} - cf}{f} \right) \times h$$

Outliers detection (IQR)

$$\text{Lower limit} = Q1 - 1.5 \times IQR$$

$$\text{Upper limit} = Q3 + 1.5 \times IQR$$



Box Plot



Mo Tu We Th Fr Sa Su

Date: / /

Min - Max Normalized Stock

Price

$X = \text{Stock Price}$

$$X = \frac{X - \text{Min}}{\text{Max} - \text{Min}}$$

Mid - Range Stock Price

$$= \frac{\text{Min} + \text{Max}}{2}$$

~~Z-score Z-Score~~ (23)

$$= \frac{X - \bar{X}}{\sigma}$$

Decimal Scaling

$$= X' = \frac{X}{10^j}$$

Mean > median (Right Skewed +)

Median > Mean (-)

$$CI = \bar{x} \pm t_{\alpha/2, df} \times \left(\frac{s}{\sqrt{n}} \right)$$

Sample Proportion

$$= \frac{\text{# Success}}{\text{Size}}$$

$$\hat{q} = 1 - \hat{p}$$

$$CI = \hat{p} \pm z \times \sqrt{\frac{\hat{p}\hat{q}}{n}}$$

Population Proportion

Test Statistic (z-Score)

$$z = \frac{\hat{p} - p_0}{\sqrt{\frac{p_0(1-p_0)}{n}}}$$

P-Value?

Test Statistic

$$T = \frac{\bar{x} - \mu_0}{\frac{s}{\sqrt{n}}}$$

Null hypothesis (H_0):

$$\begin{cases} H_0 : \mu \leq 50 \\ H_1 : \mu > 50 \end{cases}$$

The Test Statistic T measures how many Standard Errors the Sample mean is away from the hypothesized population mean (\$50)

$$T = \frac{\bar{x}_1 - \bar{x}_2}{\sqrt{\frac{s_1^2}{n_1} + \frac{s_2^2}{n_2}}}$$

Chi-Square:

$$\chi^2 = \sum \frac{(O - E)^2}{E}$$

P-Value $\leq \alpha$ reject H_0
 P-value $> \alpha$ fail to reject H_0

$$Eij = \frac{(\text{Row Total})(\text{Column Total})}{\text{Grand Total}}$$

(X)

→ One independent Variable
Simple Linear Regression and one dependent Variable

$$Y = \beta_0 + \beta_1 X + \epsilon$$

Multiple Linear Regression

$$Y = \beta_0 + \beta_1 X_1 + \beta_2 X_2 + \dots + \beta_n X_n + \epsilon$$

2 or more independent Variables and one dependent Variable.

Least Squares Estimation Method

$$Y = a + bX$$

$$b = \frac{n \sum XY - (\sum X)(\sum Y)}{n \sum X^2 - (\sum X)^2}$$

$$a = \frac{\sum Y - b \sum X}{n}$$

SSE

Mo Tu We Th Fr Sa Su

Date: / /

$$SSE = \sum (y - \hat{y})^2$$

$$SSR = \sum (\hat{y} - \bar{y})^2$$

$$SST = \sum (y - \bar{y})^2$$

\hat{y} = Put the
n value
and find
this

Coefficient of Determination

$$R^2 = \frac{SSR}{SST} = \frac{\cancel{SSR}}{\cancel{SST}} = 1 - \frac{SSE}{SST}$$

Mean Square Error

$$= MSE = \frac{SSE}{n-2}$$

$$\text{Standard Error} = \sqrt{MSE} = \sqrt{\frac{SSE}{n-2}}$$

Pearson Correlation formula:

$$r_c = \frac{n \sum xy - (\sum x)(\sum y)}{\sqrt{n \sum x^2 - (\sum x)^2}}$$

$$\sqrt{n \sum y^2 - (\sum y)^2}$$

\rightarrow Correlation

\rightarrow Coefficient $R = \sqrt{R^2}$ Date: / /

$x_{\text{max}} - x_{\text{min}}$

Normalizer

\rightarrow Coefficient
of determination

$$x' = x - x_{\text{min}}$$

$$\frac{x_{\text{max}} - x_{\text{min}}}{x_{\text{max}} - x_{\text{min}}}$$

Euclidean Distance

$$\rightarrow d = \sqrt{(x_1 - x_2)^2 + (y_1 - y_2)^2}$$

Manhattan Distance:

$$d(u, y) = \sum |x_i - y_i|$$

Minkowski Distance

$$= d(u, y) = \left\{ \sum |x_i - y_i|^k \right\}^{1/k}$$

$(k \in N)$

Discrete Metrics

$$d(u, y) = 1$$

if $x_i \neq y_i$ and $d(u, y) \geq 1$

Weighted average

$$= \frac{\sum_{i=1}^k w_i y_i}{\sum_{i=1}^k w_i}$$

Robust Scaling

$$= \frac{x - Q_2}{IQR}$$

Types of Errors in Estimation

- ① Residual errors
- ② Measurement errors
- ③ Model errors
- ④ Random errors