

## Ch : DESCRIPTIVE STATISTICS

**GOAL:** Graphical and numerical representation of data

**TOPICS :**

- 1) Basic Concepts of Data Acquisition
- 2) Analysis of Univariate Data

- 3) Analysis of Bivariate Data

### \* BASIC CONCEPTS OF DATA ACQUISITION

- **Statistical unit or entity:** A single "object" is called a stat. unit.  
e.g.: Robotics student XYZ in course statistics in WS 2023.
- **Population:** set of all objects under consideration.  
e.g.: entire class of robotics students in statistics in WS 2023.
- **Variables:** properties of a statistical unit one is interested in.  
a.k.a attributes / characteristics
- **Occurrence:** specific value a variable takes on.

#### Example :

Statistical unit	Matr. no.	Citizenship	Civil Status	Age	variable
					Population
	3313019	UK	single	21	
	3312008	India	single	18	Occurrences

→ **Key Attributes:** A minimal combination of attributes that identifies a statistical unit uniquely.

→ **Possible Occurrences:**

→ **Scale of Variables:** • **Categorical / Nominal scale:** only labels  
language

• **Ordinal scale:** sorting possible → Citizenship

→ Army rank

→ Placement in a sporting event.

• **Cardinal / metric scale:** Calculation possible

→ Stock & commodity prices

→ Age, weight, length

*NOTE: Categorical scale can be multi,  
but not the other 2 sc*

→ Multivalued Variables: Sometimes a variable can take 2/more values at the same time.

e.g.: A student  
 ↪ more than 1 citizenship  
 ↪ speaks more than 1 language

### Q

### Exercise

- | Nominally Scaled               | Ordinaly Scaled  | Cardinally Scaled        |
|--------------------------------|--|--------------------------|
| • Gender → m/w/d               | • Job → team member<br>fire fighter → junior chief<br>professor → senior chief | • Income → DE, ISOK, 300 |
| • Job → programmer<br>Tia, Ria | • Academic degree<br>Knight, Duke, Count                                       | • Age → 2, 15, 32 (yo)   |
| • First name → Rob, Tia, Ria   | • Knight, Duke, Count → title  | • Temperature (in K)     |
|                                |  | • Grades                 |



### ANALYSIS OF UNIVARIATE DATA

- raw data - all occurrences written down as a list with no specific order.  
multiset / bag
- sorted lists - raw data sorted w.r.t specific criterion  
(e.g.: in ascending order)  
at least ordinal scale data required.  
every occurrence is written down only once, plus we write frequencies.

- Univariate frequency:  $h(a_j)$  of the occurrence  $a_j$  appears in raw data.
- relative frequency:  $f(a_j) = h(a_j) / n$  of occurrence  $a_j$  tells us fraction with which  $a_j$  occurs in raw data.

raw data      sorted list  
 ex: 2, 3, 1, 1, 5, 3, 6, 3, 3, 1 → 1, 1, 1, 2, 3, 3, 3, 3, 5, 6

$a_j$	1	2	3	4	5	6	$\Sigma$
$h(a_j)$	3	1	4	0	1	1	10
$f(a_j)$	0.3	0.1	0.4	0	0.1	0.1	1

$$h_j = h(a_j) = \text{absolute frequencies}$$

$$f_j = f(a_j) = \frac{h(a_j)}{n}$$

- Cumulative frequencies
- n: number of univariate statistical units ( $n_1, \dots, n_n$ ) and k: number of different occurrences ( $a_1, \dots, a_k$ )

- Absolute cumulative frequency =  $H_n = (a_1, a_2, \dots, a_n)$
- $H_j = h_1 + h_2 + \dots + h_j$
- relative cumulative frequency

$$F_j = f_1 + f_2 + \dots + f_j = \frac{H_j}{n}$$

Generally  $H_k = n$

$$F_k = 1$$

- Cumulative frequency Distribution

$$H(n) = \begin{cases} 0 & \text{if } n < a_1 \\ h_j & \text{if } a_j \leq n < a_{j+1} \\ n & \text{if } n \geq a_k \end{cases}$$

$$F(n) = H(n) = \begin{cases} 0 & \text{if } n < a_1 \\ F_j & \text{if } a_j \leq n < a_{j+1} \\ 1 & \text{if } n \geq a_k \end{cases}$$

g) Exercise : total no. of objects produced in one hour:

No. of objects:	3	4	5	6	7	8	9	10	
Abs. frequency:	10	15	30	30	25	20	15	5	

$$F(3), F(5.5), F(10) = ?$$

How many hours did data acquisition take at least?

Statistical unit: production hours (1 hour)

Variable: number of objects

$$n = \text{sum of all abs. frequencies} = 150$$

$$a_1 = 3, a_2 = 4, \dots, a_8 = 10$$

$$F(n) = \begin{cases} 0 & \text{if } n < a_1 \\ \frac{H_j}{n} & \text{if } a_j \leq n < a_{j+1} \\ 1 & \text{if } n \geq a_8 \end{cases}$$

$$x = a_1$$

$$\rightarrow F(3) = \frac{H_1}{150} = \frac{10}{150} = \frac{1}{15} = 0.0667 \approx 0.067$$

$$\rightarrow F(5.5) = \frac{H_3}{150} = \frac{55}{150} = \frac{11}{30} = 0.367$$

$$\rightarrow F(10) = 1$$

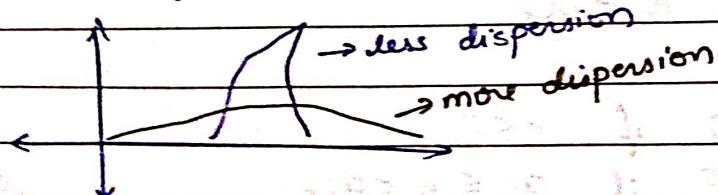
- Numerical Representation of Univariate Data.

univariate data can be numerically represented by

→ measures of central tendency

→ measures of dispersion

→ measures of concentration



- **Measures of central tendency**

- ↳ describes where data is primarily located
- ↳ whole data sample summed to a single value
- imp. measures of C.T.
- mode
- median
- mean

- **MODE**

- $\bar{n}_{\text{mod}}$  is the most frequent value.
- Adv.: mode can already be used with nominal data
- Disadv.: only makes sense if value is unique.

e.g. data sample:  $(2, 3, 1, 1, 5, 3, 6, 3, 3, 1)$

$a_j$	1	2	3	4	5	6	$\bar{n}_{\text{mod}} = 3$
$h_j$	3	1	4	0	1	1	

- **MEDIAN**

- $\bar{n}_z$  separates lower half from upper half.
- median appears as a value in the data list. (unlike mean)
- if data sample is sorted in asc. order, median  $\bar{n}_z$  is defined by

$$\bar{n}_z = \begin{cases} n_{\frac{z+1}{2}} & \text{if } n \text{ odd} \\ \frac{n_z + n_{z+1}}{2} & \text{if } n \text{ even} \end{cases}$$

Alternate:  $\bar{n}_z = \begin{cases} n_{\frac{z+1}{2}} & \text{if } n \text{ odd} \\ \frac{n_z + n_{z+1}}{2} & \text{if } n \text{ even} \end{cases}$

- Adv.: median can be used with ordinal data
- Dis.: not unique for even  $n$ .

- p Quantile

- in case of ordinal data, concept of median can be generalized to p-quantile
- at most  $p \cdot n$  of the values in data sample are smaller and at most  $(1-p) \cdot n$  of values are larger than value of the p quantile.
- $p = 0.25$  first quartile  $Q_1$   
 $p = 0.5$  second quartile  $Q_2$  = median  
 $p = 0.75$  third quartile  $Q_3$

g) **Exercise:** Determine the quartiles  $Q_1, Q_2, Q_3$  of the following data samples:

a)  $(2, 7, 8, 11, 13, 17)$

$\downarrow \quad \uparrow \quad \uparrow \quad \downarrow$   $Q_1, Q_2, Q_3$  means met at most 25% values smaller than  $Q_1$ , at most 75% values larger than  $Q_1$ ,  $\therefore 2 \neq Q_1$ , since  $2 < 7$

process

$Q_1 = 7$

similar logic for  $Q_3$

b)  $(1, 5, 66, 234, 440, 489, 500)$

$\downarrow \quad \downarrow \quad \downarrow$   $Q_1, Q_2, Q_3$

c)  $(1, 3, 5, 66, 111, 234, 440, 489, 500, 777)$

$\downarrow \quad \uparrow \quad \uparrow \quad \downarrow$   $Q_1, Q_2, Q_3$

## • MEAN

↳ computation of mean only possible for cardinal data.

→ Arithmetic mean  $\bar{x} = \frac{1}{n} \sum_{i=1}^n x_i$

→ Geometric mean  $\bar{x}_G = (x_1 \cdot x_2 \cdot \dots \cdot x_n)^{\frac{1}{n}}$  if product ≠ -ve

→ Harmonic mean  $\bar{x}_H = \frac{n}{\sum_{i=1}^n \frac{1}{x_i}}$

→ ARITHMETIC MEAN:  $\bar{x}_A = \frac{1}{n} \sum_{i=1}^n x_i$

→ for given data, arithmetic mean solves minimization problem.

→ in contrast to mode / median, arith. mean does not have to occur in data (sample).

## → GEOMETRIC MEAN:

$$\text{Area of rectangle } ab = \text{Area of square } s^2 \Rightarrow ab = s^2 \Rightarrow s = \sqrt{ab}$$

(can also do this for cubes etc)

↳ geometric mean

$$\rightarrow \bar{x}_G = \sqrt[n]{x_1 \cdot x_2 \cdot \dots \cdot x_n}$$

→ only defined for +ve values.

$$\bar{x} \geq \bar{x}_G \Rightarrow \text{AM vs. GM}$$

## → HARMONIC MEAN:

$$\bar{x}_H = \frac{n}{\sum_{i=1}^n \frac{1}{x_i}}$$

nominal data → mode

ordinal data → median

cardinal data → mean

$$\min \leq x_H \leq x_G \leq \bar{x} \leq \max$$

- **Measures of Dispersion**

- describe how data is dispersed around a central tendency
- most important dispersion parameters are:
- the range
- mean deviation
- empirical variance and standard deviation

- **Range**

- a.k.a width = largest value - smallest value  
 $a_x - a_1$

- **Mean Deviation**

- a.k.a mean absolute deviation (mad)

$$\text{m.a.d} = \frac{1}{n} \sum_{i=1}^n |x_i - m| \quad : m = \text{arithmetic mean or med./mode}$$

$$\text{m.a.d} = \frac{1}{n} \sum_{i=1}^n |x_i - m| = \frac{1}{n} \sum_{j=1}^k h_j \cdot |a_j - m| = \sum_{j=1}^k f_j \cdot |a_j - m|$$

- **Empirical variance and standard deviation**

$$\rightarrow s^2 = \frac{1}{n} \sum_{i=1}^n (x_i - \bar{x})^2 = \frac{1}{n} \left( \sum_{i=1}^n x_i^2 - \bar{x}^2 \right)$$

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

$$\rightarrow s^2 = \frac{1}{n} \sum_{i=1}^n (m_i - \bar{x})^2 = \frac{1}{n} \sum_{j=1}^k h_j \cdot (a_j - \bar{x})^2 = \sum_{j=1}^k f_j \cdot (a_j - \bar{x})^2$$

- Relative dispersion - generally defined as a ratio  
  - absolute dispersion
  - mean

$$\rightarrow \text{coefficient of variation} \quad V = \frac{s}{\bar{x}} \rightarrow \text{std. dev.} \\ \frac{s}{\bar{x}} \rightarrow \text{mean}$$

Q

Exercise	$a_j$	1	2	3	4	5	6	a) mode, median, mean, GM, HM, AM
$\sum = 10$	$b_j$	1	2	3	2	1	1	HM : $\frac{1}{1} + \frac{1}{2} + \frac{1}{3} + \dots + \frac{1}{6}$
$\sum = 1$	$f_j$	0.1	0.2	0.3	0.2	0.1	0.1	b) range, m.a.d., variance, std dev, coeff of var.

a) mode = 3 median =  $\frac{3+3}{2} = 3$  (sorted list: 1, 2, 2, 3, 3, 4, 4, 5, 6)

arithmetic mean =  $\frac{\sum n_i}{n} = \frac{33}{10} = 3.3$  HM:  $\frac{1}{1} + \frac{1}{2} + \frac{1}{3} + \dots + \frac{1}{6}$

Geometric mean =  $(1 \cdot 2 \cdot 2 \cdot 3 \cdot 3 \cdot 4 \cdot 4 \cdot 5 \cdot 6)^{\frac{1}{10}} = 2.96 = 2.59$

b) Range = 6 - 1 = 5 mean abs. dev. =  $\frac{\sum |n_i - m|}{n} = \frac{|3.3 - 1| + |3.3 - 2| + \dots + |3.3 - 6|}{10}$

$s^2 = 2.01$   $s = 1.42$

$= 2.3 + 1.3 + 1.3 + (0.3)3 + (0.7)2 + \dots$

$V = \frac{s^2}{\bar{x}^2} = 0.43$

$3.33$

10

$= \frac{11.6}{10} = 1.16$

### Merging Aggregates

→ for each group  $G_e$ , we know the size of the group  $n_e$ , the arithmetic mean  $\bar{x}_e$  and empirical variance  $s_e^2$ .

→ overall result from merging the groups:

$$n = \sum_{e=1}^m (n_e)$$

$$s^2 = \sum_{e=1}^m n_e \cdot s_e^2 + \sum_{e=1}^m n_e \cdot (\bar{x}_e - \bar{x})^2$$

$$\bar{x} = \frac{\sum_{e=1}^m n_e \cdot \bar{x}_e}{n}$$

Q Exercise - Consider a statistical survey where the results of  $G_1, G_2, G_3$  are given:  $G_1: n_1 = 10, \bar{x}_1 = 100, s_1^2 = 400$

$$G_2: n_2 = 5, \bar{x}_2 = 120, s_2^2 = 144$$

$$G_3: n_3 = 20, \bar{x}_3 = 60, s_3^2 = 100$$

Calculate overall  $\bar{x}$  and  $s^2$

$$n = 10 + 5 + 20 = 35$$

$$\bar{x} = \frac{10 \times 100}{35} + \frac{5 \times 120}{35} + \frac{20 \times 60}{35} = 80$$

$$s^2 = \frac{10 \times 400}{35} + \frac{5 \times 144}{35} + \frac{20 \times 100}{35} + \frac{10 \times 20^2}{35} + \frac{5 \times 40^2}{35} + \frac{20 \times 20^2}{35}$$

$$= 763.43$$

$$s = 27.63$$

\*

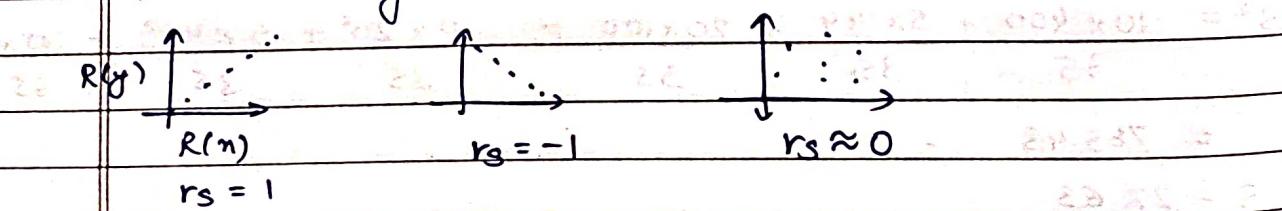
## ANALYSIS OF BIVARIATE DATA

- Bivariate or 2-D data refers to a collection of pairs  $((n_1, y_1), \dots, (n_n, y_n))$  where the values of  $n$  and  $y$  are recorded simultaneously.
- Imp concepts: correlation and regression
- Scatter plot: data is displayed as a collection of points. it is possible to visually identify correlations b/w the variables.
- Correlation:
- correlation describes a relationship b/w two variables which are of at least ordinal scale.  
eg: taller people are usually heavier.
- correlation is just a statistical observation. NOT a cause-symptom relation
- Spearman's Rank Correlation Coefficient
- Consider bivariate data  $(n_1, y_1), \dots, (n_n, y_n)$  where  $n_i, y_i$  are occurrences of ordinal scale.
- every  $n_i$  receives a rank  $R(n_i)$  and every  $y_i$  receives a rank  $R(y_i)$  w.r.t their position.

### POSSIBLE CASES

1. The ranks  $n_i$  and  $y_i$  run in the same direction  
↳ POSITIVE CORRELATION
2. The ranks  $n_i$  and  $y_i$  run in opposite direction  
↳ NEGATIVE CORRELATION
3. There is no relation b/w ranks of  $n_i$  and  $y_i$   
↳ uncorrelated variables

Spearman's rank correlation coeff :  $r_s$  measures strength & direction of correlation.



$$r_s = 1 - \frac{6 \cdot \sum_{i=1}^n (R(n_i) - R(y_i))^2}{(n-1) \cdot n \cdot (n+1)}$$

$$-1 \leq r_s \leq 1.$$

$r_s = 1$ : ranks run in same direction

$r_s = -1$ : ranks run in opp. direction

$r_s = 0$ : ranks run without relation

- Q calculate Spearman's rank correlation coeff for the variables "Age" and "Points in test".

Student No.	1	2	3	4	5	6	7	8	9	10	11
Age (years)	38 <sup>4</sup>	47 <sup>2</sup>	44 <sup>5</sup>	51 <sup>1</sup>	35 <sup>5</sup>	29 <sup>6</sup>	22 <sup>7</sup>	14 <sup>9</sup>	12 <sup>10</sup>	19 <sup>8</sup>	9 <sup>11</sup>
Points in test	39 <sup>3</sup>	34 <sup>4</sup>	31 <sup>5</sup>	48 <sup>1</sup>	46 <sup>2</sup>	23 <sup>7</sup>	17 <sup>8</sup>	12 <sup>10</sup>	16 <sup>9</sup>	28 <sup>6</sup>	10 <sup>11</sup>

$$r_s = 1 - \frac{6 \cdot \sum (R(n_i) - R(y_i))^2}{(n-1)(n)(n+1)}$$

$$\sum (R(n_i) - R(y_i))^2 = 1^2 + 2^2 + 2^2 + 0 + 3^2 + 1^2 + 1^2 + 1^2 + 1^2 + 2^2 + 0 \\ = 1 + 4 + 4 + 9 + 1 + 1 + 1 + 1 + 4 = 26$$

$$\therefore r_s = 1 - \frac{6 \times 26}{10 \times 11 \times 12} = 1 - 0.118 = 0.881$$

most imp

### • Coefficient of Correlation of Bravais Pearson:

takes into account distance of data points from their respective ranks  
and not only ranking

- for cardinal data
- it is a measure for the linear relationship b/w 2 cardinal variables.
- It measures how well the bivariate data can be approximated by a so called linear trend or regression line.

$$r = \frac{\sum_{i=1}^n (n_i - \bar{n}) \cdot (y_i - \bar{y})}{\sqrt{\sum_{i=1}^n (n_i - \bar{n})^2} \cdot \sqrt{\sum_{i=1}^n (y_i - \bar{y})^2}} \quad -1 \leq r \leq 1$$

$r=1$ : data exactly located on a line with the slope 1

$r=-1$ : data exactly located on a line with -ve slope

$r=0$ : data shows no linear relationship at all.

Q

Determine the coeff. of correlation for the variables "Inflation" and "Jobless rate":

Year	2001	2002	2003	2004	2005
Inflation (%)	2	3	3	2	5
Jobless rate (%)	4	7	2	3	4
$n_i - \bar{n}$	-1	0	0	-1	2
$y_i - \bar{y}$	0	3	-2	-1	0

$$\bar{n} = \frac{2+3+3+2+5}{5} = 3$$

$$\bar{y} = \frac{4+7+2+3+4}{5} = 4$$

$$r = \frac{(-1)(0) + (0)(3) + (0)(-2) + (-1)(-1) + (2)(0)}{\sqrt{1+0+0+1+4} \cdot \sqrt{0+9+4+1+0}}$$

$$= \frac{-1}{\sqrt{6} \cdot \sqrt{14}} = \frac{1}{9.165} \approx 0.109 \rightarrow \text{no linear correlation}$$

- **Empirical covariance (cov)**

→ is an auxiliary measure for the joint variability of two cardinally scaled variables.

(for bivariate data)

$$C.O.V(n, y) = \frac{1}{n} \sum_{i=1}^n (n_i - \bar{n})(y_i - \bar{y})$$

$$\rightarrow \text{coeff. of covariance } r = \frac{\text{cov}(n, y)}{S_x \cdot S_y} \quad \begin{matrix} \text{standard deviations} \\ \downarrow \end{matrix}$$

relationship b/w coeff. of Spearman & Brava's Pearson

$$r_s = \frac{\sum_{i=1}^n (R(n_i) - \bar{R}(n)) \cdot (R(y_i) - \bar{R}(y))}{\sqrt{\sum_{i=1}^n (R(n_i) - \bar{R}(n))^2} \cdot \sqrt{\sum_{i=1}^n (R(y_i) - \bar{R}(y))^2}}$$

$$\text{if ranks are unique } \bar{R}(n) = \bar{R}(y) = \frac{n+1}{2}$$

$m_i$	$R(m_i)$
1000	1
800	2
100	3
50	4
10	5

$\left. \begin{array}{l} \\ \\ \\ \\ \end{array} \right\} R(n) = 3$

- **Regression Analysis**

→ in contrast to correlation analysis, regression analysis assigns two different roles for  $n$  and  $y$ .

→  $n$ : cause       $y$ : symptom

$y = f(n)$  functional relationship  
↓ independent  
↓ dependent symptom

- Computing a linear trend line:

→ searching for linear function of the type:

$$y = f(n) = an + b$$

with minimal distances to the data  $(n_1, y_1), \dots, (n_n, y_n)$

$a, b$ : coefficients of regression

→ to find  $a, b$ : OPTIMIZATION PROBLEM:

$$Q(a, b) = \sum_{i=1}^n (y_i - f(n_i))^2 = \sum_{i=1}^n (y_i - an_i - b)^2 \rightarrow \text{minimize}$$

residual

$$= \sum_{i=1}^n (y_i - an_i - b)^2 \rightarrow \text{minimize}$$

$$\frac{\partial Q}{\partial a} = \sum_{i=1}^n -2n_i(y_i - an_i - b) = 0$$

$$\frac{\partial Q}{\partial b} = \sum_{i=1}^n -2(y_i - an_i - b) = 0$$

Computing, we get  $a = \frac{\sum_{i=1}^n m_i \cdot y_i - \frac{1}{n} \left( \sum_{i=1}^n m_i \right) \cdot \left( \sum_{i=1}^n y_i \right)}{\sum_{i=1}^n m_i^2 - \frac{1}{n} \left( \sum_{i=1}^n m_i \right)^2}$

(slope)

(y intercept)  $b = \frac{1}{n} \left( \sum_{i=1}^n y_i - a \sum_{i=1}^n m_i \right)$

→ Alternative expressions

$$a = \frac{\sum_{i=1}^n (x_i - \bar{x}) \cdot (y_i - \bar{y})}{\sum_{i=1}^n (x_i - \bar{x})^2} = \frac{\text{cov}(x, y)}{s_x^2} = r \frac{s_y}{s_x}$$

$$b = \bar{y} - a\bar{x}$$

- Q A distribution center runs ten shrink wrap machines on different velocity  $x$  (m/min) the variable  $y$  indicates the number of stops caused by cracks. Compute  $a$  &  $b$

Machine	1	2	3	4	5	6	7	8	9	10	$\bar{x} = 26$
$x_i$	21	22	23	24	25	26	27	28	30	34	$\bar{y} = 50$
$y_i$	30	30	30	40	50	50	60	60	70	80	
$x_i - \bar{x}$	-5	-4	-3	-2	-1	0	1	2	4	8	
$y_i - \bar{y}$	-20	-20	-20	-10	0	0	10	20	30	40	
$(x_i - \bar{x})(y_i - \bar{y})$	100	80	60	20	0	0	10	20	80	240	
$(x_i - \bar{x})^2$	25	16	9	4	1	0	1	4	16	64	

$$a = \frac{610}{140} = 4.36$$

$$b = \bar{y} - a\bar{x} = 50 - (4.36 \times 26) = -63.36$$

- Nonlinear Regression

- linear regression is not always suitable for relation b/w  $x$  and  $y$ .
- in case of monotonically increasing / decreasing regression functions, suitable substitution can reduce the problem to linear regression.

e.g:

$x_i: 1 \ 2 \ 3 \ 4$

$y_i: 8 \ 18 \ 30 \ 51$

→ To compute  $a$ ,  $b$

$\therefore \tilde{x} = 1 + \frac{1+2+3+4}{4} = 2.5$

$$a = \frac{\sum (\tilde{x}_i - \bar{\tilde{x}}) \cdot (y_i - \bar{y})}{\sum (\tilde{x}_i - \bar{\tilde{x}})^2} = 2.82$$

$$\sum (\tilde{x}_i - \bar{\tilde{x}})^2$$

$$b = \bar{y} - a\bar{\tilde{x}} = 5.6$$

$$\therefore y = f(x) = 2.82x^2 + 5.6$$

- Quality of Regression Analysis: we compare empirical values  $y_i$  against the corresponding values  $f(n_i)$  of the regression function

→ estimates  $\hat{y}_i = f(n_i)$

residual  $\hat{u}_i = y_i - \hat{y}_i$

→ coefficient of determination  $R^2 = 1 - \sum_{i=1}^n \hat{u}_i^2$

$0 \leq R^2 \leq 1$

if  $R^2 = 1$ : Perfect regression function, all residuals are 0.  
 $R^2 < 0.1$ : Obtained regression function is nonsense

→ Alternate Formula: to find  $R^2$

In case of a linear regression function  $y = f(n) = an + b$

$$R^2 = \frac{S_y^2}{S_{\hat{y}}^2} = \frac{\sum_{i=1}^n (y_i - \bar{y})^2}{\sum_{i=1}^n (\hat{y}_i - \bar{y})^2} = a^2 S_x^2 = r^2$$

↳ coeff. of correlation of Bravais Pearson

- Time Series Analysis

→ time series is a time ordered collection of sequentially observed data pairs.

$$(t_i, n_i)_{i=1}^n = ((t_1, n_1), \dots, (t_n, n_n))$$

where  $t$  represents points in time  $t_1 < t_2 < \dots < t_n$

$n$ : cardinally scaled

→ time series shows temporal patterns

- seasonal cycles
- long term trends
- random fluctuations

→ time dependent variable  $n = T + C + S + R$

T: trend: long term tendency

C: cycle: medium term periodic component

S: season: short term periodic component

R: random component: irregular fluctuations

- Trend:  $T$  - long term tendency of a time series, often in terms of a straight line (regression analysis)
- Component Cycle:  $C$  - medium term oscillations, cycle  $c$  is a periodic component with cycle length or period larger than the length of  $S$ . Standardized  $C$  shows oscillations without trend  $T$
- Component Season:  $S$  - short term oscillations with typical pattern and well known wavelengths, typical season period is one year.
- Random Component:  $R$  - irregular & unpredictable effects  
 $R = n - (T+C+S)$ , Standardized  $R$  fluctuates around 0 line

⇒ COURSE OF ACTION: Time Series

- 1) Determine  $T$  by regression analysis
- 2) Determine smooth component  $G = T + C$  by filtering time series based on moving averages.
- 3)  $C = G - T$
- 4) Find  $S$  by averaging over difference  $n - G = S + R$
- 5)  $R = n - T - C - S$

Q Determine the moving averages  $n_i^*$  of the following time series according to a time window of length 5:

$t_i$	0	1	2	3	4	5	6	7	8	9	10	11	12	13	14
$n_i$	3.00	4.07	5.13	6.20	7.27	3.33	4.40	5.47	6.53	7.60	3.67	4.73	5.80	6.87	7.93
$n_i^*$	X	X	5.13	5.2	5.27	5.33	5.4	5.47	5.53	5.6	5.67	5.73	5.8	5.87	X

Case: In case of an even time window eg:  $6 = 2k \Rightarrow k = 3$

0	1	2	3	4	5	6	7	8
3	4.07	5.13	6.20	7.27	3.33	4.40	5.47	6
X	X	X	4.95					

$$\frac{1}{6} \left( \frac{3}{2} + \frac{4.4}{2} + 4.07 + 5.13 + 6.20 + 7.27 + 3.33 \right)$$

- Seasonal Adjustment by Moving averages

$$\rightarrow n^* \approx G = T + C$$

$$\therefore n - n^* \approx S + R$$

→ Summary :

- 1) Filter  $n$  in terms of  $n^*$
- 2) Remove smooth component  $n^* \approx G = T + C$  from time series to get  $S + R$
- 3) Average  $n - n^*$  to obtain seasonal averages
- 4) Standardize seasonal averages so that they oscillate around the zero line to obtain seasonal pattern  $S$ .
- 5) determine seasonally adjusted time series  $n - S$  by removing  $S$  from  $n$ .

$$n = T + C + S + R$$

$$G = T + C$$

Q Exercise : In a tide station, the water level is measured on four consecutive days at 0:00 am, 8:00am and 4:00 pm respectively. The results of the measurements in cm are as follows -

	Day 1	Day 2	Day 3	Day 4
	603   723   480	606   720   420	600   660   420	537   660   600

Apply appropriate moving averages to ~~get~~ obtain the seasonal pattern and the seasonally adjusted time series.

$t_i$	0:00	8:00	4:00	0:00	8:00	4:00	0:00	8:00	4:00	0:00	8:00	4:00
$n_i$	1	2	3	4	5	6	7	8	9	10	11	12
$G = n_i^*$	603	723	480	606	720	420	600	660	420	537	660	600
$n_i - n_i^*$	X	602	603	602	582	580	560	560	539	538	599	X
$\bar{n}_i - \bar{n}_i^*$	X	(121)	-123	4	(138)	-160	40	(100)	-119	-2	(61)	X = S + R
$\bar{n}_i - S$	584	613	609	587	610	549	581	550	549	518	550	729

Seasonally adjusted time series	0:00 am	8:00 am	4:00 pm	
	0:00 am	14	$\frac{4+40-2}{3}$	+ 5
	8:00 am	$105 \rightarrow \frac{121+138+61}{3}$	+ 5	$110 = \begin{pmatrix} S_1 \\ S_2 \\ S_3 \end{pmatrix} = S$
	4:00 pm	-134	+ 5	-129

- Time Series Forecasting by Extrapolation
  - 1) Acquire time series data  $((t_i, n_i))_{i=1}^n$
  - 2) Check preconditions
    - one full cycle is contained in time series data
    - atleast a couple of seasons
    - seasonal wavelength is well known
  - 3) Determine  $T, C, S$
  - 4) Forecast the time series by forwarding  $T$  and repeating periodic components  $C$  and  $S$ .
- Alternative methods of filtering

- Filtering by smoothing
- a.k.a exponential weighted moving average
- smoothing factor  $\alpha$  where  $0 \leq \alpha \leq 1$   
larger values of  $\alpha$  reduce extent of smoothing

$$\rightarrow n_i^* = n_i, \quad n_{j+1}^* = \alpha \cdot n_j + (1-\alpha) \cdot n_j^*$$

→ smoothed once - one level smoothing

smoothed again - 2 level smoothing

smoothed  $n$  times -  $n$  level smoothing

$$\text{for forecast: } n_i^* = n_i, \quad n_{j+1}^* = \alpha n_j + (1-\alpha) n_j^*$$

8 Quarterly production volumes of a garment factory are available for three consecutive years 2021, 2022 and 2023.

	2021	2022	2023
	10   12   8   14	12   14   16   12   18	18   20   16   22

Apply exponential smoothing to the time series with  $\alpha = 0.2$  in order to give a forecast on the production volume of the first quarter in 2024.

$t_i$	1	2	3	4	5	6	7	8	9	10	11	12	13
$n_i$	10	12	8	14	14	16	12	18	18	20	16	22	?
$n_i^*$	10	10.4	9.92	10.73	11.38	12.3	12.24	10.38	12.21	15.55	14.31	16.84	

Using  $n_i^* = n_i$ ,

$$\text{and } \left. \begin{aligned} & n_{j+1}^* = \alpha n_j + (1-\alpha) n_j^* \\ & \end{aligned} \right\}$$

$$\text{and } n_{j+1}^* = \alpha n_j + (1-\alpha) n_j^*$$

} forecast