

# Big Data Analytics

22 October 2024 11:27

\* Empirical formulae

$$2 \times \text{mean} + \text{mode} = 3 \times \text{median}$$

or

$$3 \times \text{median} - 2 \times \text{mean} = \text{mode}$$

$$\text{Median} \rightarrow \text{grouped data} = l + \left\{ \frac{n/2 - cf}{f} \right\} \times h$$

$$\text{mode} \rightarrow = l + \left\{ \frac{f_1 - f_0}{2f_1 - f_0 - f_2} \right\} \times h$$

Measures of

Dispersion =



most scattered  
and dispersed value.

→ Dispersion in Statistics:

It is the way of scattering the data through which we can easily find that how far does the data is from the average value, or how close.

It shows variability & consistency of data.

Measures of Dispersion: Eg. - Range, variance & Standard deviation

↳ Dispersion is of 2 types:

1) Absolute measure of Dispersion.

2) Relative measure of Dispersion.

→ Absolute

Range, Variance, Standard Deviation, Quantile Deviation  
& Mean Deviation.

→ Relative → coeff of range, coeff of variance, coeff of mean Dev.

\* Unit of data → expressed in Imperial

coeff of range, coeff of variance, coeff of mean Dev.

↓  
Absolute measure of Dispersion.

... Relative → 2 values vary from that we express 1, 2, ... data

Above:

\* Relative  $\rightarrow$  2 values vary 8 from their mean  
& measure by taking reference of 2 data  
& also helps in scattering of data.

\* Range  $\rightarrow$  10, 20, 15, 0, 100  
 $100 - 0 = 100$   
↳ ungrouped data.

e.g.: 20, 24, 31, 17, 45, 39, 51, 61  
Range = 61 - 17 = 44  $\rightarrow$  Range.

\* Grouped data!:-

	f
0 - 10	5
10 - 20	8
20 - 30	15
30 - 40	9

Range = 0 - lowest  
40 - highest  
 $\rightarrow$  highest - lowest  
 $\rightarrow 40 - 0 = 40$

\* Mean deviation:-

$$M.D = \frac{|x_1 - \bar{x}| + |x_2 - \bar{x}| + \dots + |x_n - \bar{x}|}{n}$$

(for ungrouped data)

$\bar{x} \rightarrow$  mean

e.g. - -5, 10, 25

$$\Rightarrow \bar{x} = -5 + 10 + 25 = 10$$

$$\bar{x} = 10$$

$$M.D = \frac{|-5 - 10| + |10 - 10| + |25 - 10|}{3}$$

$$\Rightarrow \frac{-15 + 0 + 15}{3} = 0 \rightarrow M.D = 0$$

\* e.g. - mean deviation for 2, 4, 6, 8, 10  
 $\bar{x} = 2 + 4 + 6 + 8 + 10 / 5 = 6$   $\bar{x} = 6$   $\rightarrow$  mean

$$M.D = \frac{\sum |x_i - \bar{x}|}{n} = \frac{|x_1 - \bar{x}| + |x_2 - \bar{x}| + \dots + |x_n - \bar{x}|}{n}$$

$$\dots n - |2 - 6| + |4 - 6| + |6 - 6| + |8 - 6| + |10 - 6|$$

$$M.D = \frac{1}{5} |2-6| + |4-6| + |6-6| + |8-6| + |10-6|$$

$$\therefore \frac{4+2+0+2+4}{5} = \frac{12}{5} = 2.4$$

M.D = 2.4

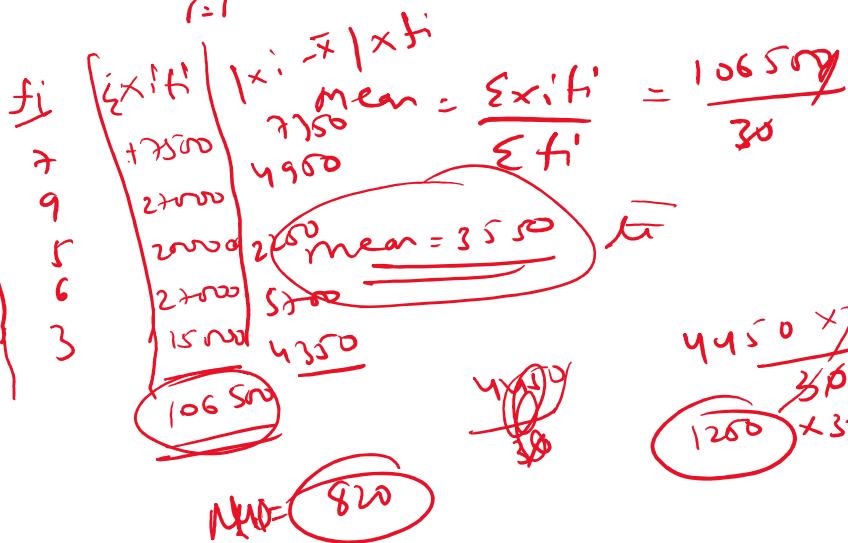
\* Mean Deviation for discrete frequency distribution! -

$$M.D = \frac{\sum_{i=1}^n f_i |x_i - \bar{x}|}{\sum_{i=1}^n f_i}$$

$$\frac{7350 + 4900 + 2100 + 5200 + 3500}{50}$$

Eg:- x: wages

-1050	2500 - 3650
-550	3000 - 3500
+450	4000
950	4500
1450	5000
<u>1250</u>	<u>1250</u>



\* Continuous  $\rightarrow M.D$

$$\begin{aligned} \text{Eg:- } & 10-20 = \frac{10+20}{2} = 15 \\ & 20-30 = 50\% \cdot 20 = 10 \\ & 30-40 = 20\% \cdot 25 = 13 \\ & 40-50 = 10\% \cdot 25 = 12 \end{aligned}$$

$x_i$	$f_i$	$x_i f_i$
15	15	225
25	10	250
35	13	455
45	12	540

$$\mu = \frac{\sum x_i f_i}{\sum f_i}$$

$$\mu = \frac{225 + 250 + 455 + 540}{50} = \frac{1470}{50} = 29.4$$

$$M.D = \frac{|15-29.4|, |25-29.4|, |35-29.4|, |45-29.4|}{50}, \underline{14.4, 4.4, 5.6, 15.6}$$

$$M.D = \frac{110 - (11.71) | 125 - 9.4 | + 120 - 15.6}{14.7, 4.4, 5.6, 15.6}$$

$(x_i - \bar{x})$	$f_i$	$\frac{\sum (x_i - \bar{x}) f_i}{\sum f_i}$
14.7	15	21.6
4.4	10	4.4
5.6	13	72.8
15.6	12	187.2
		<u><u>520</u></u>

$$M.D = \frac{\sum f_i (x_i - \bar{x})}{\sum f_i} = \frac{520}{50} = 10.4$$

\*  $S.D = \sqrt{\frac{\sum (x_i - \bar{x})^2}{N}}$

$$\text{variance}(\sigma^2) = \frac{\sum (x_i - \bar{x})^2}{N}$$

\* SD of ungrouped data:-

$$S.D \text{ of } = \sqrt{\frac{\sum (x_i - \mu)^2}{N}}$$

Q:-  $x = \{2, 3, 4, 5, 6\}$

$x_i$	$\frac{x_i - \mu}{N}$	$\frac{(x_i - \mu)^2}{N}$
2	2 - 4 = -2	4
3	-1	0
4	0	1
5	1	4
6	2	10

$$S.D = \sqrt{\sum_{i=1}^n \left( \frac{x_i - \mu}{N} \right)^2}$$

$$\mu = \frac{2+3+4+5+6}{5} = 4$$

$$\sqrt{\frac{10}{5}} = \sqrt{2} = \underline{\underline{1.414}}$$

\* SD by assumed mean method:-

$$\sigma = \sqrt{\frac{\sum (x_i - A)^2}{N}}$$

\* Standard deviation for discrete ungrouped data.

1) Actual mean method

2) Assumed mean method.

$$\text{Formulae} = \sqrt{\frac{\sum_{i=1}^n (x_i - \mu)^2 f_i}{N}}$$

formulae =  $\sum_{i=1}^N$

Eg:-

$x_i$	$f_i$
10	1
9	3
6	5
8	1

Standard deviation for Discrete.

$$SD = \sqrt{\frac{\sum f_i(x_i - \mu)^2}{N}}$$

$$\mu = \frac{\sum x_i f_i}{\sum f_i} = \frac{10 + 12 + 30 + 8}{10} = \frac{60}{10} = 6$$

$x_i$	$f_i$	$(x_i - \mu)$	$f_i(x_i - \mu)^2$	$f_i \sum_{i=1}^N (x_i - \mu)^2$
10	1	4	16	16
9	3	-2	4	12
6	5	0	0	0
8	1	2	4	4

$$SD = \sqrt{\frac{\sum f_i(x_i - \mu)^2}{N}} = \sqrt{\frac{32}{10}} = \boxed{1.78}$$

\* for continuous Grouped distribution :-

$$SP = \sqrt{\frac{\sum f_i(x_i - \mu)^2}{N}}$$

Eg:-

$x_i$	$f_i$	$\bar{x}_i$	$f_i \bar{x}_i$	$x_i - \mu$	$(x_i - \mu)^2$	$f_i(x_i - \mu)^2$
0-10	2	5	10	-14	196	392
10-20	4	15	60	-4	16	64
20-30	2	25	50	6	36	72
30-40	2	35	70	16	256	512

$$\text{Mean} = \frac{\sum x_i f_i}{\sum f_i} = \frac{190}{10} = 19 \quad \boxed{\mu = 19}$$

$$SD = \sqrt{\frac{\sum (x_i - \mu)^2}{N}} = \sqrt{\frac{1040}{10}} = \boxed{10.19}$$

\* Quartile deviation:-

Upper quartile deviation - Lower quartile deviation

Quartile -

$$\frac{\text{Highest deviation} - \text{lowest deviation}}{2}$$

$$Q_D = Q_3 - Q_1$$

$Q_3$  = highest deviation,  $Q_1$  = lowest deviation.

$$Q_1 = \text{size} \left[ \frac{n+1}{4} \right]^{\text{th}}$$

$$Q_3 = \text{size} \left[ 3 \left( \frac{n+1}{4} \right) \right]^{\text{th}}$$

$$\text{coff of } Q_D = \frac{Q_3 - Q_1}{Q_3 + Q_1}$$

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

Ex1 - for ungrouped data = Discrete.

$Q = 150, 150, 268, 280, 195, 140, 200$ .

Soln - Arrange in ascending order.

100 140 150 150 195 200 268 280

Find Interquartile Range, quartile deviation & coff of Quartile Deviation.

Soln - Interquartile Range =  $Q_3 - Q_1$

$$Q_1 = \left[ \frac{n+1}{4} \right]^{\text{th}} \quad n=7 \quad \left[ \frac{7+1}{4} \right]^{\text{th}} = \left[ \frac{8}{4} \right] = \text{2nd}$$

$$Q_1 = 2^{\text{nd}} = 140$$

$$Q_3 = 3 \times 2 = 6^{\text{th}} = 268$$

$$\text{Interquartile Range} = 268 - 140 = 128$$

$$\text{Quartile Deviation} = \frac{128}{2} = 64 ; \quad \left( \frac{Q_3 - Q_1}{2} \right)$$

$$\text{coff of Quartile Deviation} = \frac{Q_3 - Q_1}{Q_3 + Q_1} = \frac{128}{408}$$

$$= 0.31$$

Q

$x_i$	$f_i$	$cf$
60	25	25
62	21	46

Interquartile :-  $Q_3 - Q_1$

$$Q_1 = \left[ \frac{n+1}{4} \right]^{\text{th}}$$

$$= \frac{n=199}{4} = 50^{\text{th}}$$

60	25	188	$Q_1 = \left[ \frac{n}{4} \right]$	$\bar{=}$
62	21	46	$Q_1 = 68$	
68	20	77		
70	18	92		
75	24	112		
80	20	136		
88	21	160		
90	17	177		
97	22	197		

$$Q_3 = 3 \times 50 = 150$$

$$\text{Inter} = Q_3 - Q_1$$

$$88 - 68 = 20$$

$$\text{Quartile Deviation} = \frac{Q_3 - Q_1}{2} = \frac{20}{2} = 10$$

$$\text{Coff of Quartile Deviation} = \frac{20}{150} = 0.128$$

\* Quantile Deviation for Continuous frequency  
(Grouped)

Size	f	Cf
0 - 10	9	9
10 - 20	16	24
20 - 30	29	53
30 - 40	27	77
40 - 50	3	80
50 - 60	20	100

$$Q_1 = \left\{ \frac{N}{4} \right\}^{\text{th}} = 25^{\text{th}}$$

$$l = 20, C_f = 24, f = 29, i = 10, Q_3 = \left\{ \frac{3N}{4} \right\}^{\text{th}} = 75^{\text{th}}, l = 30, C_f = 53, f = 27$$

$$Q_1 = 20 + \left( \frac{25 - 24}{29} \right) \times 10$$

$$Q_3 = 30 + \left( \frac{75 - 53}{27} \right) \times 10$$

$$IQR = 20 + \frac{1}{4} \times 10 \\ IQR = 20.25$$

$$Q_3 = 39.16$$

$$\text{Interquartile Range} = Q_3 - Q_1 = 39.16 - 20.25 = 18.82$$

$$\text{Quartile Deviation} = \frac{Q_3 - Q_1}{2} = \frac{18.82}{2} = 9.41$$

$$\text{Coff} = \frac{18.82}{59.5} = 0.31$$

$$Q_1 = (25)^{\text{th}} = 50^{\text{th}}$$

$$Q_3 = \frac{150^{\text{th}}}{20 - 25} = 10$$

$Q$	Size	$f$	$Cf$	$\overline{Q_1} = \frac{(Cf)^m}{f_i} = \underline{50}$	$20 - 25$
0-5	6	6	6		$l=20, Cf=100$
5-10	18	24	54		$f=60$
10-15	30	100	100	$Q_1 = 10 + \left(\frac{50-24}{30}\right) \times 5 = 14.3$	$Q_3 = 20 + 3\left(\frac{150-100}{60}\right) \times 5$
15-20	40	140			$Q_3 = 20 + 12.5$
20-25	60	200			$= 32.5$
25-30	40				

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

$$32.5 - 14.3 = \underline{18.2}$$

$$\text{Quartile Deviation} = \frac{18.2}{2} = \underline{9.1}$$

$$\text{coff} = \frac{18.2}{46.8} = \underline{0.38}$$

\* coff of Range:-  $\left\{ \frac{\text{largest} - \text{smallest}}{\text{largest} + \text{smallest}} \right\}$

$$Q = 1575 = \text{largest}$$

$$1380 = \text{smallest}$$

$$\text{Range} = 1575 - 1380 = 195$$

$$\text{coff of Range} = \frac{195}{2955} = 0.065 = \frac{\text{largest} - \text{smallest}}{\text{largest} + \text{smallest}}$$

\* coff of Mean deviation:-

$$M.D = \frac{\sum f_i |x_i - \bar{x}|}{N}$$

$$\text{coff} = \frac{MD}{\bar{x}}$$

$$\frac{MD \times \bar{x}}{N}$$

Ex:- 2, 7, 6, 8, 10

$$\begin{aligned} x_i & x_i - \bar{x} & f(x_i - \bar{x}) \\ 2 - 6 & = 4 & = \frac{\sum f_i x_i}{N} = 30/5 = 6 \\ 7 - 6 & = 1 & \\ 6 - 6 & = 0 & \\ 8 - 6 & = 2 & \\ 10 - 6 & = 4 & \\ \hline & \sum f_i = 12 & \end{aligned}$$

$$\frac{\sum f_i x_i}{N} = \bar{x}$$

$$\bar{x} = 6$$

$$MD = \frac{\sum f_i |x_i - \bar{x}|}{N}$$

$\therefore$	$x_i$	$f$	(12)	$\sum f_i = n$	$\therefore$	$N$
	2500	7				$\Rightarrow 820$
	3000	9				$\text{coff} = \frac{MD}{\mu} = \frac{820}{380} = 0.22$
	4000	5				
	4500	3				
	5000	3				
				$\sum f_i = 3883$		
				$\frac{11686}{3883} = 3$		

\* off of mean deviation in continuous (grouped) :-

$x_i$	$x_i - \bar{x}$	$f$	$\sum f_i = 30$	$\sum f_i x_i = 225 + 250 + 455 + 540$
15-20	15-29.4	15		
25-30	25-29.4	10		
35-40	35-29.4	13		
45-50	45-29.4	12		
			$\bar{x} = 29.4$	
			$\sum f_i = 30$	

$x_i - \bar{x}$	$ x_i - \bar{x} $	$f$	$\sum f_i  x_i - \bar{x}  = 524$	$MD = \frac{524}{30} = 10.4$
15-29.4	15.4	15		
25-29.4	4.4	10		
35-29.4	5.6	13		
45-29.4	15.6	12		
			$\sum f_i = 30$	
			$MD = 10.4$	
			$coff = \frac{10.4}{29.4} = 0.35$	

\* coefficient of variance :-

= Standard Deviation  $\times 100$

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

$$\therefore 2+3+4+5+6/5 = 4 \quad (\bar{x} = 4)$$

$$2-4, 3-4, 4-4, 5-4, 6-4$$

$$\rightarrow 2 \quad 1 \quad 0 \quad 1 \quad 2$$

$$SD = \sqrt{\frac{30}{30}} = 1.5$$

$$SD = \sqrt{1.414}$$

$$4 \quad 1 \quad 0 \quad 1 \quad 2$$

$$\therefore SD/\text{mean} = \frac{1.414}{4} = 0.3535$$

$x_i$	$f_i$	$16$	$16$	$SD = \sqrt{\frac{\sum f_i (x_i - \bar{x})^2}{n}}$
10-14	1			
9-13	4			
8-12	0			
7-11	4			
				$\sqrt{\frac{32}{4}} = \sqrt{8} = 2\sqrt{2}$

$$\frac{\sum f_i}{n} = \frac{10+12+30+8}{40}$$

$$\bar{x} = 6$$

$$\begin{array}{c|c|c} 4 & 3 \\ \hline 5 & 5 \\ \hline 6 & 1 \end{array}$$

$$0$$

$$0$$

$$\sqrt{\frac{32}{4}} = \sqrt{8} = 2\sqrt{2}$$

$$\begin{array}{l} Q = 5 = 0 - 10 \\ 15 = 10 - 20 \\ 25 = 20 - 30 \\ 35 = 30 - 40 \end{array}$$

$f_i$	$x_i$	$f_i \cdot x_i$
2	10	20
4	20	80
2	30	60
2	40	80
	10	70
		190
$\mu = 19$		

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

$$SD = 10 \cdot 178$$

$$coff = SD / mean = 10 \cdot 178 / 19 = 0.92$$

$$\begin{array}{l} Q = SD = 8.5 \\ \mu = 14.5 \\ coff = \frac{8.5}{14.5} \Rightarrow 100 = 58.6 \% \end{array}$$

$$\begin{array}{l} Q = SD = 1.7 \\ \sigma = 26.5 \\ coff = \frac{SD}{mean} = \frac{1.7}{26.5} = 0.052 \times 100 = 52.8 \% \end{array}$$

$$\begin{array}{l} Q = \mu = 13 \\ coff = 38 \\ \cancel{coff} = \frac{SD}{mean} \times 100 = \frac{38}{13} \times 100 \\ m = \frac{13 \times 100}{38} = 34.2 \end{array}$$

$$\underline{SD = 4.94}$$

$$\begin{array}{l} Q = \frac{SD}{mean} = coff \\ \frac{10}{65} = 15 \% \rightarrow mean \\ \frac{12}{70} = 17.14 \% \end{array}$$

$$\frac{12}{60} = 20\%$$

$$\frac{14}{57} = 24\%$$

$\rightarrow$  Economics

\* Null hypothesis:-

$$H_0$$

opposite condition

$$M_0 = M_1$$

\* Alternative hypothesis:-

$$M_1 \neq M_0$$

$$H_a$$

\* One tailed test  $\rightarrow$  left  $\rightarrow H_0: M \geq 50, H_a: M < 50$

$\rightarrow$  right  $\rightarrow H_0: M \leq 50, H_a: M > 50$

\* Two tailed test:- In which we include both hypothesis & assuming them

\* Type I & II errors:-

i) Type I  $\rightarrow$  null hypothesis  $\checkmark$  but we rejected ( $\alpha$ )

ii) Type II  $\rightarrow$  null hypothesis  $\times$  but we accepted ( $\beta$ )

\* Z testing:-

$$Z = \frac{\bar{x} - \mu}{\frac{\sigma}{\sqrt{n}}}$$

$\left. \begin{array}{l} \bar{x} = \text{sample mean} \\ \mu = \text{population mean} \\ \sigma = \text{standard deviation} \\ n = \text{size} \end{array} \right\}$

Q Data  $n=25, M=200, \sigma = 5 \text{ mg/1DL}$

Sol:- firstly we accept null hypothesis.

$H_0$ : Average cholesterol is 200

$H_a$ :  $\bar{x}$  is not 200

$$\text{Significance} = 0.05$$

$$Z = \frac{203.8 - 200}{\frac{5}{\sqrt{25}}} = 3.8$$

$\frac{5}{\sqrt{25}}$  rejected null hypothesis

$$S = \sqrt{\frac{\sum (x_i - \bar{x})^2}{n}}$$

$$Z_{\text{tabulated}} = 0.05$$

$$-1.96 \text{ to } 1.96$$

$Z = 2.98$   $\rightarrow$  2nd compart.

\* Standard deviation value when not given in t-test.

$$t = \frac{\bar{x} - \mu}{\sigma}$$

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

$$S = \sqrt{\frac{1}{n-1} \left\{ \sum x^2 - \frac{(\sum x)^2}{n} \right\}}$$

→ standard deviation

Now + test formulas:

$$t_{\text{critical}} = \frac{\bar{x} - \mu}{\frac{s}{\sqrt{n}}} \quad \left\{ \begin{array}{l} \text{when } s \text{ is not given} \\ \text{when } s \text{ is already given.} \end{array} \right.$$

or

$$\frac{\bar{x} - \mu}{\frac{s}{\sqrt{n-1}}} \quad \left\{ \begin{array}{l} \text{when } s \text{ is already given.} \\ \text{when } s \text{ is not given} \end{array} \right.$$

122.1

$$\begin{aligned} m \rightarrow x_{\text{after}} - x_{\text{before}} \\ \approx 118.2 - 122.1 = -3.9 \end{aligned}$$

$$s = 1.37$$

$$\boxed{m = -3.9}$$

$$t_{\text{critical}} = \frac{\bar{x} - \mu}{\frac{s}{\sqrt{n}}} \quad \frac{3.9}{1.37 / \sqrt{10}}$$

$$S = \sqrt{\frac{1}{n-1} \left[ \sum x^2 - \frac{(\sum x)^2}{n} \right]}$$

For 2 values

$$S = \sqrt{\frac{1}{n_A + n_B - 2} \left[ \sum x^2 - \frac{(\sum x)^2}{n_A} - \frac{(\sum x_B)^2}{n_B} \right]}$$

\* chi square test:-

$$\chi^2 = \sum \frac{(O_i - E_i)^2}{E_i}$$

	water/air	γ	N	foot
m	40	44	40	165
F	178	38	216	

$$E_i = \frac{RT \times CT}{M \rightarrow \gamma}$$

$$\frac{184 \times 318}{400} = 0.46$$

	F	$\sum_{i=1}^n x_i = 178$	$\sum_{i=1}^n \frac{x_i - \bar{x}}{E_i} = 38$	$\sum_{i=1}^n \frac{(x_i - \bar{x})^2}{E_i} = 216$
M	Y	140	146	36
M	N	44	38	36
F	Y	178	172	36
F	N	38	44	36

$$\frac{184 \times 310}{400} =$$

$M \rightarrow N$

$$\frac{82 \times 184}{400} = 36$$

$$F \rightarrow Y \quad \frac{318 \times 216}{400} = 172$$

$$F \rightarrow N \quad \frac{82 \times 216}{400} = 44$$

$$\frac{\sum (x_i - \bar{x})^2 / E_i}{\sum (x_i - \bar{x})^2}$$

\* F test hypothesis

$$F = \frac{s_1^2}{s_2^2}$$

$$\left\{ \begin{array}{l} s_1^2 = \frac{\sum (x_A - \bar{x})^2}{n_A - 1} \\ s_2^2 = \frac{\sum (x_B - \bar{x})^2}{n_B - 1} \end{array} \right\}$$

$$Q = \frac{x_1}{20} \quad \left| \quad \begin{array}{l} x_1 \\ 16 \\ 26 \\ 28 \\ 23 \\ 22 \end{array} \right.$$

when  $(\bar{x})$  is not given

$$s_1^2 = \frac{n_1}{n_1 - 1} (s_1)^2, \quad s_2^2 = \frac{n_2}{n_2 - 1} (s_2)^2$$