



Vidyalankar
Institute of
Technology

(Accredited A+ by NAAC)
(Autonomous Institute Affiliated to University of Mumbai)

Prof. Rancharan Dhusi

Mid Semester Examination

Branch	Date	Sem.	Roll No. / Exam Seat No.	Subject	Student's Signature	Junior Supervisor's Name and Sign
AI	1/4/24	6		SL for DS		

Question No.	A	B	C	D	E	F	G	H	Total	Total out of (20 / 30 / 40)
1										
2										
3										
4										

Examiners Signature	Student's Sign (After receiving the assessed answer sheet)
---------------------	---

Q1. What are the types of statistics? Give examples.
Also explain classification of data & what are the different scales of measurement for data?

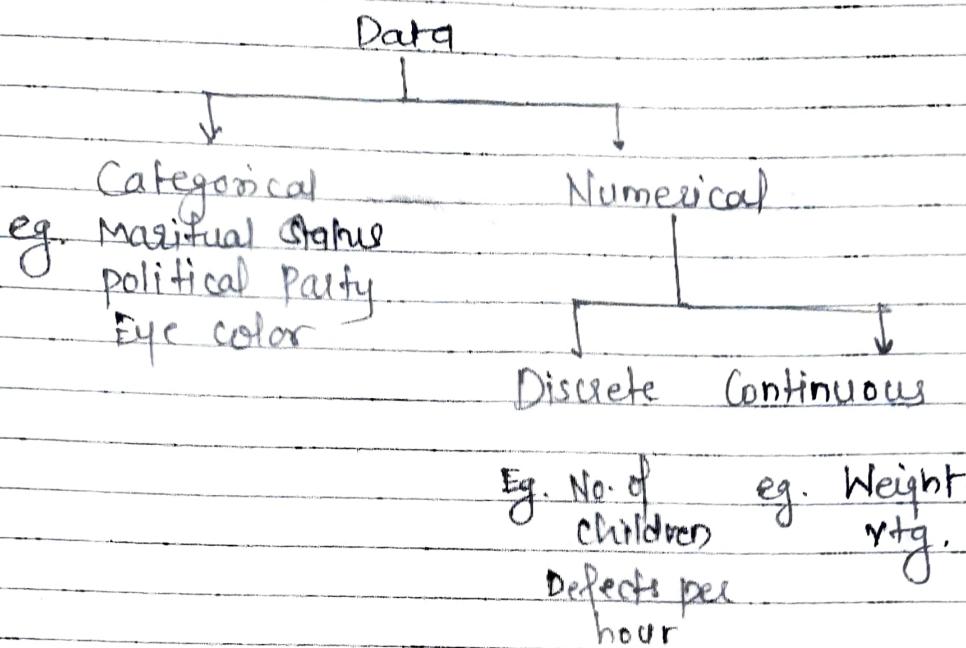
→ A branch of mathematics taking & transforming numbers into useful information for decision makers.

Types of statistics.

→ Descriptive statistics :- collecting, summarizing & describing data. called descriptive statistics.

→ Inferential Statistics:- Drawing conclusions and/or making decisions concerning a population based only on sample data.

Classification of Data



Scales of Measurement:-

- The scale determines the amount of information contained in the data!
- The scale indicates the data summarization & statistical analyses that are most appropriate.

Scales of measurement include:

Nominal	Interval
ordinal	Ratio .

→ Nominal :-

Data are labels or names used to identify an attribute of the element. A non-numeric label or numeric code may be used.

eg. Students of a university are classified by the school in which they are enrolled using a non-numeric label

such as Business, Humanities, Education

I. Business & Humanities & Education

2) Ordinal :-

The data have the properties of nominal data & the order or rank of the data is meaningful.

A non numeric label or numeric code may be used.

3) Interval :- The data have the properties of ordinal data & the interval between observations is expressed in terms of a fixed unit of measure

- Interval data are always numeric.
eg. Melissa has an SAT score of 1205 while Kevin has an SAT score of 1090. Melissa scored 115 points more than Kevin.

4) Ratio :-

The data have all the properties of interval data the ratio of two values is meaningful

- Variables such as distance, height, weight & time use the ratio scale.

- This scale must contain a zero value that indicates that nothing exists for the variable at the zero point.

eg. Melissa's college record shows 36 credit hours earned, while Kevin's record shows 32 credit hours earned. Kevin has twice as many credit hours earned as Melissa.

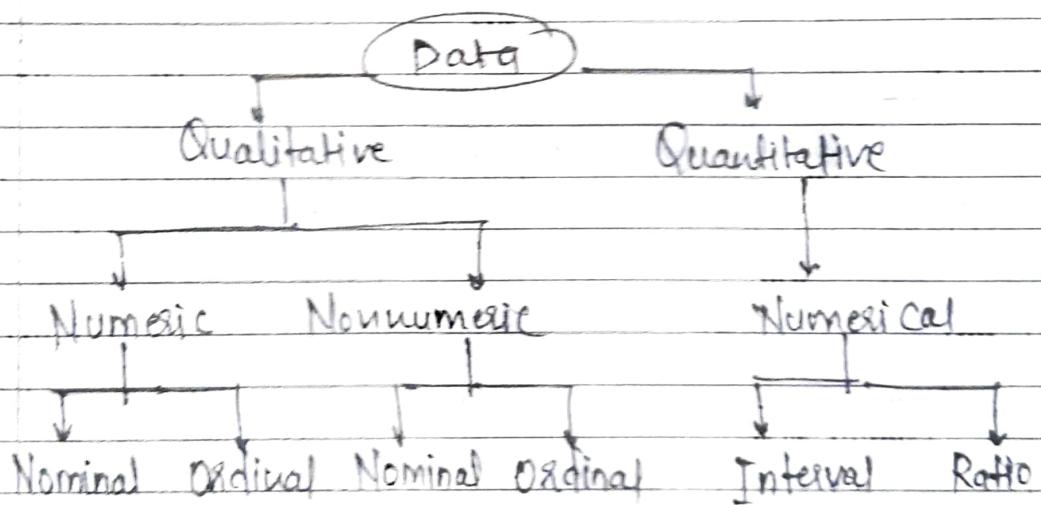
Qualitative & Quantitative data.

- i) Data can be further classified as being qualitative or quantitative.
↳ Qualitative Data:-

- Labels or names used to identify an attribute of each element.
- Often referred to as categorical data.
- Use either the nominal or ordinal scale of measurement
- Can be either numeric or nonnumeric.

Quantitative Data:-

Discrete, if measuring how many
Continuous, if measuring how much
Quantitative data are always numeric.
Ordinary arithmetic operations are meaningful for quantitative data.



(Q2) What are the various ways to present summary of Categorical data & Numerical data.

	Class	class mid pt	Frequency
10 but less than 20		15	3
20 but less than 30		25	6
30 but less than 40		35	5
40 but less than 50		45	4
50 but less than 60		55	2
60 but less than			

Plot polygon, histogram & ogive graph for above mentioned data.



Categorical Data

↓
Tabulating Data
↓
Summary Table

↓
Graphing Data
↓
Barcharts Piecharts Pareto Charts

Numerical Data

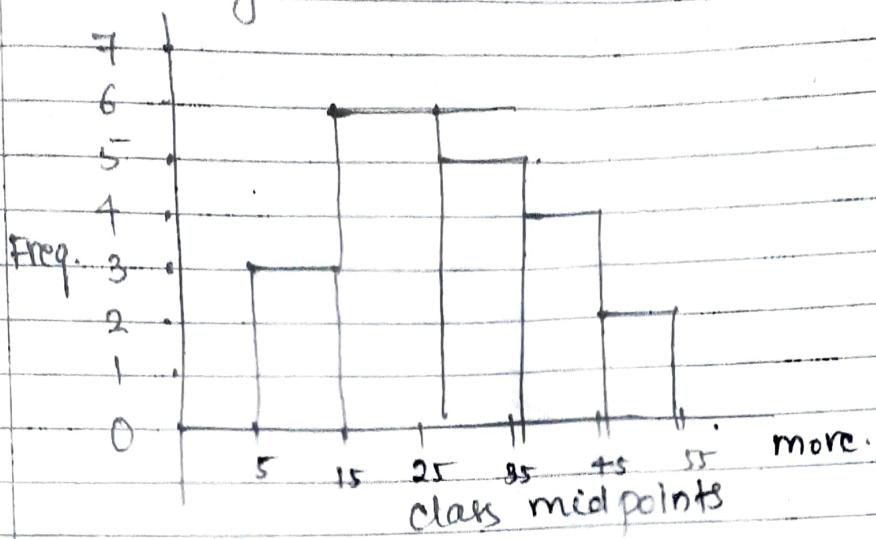
↓
Ordered Array
↓
Stem & Leaf Display

↓
Freq. Distributions &
Cumulative Distributions

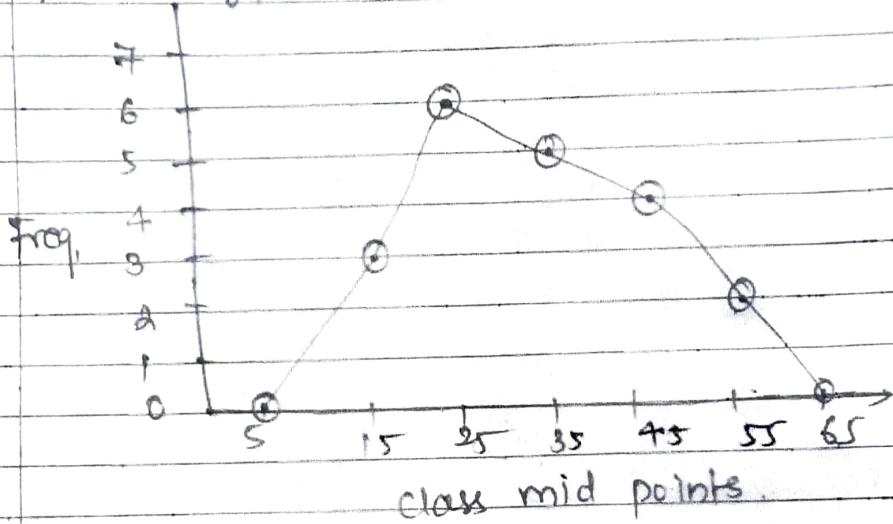
↓
Histogram Polygon Ogive

Class	Freq.	mid pt	Relative Freq.	%	Cum. Freq	Cum %
$10 \geq 20$	3	15	0.15	15	3	15
$20 \geq 30$	6	25	0.30	30	9	45
$30 \geq 40$	5	35	0.25	25	14	70
$40 \geq 50$	4	45	0.20	20	18	90
$50 \geq 60$	2	55	0.10	10	20	100
Total	20	—	1.00	100		

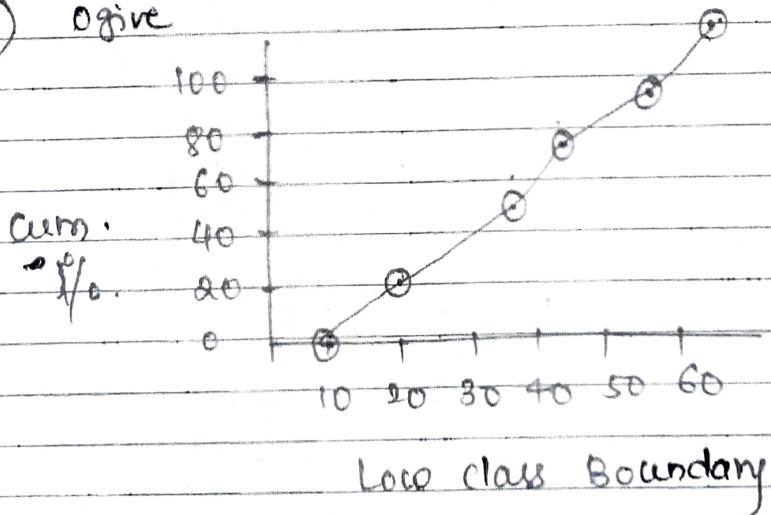
a) Histogram :-



b) Polygon .



c) ogive



Q3) Define discrete prob. distribution Function.
Explain Binomial & Poisson's Distribution in details (eg. formula for expected mean, Variance & Standard deviation)

⇒ Discrete Prob. distributions:-

The prob. distribution for a random variable describes how probabilities are distributed over the values of the random variable.

- we can describe a discrete prob. distribution with a table, graph or equation.
- The prob. distribution is defined by a prob. function, denoted by $P(x)$, which provides the prob. for each value of the random variable.
- The required conditions for a discrete probability functions are:

$$P(x) \geq 0$$

$$\sum P(x) = 1.$$

a) Binomial Distribution:- Describes discrete, not continuous, data resulting from an experiment known as Bernoulli process.

Four properties of a Binomial Exp:

1. The experiment consists of a sequence of n identical trials.
2. Two outcomes, success & failure are possible on each trial.
3. The prob. of a success, denoted by p , does not change from trial to trial.
4. The trials are independent.

- The number of successes occurring in the ~~n~~ trials.
- We let x denote the no. of successes occurring in the n trials.

Binomial prob. Functions :-

$$P(x) = \frac{n!}{x!(n-x)!} p^x (1-p)^{(n-x)}$$

where :

$P(x)$ = the prob. of x success in n trials

n = the no. of trials.

p = the prob. of success in any one trial

Expected value :- $E(x) = \mu = np$

Variance (σ^2) = $\sigma^2 = np(1-p)$

Standard Deviation (σ) = $\sqrt{np(1-p)}$

2] Poisson Distribution:-

- The poisson distribution & the binomial distribution have some similarities but also several differences.
- The binomial distribution describes a distribution of two possible outcomes designated as success & failure from a given number of trials.
- The poisson distribution focuses only on the number of ~~trials~~ discrete occurrences over some interval or continuum.
- The poisson distribution describes the

Occurrence of rare events. In fact, the Poisson formula has been referred to as the law of improbable events.

$$f(x) = \frac{e^{\mu} e^{-\mu}}{x!}$$

$f(x)$ = prob. of x occurrences in an interval
 μ = mean number of occurrences in an interval

$$\epsilon = 2.71828$$

$$[\mu = 6^2]$$

A property of the Poisson distribution is that the mean & variance are equal.

Q2)

A	Class	Freq(f)	mid pt (x)	$f x^2$	$x_i - \bar{x}$
	10.0 - 10.9	1	10.5	10.5	4.07 - 4.07
	11.0 - 11.9	4	11.5	46.0	-3.21
	12.0 - 12.9	6	12.5	75.0	-2.21
	13.0 - 13.9	8	13.5	108.0	-1.21
	14.0 - 14.9	12	14.5	174.0	-0.21
	15.0 - 15.9	11	15.5	170.5	0.29
	16.0 - 16.9	8	16.5	132.0	1.29
	17.0 - 17.9	7	17.5	122.5	2.29
	18.0 - 18.9	6	18.5	111.0	3.29
	19.0 - 19.9	2	19.5	89.0	4.29
	Total.	65		988.5	

a) $\bar{x} = \frac{\sum (fx^2)}{n} = \frac{988.5}{65} = 15.2077$ pounds.

b) Stand. dev. = $\sqrt{\frac{\sum (x_i - \bar{x})^2}{n-1}} = \sqrt{\frac{82.941}{9}} = \sqrt{9.215} = 3.035$

c) Variance = $s^2 = 0.215 = 9.215$

Q2 B) Find outlier ?

850	875	9000	9900	5300	5700	6700	7300	7700	8100
8300	8400	8700	8700	8900	9300	9500	9500	9700	10000
10300	10500	10700	10800	11000	11300	11300	11800	12700	12900
13100	13500	13800	14900	16300	17200	18500	20300	21310	213

$$Q_1 = 8100.$$

$$\begin{aligned} Q_1 &= \frac{1}{4} \times n = 10^{\text{th}} = \\ &= \frac{40}{+} \rightarrow 10 = \underline{\underline{8100}} \end{aligned}$$

$$Q_3 = \frac{3}{4} \times n = \frac{3}{4} \times 40 = 30$$

$$[Q_3 = 12900]$$

$$IQR = 12900 - 8100 = 4800$$

$$1.5 \times IQR = 7200.$$

$$1.5 \times 4800 = 7200$$

$$\text{outliers} = Q_1 - 7200 = 8100 - 7200 = 900$$

$$\text{outliers} = Q_3 + 7200 = 12900 + 7200 = 20100$$

Any data point below 900 & above 20100
are outliers.

→ Outlier :- An outlier is a single data point that goes far outside the average value of a group of statistics.
Outlier may be exceptions that stand outside individual samples of populations as well.

Locating quartiles ..

Find a quartile by determining the value in

the appropriate position in the ranked data.

$$Q_1 = \text{first quartile position} = Q_1 = \frac{(n+1)}{4} \text{ ranked value}$$

$$Q_2 = \frac{n+1}{2}$$

$$Q_3 = \frac{3(n+1)}{4} \quad \text{where } n = \text{number of observed values}$$

$$\text{Interquartile range} = Q_3 - Q_1$$

five number summary

x_{smallest}

First quartile (Q_1)

Median (Q_2)

Third quartile (Q_3)

x_{largest}

(Q2c)	Gender	First time offender	Repeat offender
	Male	60	70
	Female	44	76
	Total	104	146

a) The prob. that the shoplifter is male

$$P(M) = \frac{60+70}{250} = 0.520$$

b) The prob. that the shoplifter is first time offender, given that the shoplifter is male

$$P(F|M) = \frac{P(F \cap M)}{P(M)} = \frac{60/250}{130/250} = 0.465$$

c) $P(W|R)$ The prob. that the shoplifter is female, given that the shoplifter is a repeat offender.

$$P(W|R) = \frac{P(W \cap R)}{P(R)} = \frac{\left(\frac{76}{250}\right)}{\left(\frac{146}{250}\right)} = 0.521.$$

d) The prob. that the shoplifter is female, given that the shoplifter is a first-time offender.

$$P(W|F) = \frac{P(W \cap F)}{P(F)} = \frac{\left(\frac{44}{250}\right)}{\left(\frac{104}{250}\right)} = 0.423$$

e) The prob. that the shoplifter is both male & a repeat offender.

$$P(M \text{ and } R) = \frac{70}{250} = 0.280.$$

Q.3 A) Three types of probabilities

- 1] Classical approach
- 2] Relative frequency approach
- 3] Subjective approach

i] Classical approach:-

prob. of an event = (no. of outcomes where the event occurs)

(Total no. of possible outcomes)

$$P(H) = \frac{1}{(1+1)} = \frac{1}{2}. \quad \begin{matrix} \text{Prob of getting} \\ \text{head when} \end{matrix}$$

coin is tossed.

VITVidyalankar
Institute of
Technology

(Accredited A+ by NAAC)

(Autonomous Institute Affiliated to University of Mumbai)

Mid Semester Examination

Branch	Date	Sem.	Roll No./ Exam Seat No.	Subject	Student's Signature	Junior Supervisor's Name and Sign

Question No.	A	B	C	D	E	F	G	H	Total	Total out of (20 / 30 / 40)
1										
2										
3										
4										

Examiners Signature	Student's Sign (After receiving the assessed answer sheet)

$$P(5) = 1/6 \text{ for the dice rolling}$$

Classical prob is also called as a priori prob. because we don't need to perform experiments.

2] Relative Frequency (RF):-

Live upto 85 years, plant near river will substantially kill the fish.

We need experiment to answer these. This method uses the relative frequencies of past occurrences as probabilities.

- How often something has happened in past - we predict future.
- more trials greater accuracy.

3] Subjective probability :- Based on belief, experience, when event has occurred once

or few times.

Because most higher-level & social & managerial decisions are concerned with specific, unique situations, rather than with a long series of identical situations, decision makers use this prob.

probability under conditions of Statistical independence :-

→ Statistical independence : The occurrence of one event has no effect on the prob. of occurrence of any other event.

1. Marginal prob. under statistical independence :-

Tossing of fair coin. Outcome of second toss is independent of outcome of first toss.

2. Joint prob. under statistical independence :-

The prob. of two or more independent events occurring together or in succession is the product of their marginal probabilities.

Joint prob. of two independent events :

$$P(AB) = P(A) \times P(B)$$

$P(AB)$ = prob. of events A & B occurring together, this is known as a joint prob.

$P(A)$ = marginal prob. of event A occurring

$P(B)$ = marginal prob. of event B occurring

$$P(H_1 H_2) = 0.5 \times 0.5 = 0.25.$$

This is the prob. of heads in two successive tosses)

- 3) Conditional probabilities under statistical independence :-

It is written as $P(B/A)$. The prob. of event B given A has occurred.

$$P(B/A) = P(B)$$

What is the prob. that the second toss of a fair coin will result in heads, given that heads resulted in first toss.

$P(H_2/H_1)$, we know that independence means the first toss's result would not affect the result of second toss.

$$P(H_2/H_1) = 0.5$$

Summary

	Types of prob. & symbol	Formula
Prob. under statistical independence	marginal	$P(A)$
	Joint	$P(AB)$
	Conditional	$P(B/A)$
		$P(B)$

Probabilities under conditions of statistical Dependence :

- When prob. of some event is dependent on or affected by the occurrence of some other event.

i) Conditional Prob. under statistical Dependence:
Assume a box has 10 balls as follows:

Event	Prob. of event	
1	0.1	colored & dotted
2	0.1	
3	0.1	
4	0.1	colored & striped
5	0.1	gray & dotted
6	0.1	gray & striped
7	0.1	
8	0.1	gray & striped
9	0.1	
10	0.1	

e.g. If a colored ball is drawn.

1. what is the prob. that it is dotted

$$P(D/c) = 0.3/0.4$$

2. What is the prob. that it is striped

$$P(S/c) = 0.1/0.4$$

conditional probabilities for statistical dependent events.

2. Joint prob. under statistical dependence:

We know that conditional prob. under statistical dependence:

$$P(B/A) = P(BA) / P(A)$$

$$P(BA) = P(B/A) * P(A)$$

Joint prob
of events B
& A

Prob. of
event B
given that A

has happened

Prob. of
events A.

$$\text{What is } P(D|G) = P(DG) / P(G)$$

$$= 0.2 / 0.6$$

$$= \underline{\underline{1/3}}.$$

- 3) Marginal probabilities under statistical dependence:
Are computed by summing up the prob. of all the joint events in which the simple event occurs.

$$P(c) = P(cd) + P(cs) = 0.3 + 0.1 = 0.4$$

- (Q3) b) Explain covariance & significance of coefficient of correlation method.
for following data a) Compute covariance.
b) Compute coefficient of correlation
c) Which is valuable in expressing relationship?
d) What conclusion can you reach about relationship?

1) \rightarrow Relationship between two numerical variable
1) Covariance.
2) Coefficient of Correlation.

The covariance measures the strength of the linear relationship between two numerical variables (x & y).

The sample covariance :

$$\text{cov}(x, y) = \frac{\sum_{i=1}^n (x_i - \bar{x})(y_i - \bar{y})}{n-1}$$

- only concerned with the strength of the relationship
- No causal effect is implied.

Interpreting Covariance :-

Covariance between two variables :

$\text{cov}(x, y) > 0 \rightarrow x \text{ & } y \text{ tend to move in the same direction}$

$\text{cov}(x, y) > 0 \rightarrow x \text{ & } y \text{ tend to move in opposite direction.}$

$\text{cov}(x, y) = 0 \rightarrow x \text{ & } y \text{ are independent}$

The covariance has a major flaw :

It is not possible to determine the relative strength of the relationship from the size of the covariance.

Coefficient of Correlation

- Measures the relative strength of the linear relationship between two numerical variables
- Sample coefficient of correlation :

$$r = \frac{\text{cov}(x, y)}{S_x S_y}$$

Where,

$$\text{cov}(x, y) = \frac{\sum_{i=1}^n (x_i - \bar{x})(y_i - \bar{y})}{n-1}$$

$$S_x = \sqrt{\frac{\sum_{i=1}^n (x_i - \bar{x})^2}{n-1}}$$

$$S_y = \sqrt{\frac{\sum_{i=1}^n (y_i - \bar{y})^2}{n-1}}$$

- The population coefficient of correlation is referred to as ρ .
- The sample coefficient of correlation is referred to as r .
- Either ρ or r have the following features:

Ranges between $-1 \leq r \leq 1$

- The closer to -1 , the stronger the negative linear relationship.
- The closer to 1 , the stronger the positive linear relationship.
- The closer to 0 , the weaker the linear relationship.

x (calories)	y (fat)	$x_i - \bar{x}$	$y_i - \bar{y}$	$a \times b$
1 240	8	-140	-7.78	1089.2
2 260	3.5	-120	-12.28	1473.6
3 350	22	-30	6.22	-186.6
4 350	20	-30	4.22	-126.6
5 420	16	40	0.22	8.8
6 510	22	130	6.22	808.6
7 530	19	150	3.22	483

$$\bar{x} = 380$$

$$\bar{y} = 15.78$$

$$\text{Cov}(x, y) = \frac{\sum_{i=1}^n (x_i - \bar{x})(y_i - \bar{y})}{n-1}$$

$$[\text{Cov}(x, y) = 591.66]$$

$$S_x = \sqrt{\frac{\sum_{i=1}^n (x_i - \bar{x})^2}{n-1}} = 113.14$$

$$s_y = \sqrt{\frac{\sum_{i=1}^n (y_i - \bar{y})^2}{n-1}}$$

$$= 7.26,$$

$$r = \frac{\text{cov}(x, y)}{s_x s_y}$$

$$= \frac{591.66}{113.14 \times 7.26}$$

$$\boxed{r = 0.72}$$

c) which is valuable in expressing relationship
 r = Correlation

d) $r = 0.72$ is more closer to '+1' hence
it is strong positive relationship