

Python Data Analytics

→ Analysis means evaluating the past data to figure out why/ how it happened whereas analytic is exploring possible future by using past.

Descriptive analysis → What happened → eg. company yearly report

Diagnostic analysis → Why it happened

Predictive analytics → What will happen

Prescriptive analysis → What to do for better results.

• Level of Data measurement

1. Nominal Data → classifies into distinct categories, no ranking implied
eg. gender, marital status

2. Ordinal scale → Distinct categories in which ranking is implied.
eg. student grades (A, B, C, D...).

3. Interval scale → ordered scale in which difference b/w measurements have meaningful value but

measurements don't have true zero pt.

eg. temperature, potential etc.

4. Ratio scale \rightarrow difference b/w measurement is meaningful quantity and they have true zero point

eg. weight, age, salary,

usage potential max \rightarrow ~~nominal~~ Ratio
min \rightarrow nominal

Data ~~data~~ Visualisation

1. \rightarrow Ogive Curve

\rightarrow joining the mid-points of cumulative freq ~~histogram~~ bar graph.

2. Frequency polygon

\rightarrow join mid-points of histogram to observe data trend

3. Relative ogive

\rightarrow normalised ogive curve

4. Pareto Chart

\rightarrow Used to identify importance of a variable

5. Scatter plot

\rightarrow To analyse trend of a variable wrt another variable

Central tendencies → Arithmetic mean
 → Weighted mean
 → Median
 → Percentile

Dispersion → skewness
 → kurtosis
 → Range
 → Interquartile range
 → Variance
 → Standard score
 → Coefficient of ~~variation~~ ^{variation}

1. Central Tendencies

→ mean of grouped data = $\frac{\sum f_i M_i}{\sum f_i} = \mu$

M_i = mid-point of class interval

→ eg. interval = 20-30
 $M_i = 25$

→ Median of grouped data
 = $\frac{N/2 - cfp(w)}{f_{med}} + L$

L → Lower limit of median class

cfp → cumulative freq of prev. class

f_{med} → freq. of median class.

w → width of med class

N → total of frequencies

median class is one containing $N/2$
 cumulative freq.

$$\rightarrow \text{mode of grouped} = L_{mo} + \left(\frac{d_1}{d_1 + d_2} \right) w$$

modal class has highest freq.

$d_1 \rightarrow$ diff. b/w modal class freq. and prev class freq.

$d_2 \rightarrow$ diff. b/w modal and next class freq.

$w \rightarrow$ class width.

Note

\rightarrow If data Left / Right skewed use median but if symmetric use whatever u want.

$$\rightarrow i = \frac{P(n)}{100} \quad \text{pth percentile location.}$$

2. Dispersion.

• Quartiles eg (106, 109, 114, 116, 121, 122, 125, 129)

$$i_1 = \frac{25(n)}{100} \quad \rightarrow \text{scribbled out}$$

$$Q_1 = \frac{109 + 114}{2} = 111.5$$

$$i_2 = \frac{50(n)}{100} = \frac{50(8)}{100} = 4; \quad Q_2 = \frac{116 + 121}{2} = 118.5$$

$$i_3 = \frac{75(n)}{100} = 6; \quad Q_3 = \frac{122 + 125}{2} = 123.5$$

• interquartile range

$$= Q_3 - Q_1$$

• Population variance

$$\sigma^2 = \frac{\sum (X_i - \mu)^2}{N}$$

$X \rightarrow$ ~~sample~~ ^{ith sample}

$\mu \rightarrow$ mean

$N \rightarrow$ no. of samples

$$\rightarrow \sqrt{\sigma^2} = \sigma = \text{std. deviation}$$

Spiral

◦ Sample variance

$$s^2 = \frac{\sum (x_i - \mu)^2}{n-1}$$

S = sample standard deviation.

~~Variance~~ Note

→ Variance \propto Risk

→ 68% value lie within 1 std. deviation of mean in bell shaped curve.

◦ Chebysheff's theorem

→ ~~of all data~~ : (For $k > 1$) fraction
of all data will lie

→ The fraction of data lying within the k std deviations of the mean is at least

$$1 - \frac{1}{k^2} \quad \text{for } k > 1$$

→ This is the lower bound

◦ Coeff of variance

→ relative dispersion

$$C.V. = \frac{\sigma}{\mu} \times 100 \%$$

→ lower CV \Rightarrow better option

- Variance std for grouped data

1. Population

$$\sigma^2 = \frac{\sum f(M - \mu)^2}{N}$$

$$\sigma = \sqrt{\sigma^2}$$

2. Sample

$$s^2 = \frac{\sum f(M - \bar{x})^2}{n-1}$$

$$s = \sqrt{s^2}$$

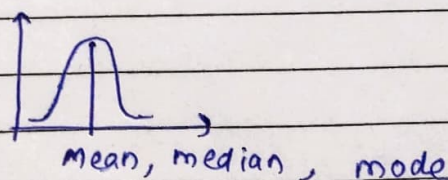
Note ~~in~~ ~~sa~~

→ In sample we divide by $n-1$ to provide more unbiased estimate of population variance (unbiased to sample)

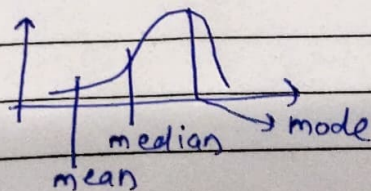
→ Above is called Bessel correction

- Skewness

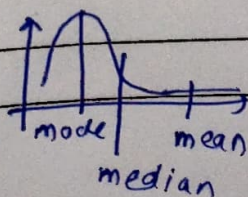
1. Symmetrical



2. Left skewed (-ve skewed)



3. Right skewed



Date _____

• coefficient^t of skewness

$$S = \frac{3(\mu - Md)}{\sigma}$$

$S < 0 \rightarrow$ -ve skew ; $S = 0 \rightarrow$ symmetric

$S > 0$ +ve skewed

• Kurtosis

Leptokurtic \rightarrow high thin

Meso kurtic \rightarrow normal in shape

Platykurtic \rightarrow flat and spread out



Note

Median = Q_2

