

1.

a) Total number of customers = $(31 + 78 + 49 + 81 + 117 + 13) = 369$

So, a total of 369 customers were served on this day at this gas station.

b) Table for class width and class midpoints,

Class boundary	Class width	Class midpoint
0 to <4	4	2
4 to <8	4	6
8 to <12	4	10
12 to <16	4	14
16 to <20	4	18
20 to <24	4	22

c) Table for relative frequency and percentage,

Class Boundaries	f	Relative frequency	Percentage
0 to <4	31	0.08	8.40
4 to <8	78	0.21	21.14
8 to <12	49	0.13	13.28
12 to <16	81	0.22	21.95
16 to <20	117	0.32	31.71
20 to <24	13	0.04	3.52
	Σf	Σ	Σ
	= 369	= 1	= 100

d) Total number of customers purchased 12 gallons or more = $(81 + 117 + 13) = 211$

So, the percentage =

$$\frac{211}{\Sigma f} \times 100\% = \frac{211}{369} \times 100\% = 57.18\%$$

e) Since, 10 is not a boundary value, therefore, we cannot determine exactly how many customers purchased 10 gallons or less.

2.

a) Table for class limits,

Class limit
1 to 200
201 to 400
401 to 600
601 to 800
801 to 1000
1001 to 1200

b) Table for class boundary and class midpoints,

Class	Class boundary	Class midpoint
1 to 200	0.5 to <200.5	100.5
201 to 400	200.5 to <400.5	300.5
401 to 600	400.5 to <600.5	500.5
601 to 800	600.5 to <800.5	700.5
801 to 1000	800.5 to <1000.5	900.5
1001 to 1200	1000.5 to <1200.5	1100.5

3.

a) Minimum value of data = 2.2

Maximum value of data = 22.8

Let's consider, class width = 5

So, number of classes =

$$\frac{22.8 - 2.2}{5} = 4.12 \approx 5$$

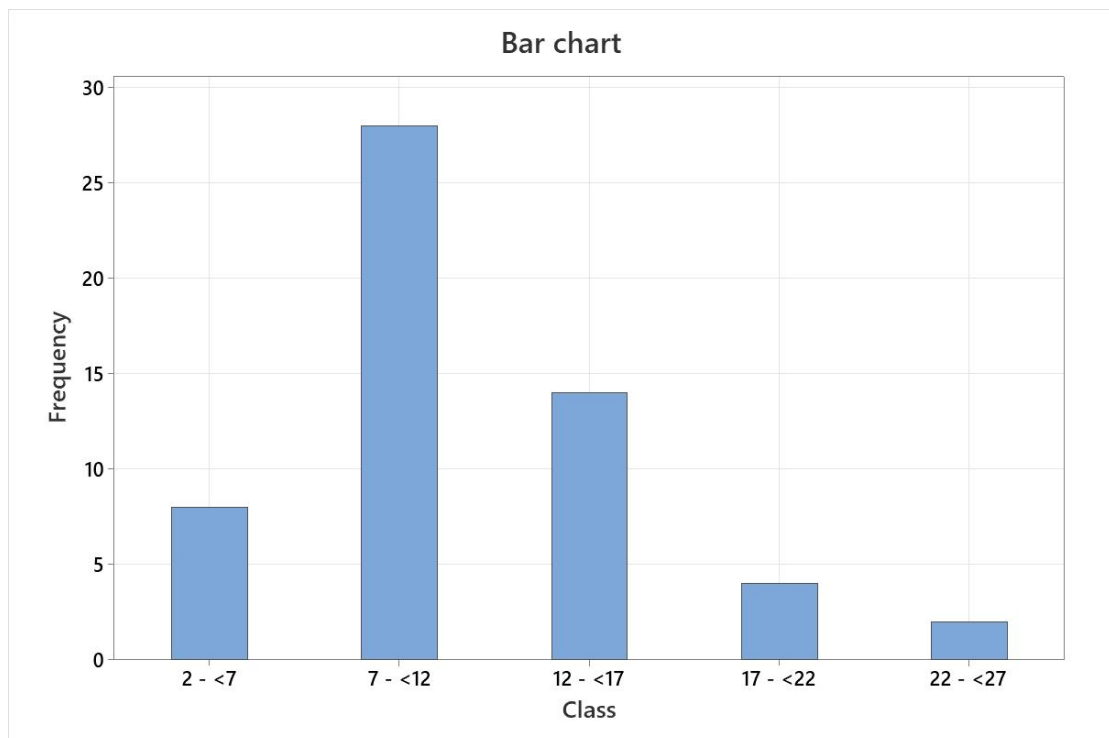
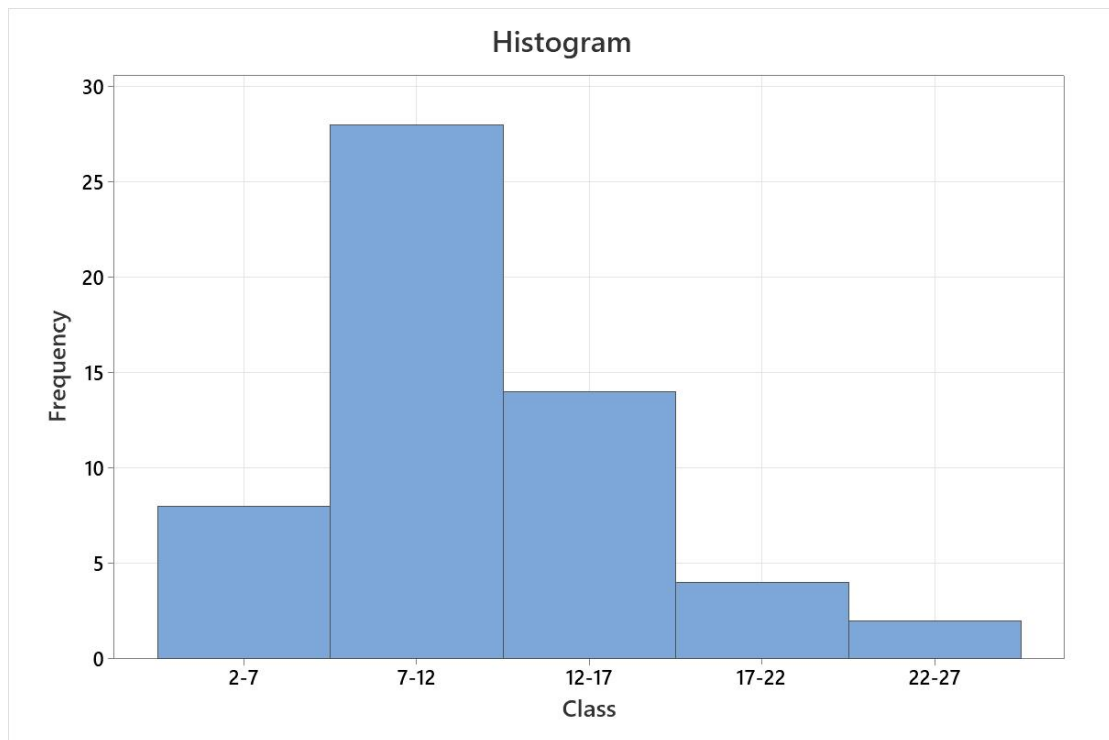
Table for frequency distribution,

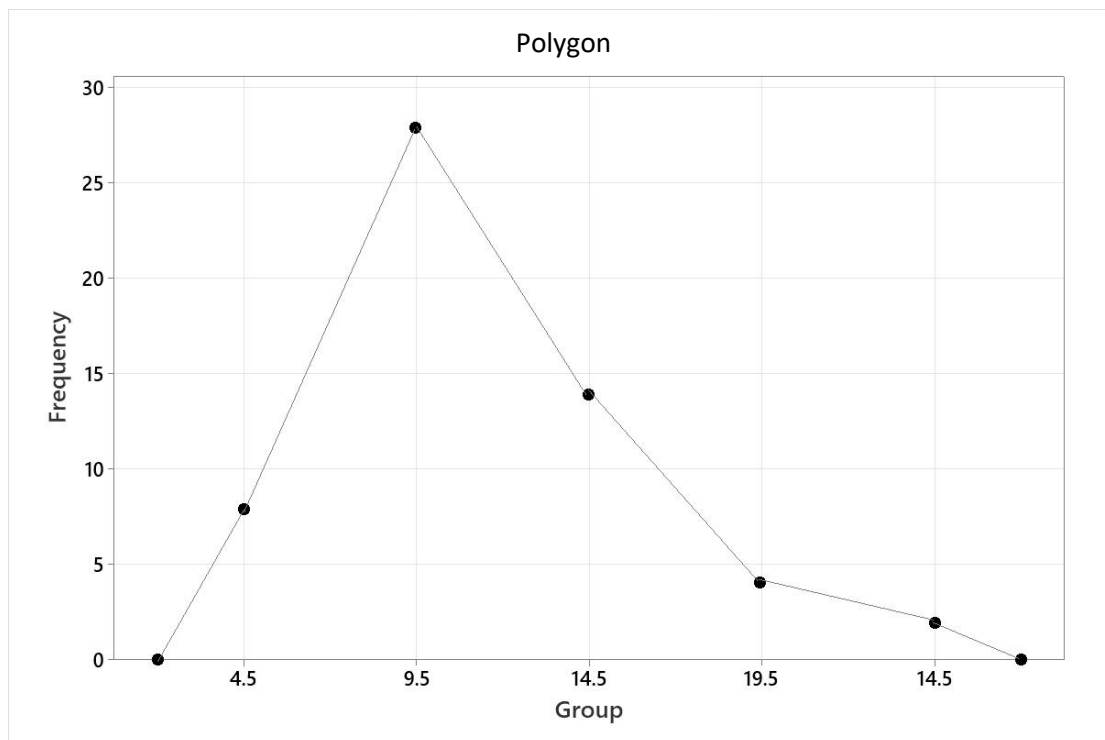
Class boundary	Tally	f
2 to <7		8
7 to <12		28
12 to <17		14
17 to <22		4
22 to <27	 	2
		Σf = 56

b) Table for relative frequency and percentage,

Class boundary	f	Relative frequency	Percentage
2 to <7	8	0.14	14.29
7 to <12	28	0.50	50.00
12 to <17	14	0.25	25.00
17 to <22	4	0.07	7.14
22 to <27	2	0.04	3.57
Σf = 56		Σ = 1	Σ = 100

c) Histogram, bar diagram and a polygon for the birth rate percentage distribution,





d) Data less than 11 =

2.2 3.1 4.6 5.4 5.6 5.6 6.5 6.6 7.4 7.7 7.8 8.1 8.2 8.5 8.6 8.8 8.9 9.3
 9.4 9.6 9.7 9.7 9.8 9.8 9.9 10.1 10.2 10.2 10.3 10.5 10.5 10.8 10.9
 10.9

Number of data less than 11 = 34

Total frequency, $\Sigma f = 56$

percentage of the countries had a birth rate of less than 11 births per 1000 people =

$$\frac{34}{56} \times 100\% = 60.71\%$$

4.

a) Smallest data = 1039

Largest data = 5490

Let's consider, class width = 1000

In one possible table,

Then 5 classes will hold all the data values if started from 1001.

Table for frequency (f) and cumulative frequency (F),

Class	f	F
1001-2000	8	8
2001-3000	3	11
3001-4000	0	11
4001-5000	0	11
5001-6000	1	12
Σf		
= 12		

b)

$$mean = \frac{\sum f}{number\ of\ states} = \frac{21636}{12} = 1803$$

Data in ascending order =

1039 1086 1113 1137 1139 1166 1374 1673 2009 2110 2300 5490

Number of states, n = 12 , an even number

$$\begin{aligned}
 median &= \frac{\left(\frac{n}{2}\right)^{th} term + \left(\frac{n}{2} + 1\right)^{th} term}{2} \\
 &= \frac{\left(\frac{12}{2}\right)^{th} term + \left(\frac{12}{2} + 1\right)^{th} term}{2} \\
 &= \frac{(6)^{th} term + (7)^{th} term}{2} \\
 &= \frac{1166}{2} \\
 &= 1270
 \end{aligned}$$

mode = None

c) From (b),

mean = 1803

median = 1270

The outlier in the data set = 5490

Data in ascending order after dropping the outlier =

1039 1086 1113 1137 1139 1166 1374 1673 2009 2110 2300

mean and median without the outlier;

$$mean = \frac{\sum f}{number\ of\ states} = \frac{16146}{11} = 1467.818$$

$$\begin{aligned} median &= \left(\frac{n+1}{2}\right)^{th} term \\ &= \left(\frac{11+1}{2}\right)^{th} term \\ &= (6)^{th} term \\ &= 1166 \end{aligned}$$

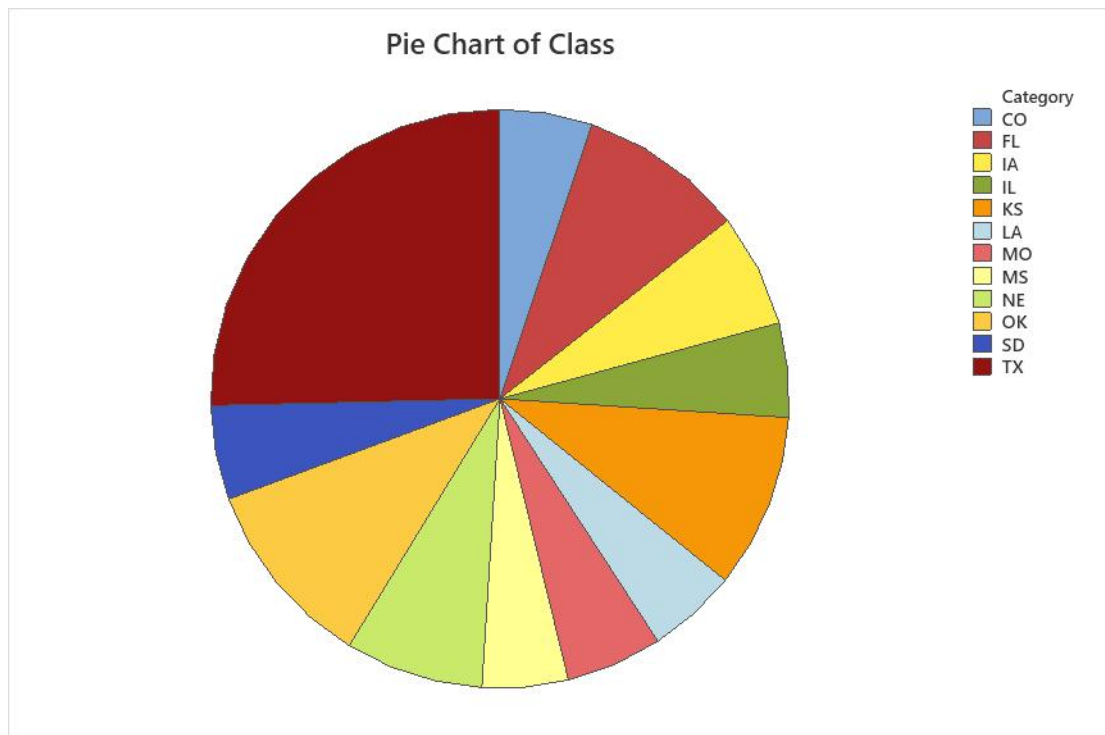
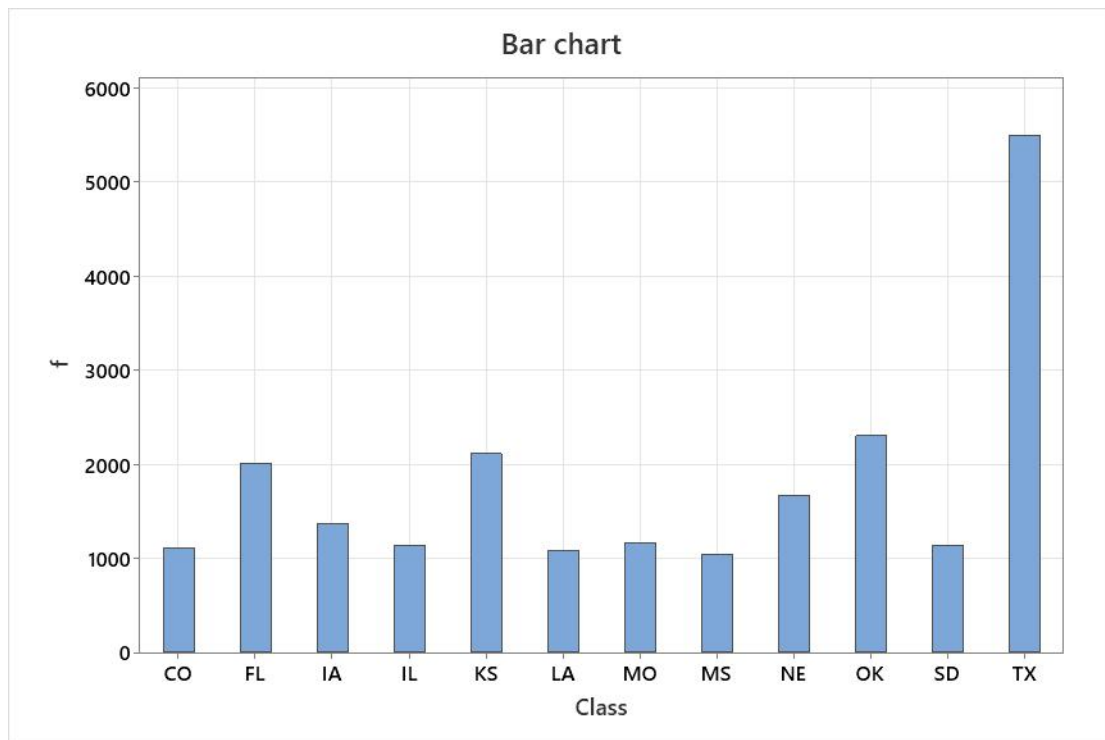
Difference in mean = 1803 - 1467.818 = 335.182

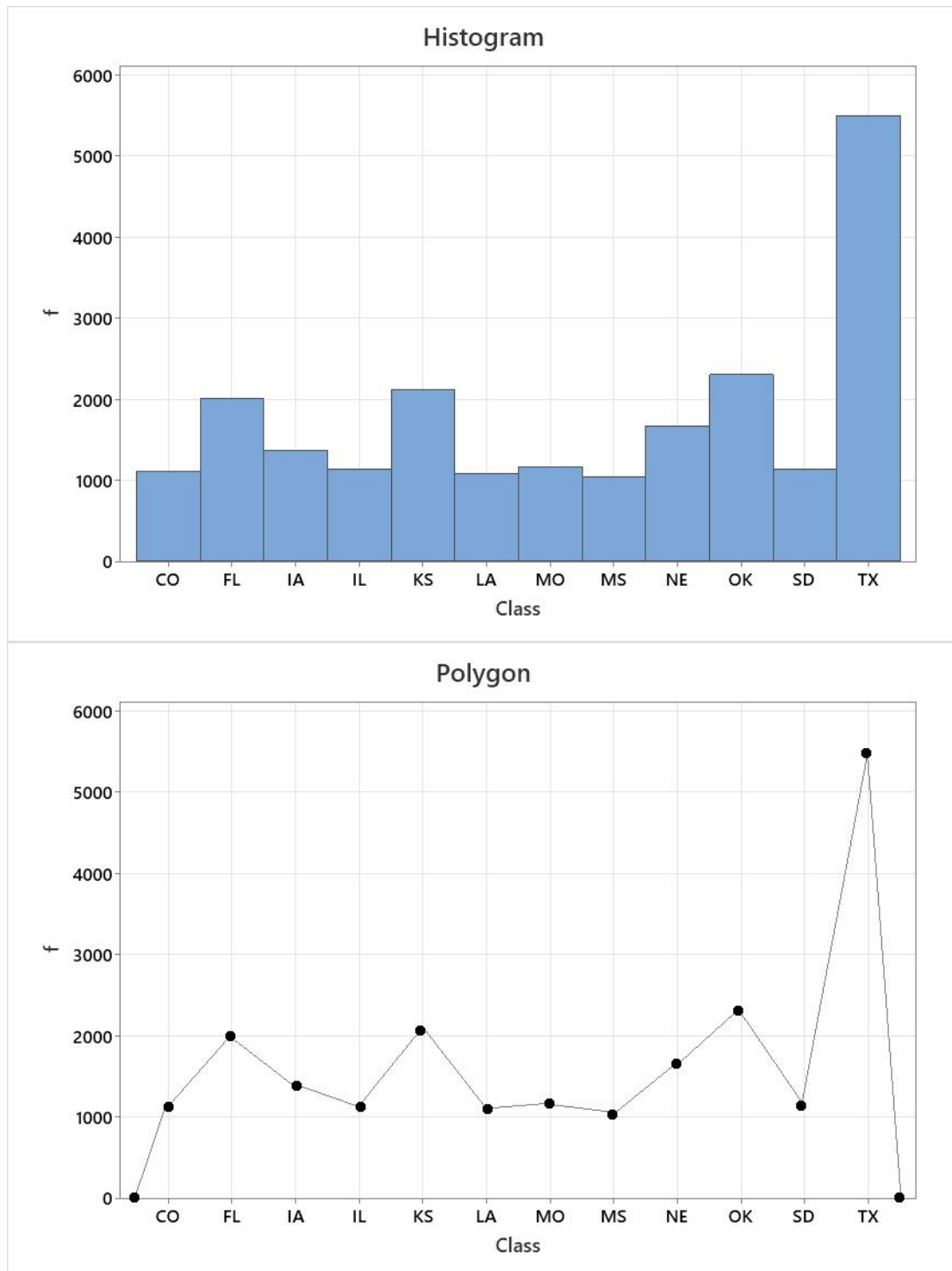
Difference in median = 1270 - 1166 = 104

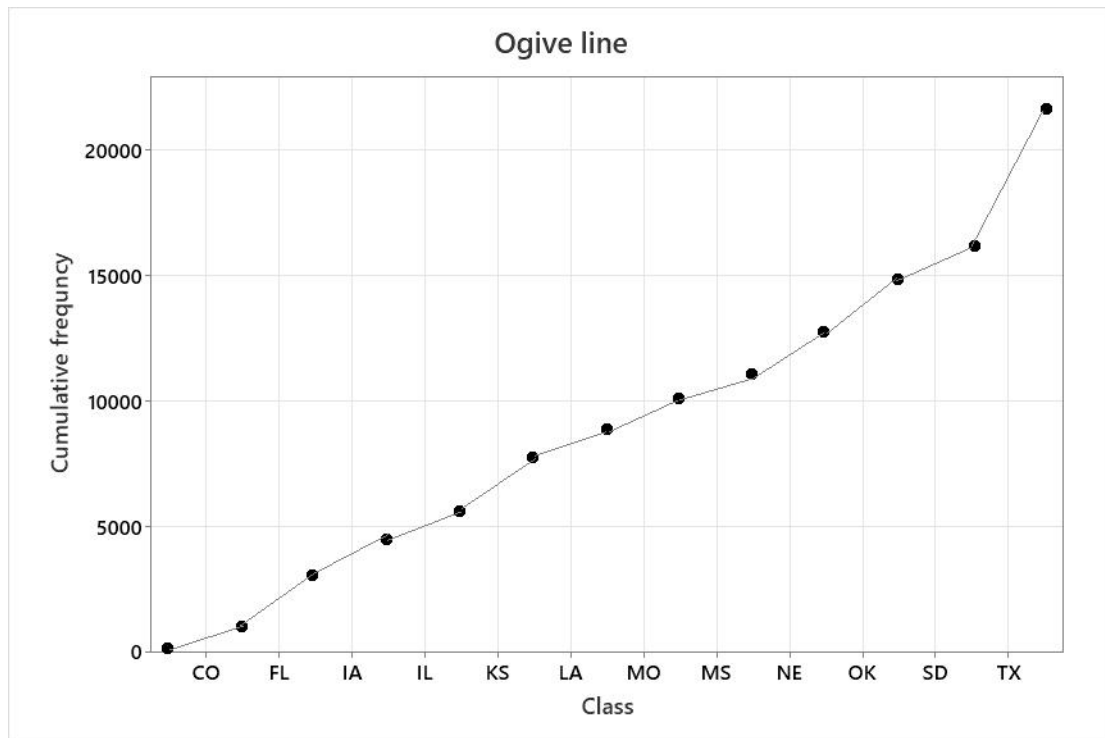
Therefore, mean is the summary measures changes by a larger amount when dropped the outlier.

d) Since there is a significant outlier in the data set, therefore median is the better summary measure for these data than the mean and the mode.

e) bar diagram, pie chart, histogram, and polygon, ogive line







5.

a)

$$\text{mean, } \bar{x} = \frac{\sum x_i}{n} = \frac{7 + 10 + 8 + 3 + 15 + 12 + 6 + 11}{8} = 9$$

Table for deviations of the data values from the mean ($x_i - \bar{x}$),

x	$x - \bar{x}$
7	-2
10	1
8	-1
3	-6
15	6
12	3
6	-3
11	2
	Σ
	= 0

Yes, the sum of these deviations is zero.

b) Smallest data value = 3

Largest data value = 15

Range = Largest - Smallest = 15 - 3 = 12

Table for calculating variance,

x	x^2
7	49
10	100
8	64
3	9
15	225
12	144
6	36
11	121
Σx	Σx^2
= 72	= 748

Number of sample, n = 8

$$\begin{aligned}
 \text{variance, } s^2 &= \frac{\sum x^2 - \frac{(\sum x)^2}{n}}{n-1} \\
 &= \frac{748 - \frac{(72)^2}{8}}{8-1} \\
 &= 14.29
 \end{aligned}$$

$$\begin{aligned}
 \text{standard deviation, } s &= \sqrt{s^2} \\
 &= \sqrt{14.29} \\
 &= 3.78
 \end{aligned}$$

Data in ascending order,

3 6 7 8 10 11 12 15

Since, $n = 8$, which is divisible by 4

$$\begin{aligned}
 \text{First quartile, } Q_1 &= \frac{(n/4)^{th} \text{ term} + ((n/4) + 1)^{th} \text{ term}}{2} \\
 &= \frac{(8/4)^{th} \text{ term} + ((8/4) + 1)^{th} \text{ term}}{2} \\
 &= \frac{(2)^{th} \text{ term} + (3)^{th} \text{ term}}{2} \\
 &= \frac{6 + 7}{2} \\
 &= 6.5
 \end{aligned}$$

$$\text{Second quartile, } Q_2 = \frac{(2n/4)^{th} \text{ term} + ((2n/4) + 1)^{th} \text{ term}}{2}$$

$$\begin{aligned}
 &= \frac{(2 \times 8/4)^{th} \text{ term} + ((2 \times 8/4) + 1)^{th} \text{ term}}{2} \\
 &= \frac{(4)^{th} \text{ term} + (5)^{th} \text{ term}}{2}
 \end{aligned}$$

$$= \frac{8 + 10}{2}$$

$$= 9$$

$$\text{Third quartile, } Q_2 = \frac{(3n/4)^{th} \text{ term} + ((3n/4) + 1)^{th} \text{ term}}{2}$$

$$= \frac{(3 \times 8/4)^{th} \text{ term} + ((3 \times 8/4) + 1)^{th} \text{ term}}{2}$$

$$= \frac{(6)^{th} \text{ term} + (7)^{th} \text{ term}}{2}$$

$$= \frac{11 + 12}{2}$$

$$= 11.5$$

Since $50 \times n = 50 \times 8 = 400$; is divisible by 100

$$50^{\text{th}} \text{ percentile, } P_{50} = \frac{(50n/100)^{th} \text{ term} + ((50n/100) + 1)^{th} \text{ term}}{2}$$

$$= \frac{(50 \times 8/100)^{th} \text{ term} + ((50 \times 8/100) + 1)^{th} \text{ term}}{2}$$

$$= \frac{(4)^{th} \text{ term} + (5)^{th} \text{ term}}{2}$$

$$= \frac{8 + 10}{2}$$

$$= 9$$

6.

a) Table for determining mean,

Class	m	f	mf
10-20	15	5	75
20-30	25	8	200
30-40	35	12	420
40-50	45	7	315
50-60	55	4	220
		Σf	Σmf
		= 36	= 1230

$$\text{mean} = \frac{\sum mf}{\sum f}$$

$$= \frac{1230}{36}$$

$$= 34.1667$$

b) Table for determining median,

Class	f	F
10-20	5	5
20-30	8	13
30-40	12	25
40-50	7	32
50-60	4	36
Σf		
= 36		

$$\begin{aligned} \text{median} &= l + \frac{(n/2) - F}{f} \times h \\ &= 30 + \frac{(36/2) - 13}{12} \times 10 \\ &= 34.1667 \end{aligned}$$

Here,
 $l = 30$
 $n = \Sigma f = 36$
 $n/2 = 36/2 = 18$
 $F = 13$
 $f = 12$
 $h = 20 - 10 = 10$

c) Table for determining mode,

Class	f
10-20	5
20-30	8
30-40	12
40-50	7
50-60	4
	Σf
	= 36

$$\begin{aligned}
 \text{mode} &= l + \frac{f_1 - f_0}{2f_1 - f_0 - f_2} \times h \\
 &= 30 + \frac{12 - 8}{2 \times 12 - 8 - 7} \times 10 \\
 &= 34.4444
 \end{aligned}$$

Here,
 $l = 30$
 $n = \Sigma f = 36$
 $n/2 = 36/2 = 18$
 $F = 13$
 $f = 12$
 $h = 20 - 10 = 10$

$$\begin{aligned}
 \text{d) First quartile, } Q_1 &= l_1 + \frac{h}{f_1} \left(\frac{n}{4} - F_{(1)-1} \right) \\
 &= 20 + \frac{10}{8} \left(\frac{36}{4} - 5 \right) \\
 &= 25
 \end{aligned}$$

$$\begin{aligned}
 \text{Second quartile, } Q_2 &= l_2 + \frac{h}{f_2} \left(\frac{2n}{4} - F_{(2)-1} \right) \\
 &= 30 + \frac{10}{12} \left(\frac{2 \times 36}{4} - 13 \right) \\
 &= 34.1667
 \end{aligned}$$

$$\begin{aligned}
 \text{Third quartile, } Q_3 &= l_3 + \frac{h}{f_3} \left(\frac{3n}{4} - F_{(3)-1} \right) \\
 &= 40 + \frac{10}{7} \left(\frac{3 \times 36}{4} - 25 \right) \\
 &= 42.8571
 \end{aligned}$$

$$\begin{aligned}
 \text{25th percentile, } P_{25} &= l_{25} + \frac{h}{f_{25}} \left(\frac{25n}{100} - F_{(25)-1} \right) \\
 &= 20 + \frac{10}{8} \left(\frac{25 \times 36}{100} - 5 \right) \\
 &= 25
 \end{aligned}$$

$$\begin{aligned}
 \text{75th percentile, } P_{75} &= l_{75} + \frac{h}{f_{75}} \left(\frac{75n}{100} - F_{(75)-1} \right) \\
 &= 40 + \frac{10}{7} \left(\frac{75 \times 36}{100} - 25 \right) \\
 &= 42.8571
 \end{aligned}$$

e) Table for mean deviation and coefficient of mean deviation, range coefficient, quartiles deviation, variance, standard deviation,

Class	x	x-mean	f	f x-mean	f(x-mean) ²
10-20	15	19.17	5	95.85	1837.44
20-30	25	9.17	8	73.36	672.71
30-40	35	0.83	12	9.96	8.27
40-50	45	10.83	7	75.81	821.02
50-60	55	20.83	4	83.32	1735.56
			Σf	Σ	Σ
			= 36	= 338.3	= 5075

Here, mean = 34.17

$$\begin{aligned} \text{mean deviation} &= \frac{\sum f|x - \text{mean}|}{\sum f} \\ &= \frac{338.3}{36} \\ &= 9.3972 \end{aligned}$$

$$\begin{aligned} \text{Coefficient of mean deviation} &= (\text{mean deviation}) / \text{mean} \\ &= 9.3972 / 34.17 \\ &= 0.28 \end{aligned}$$

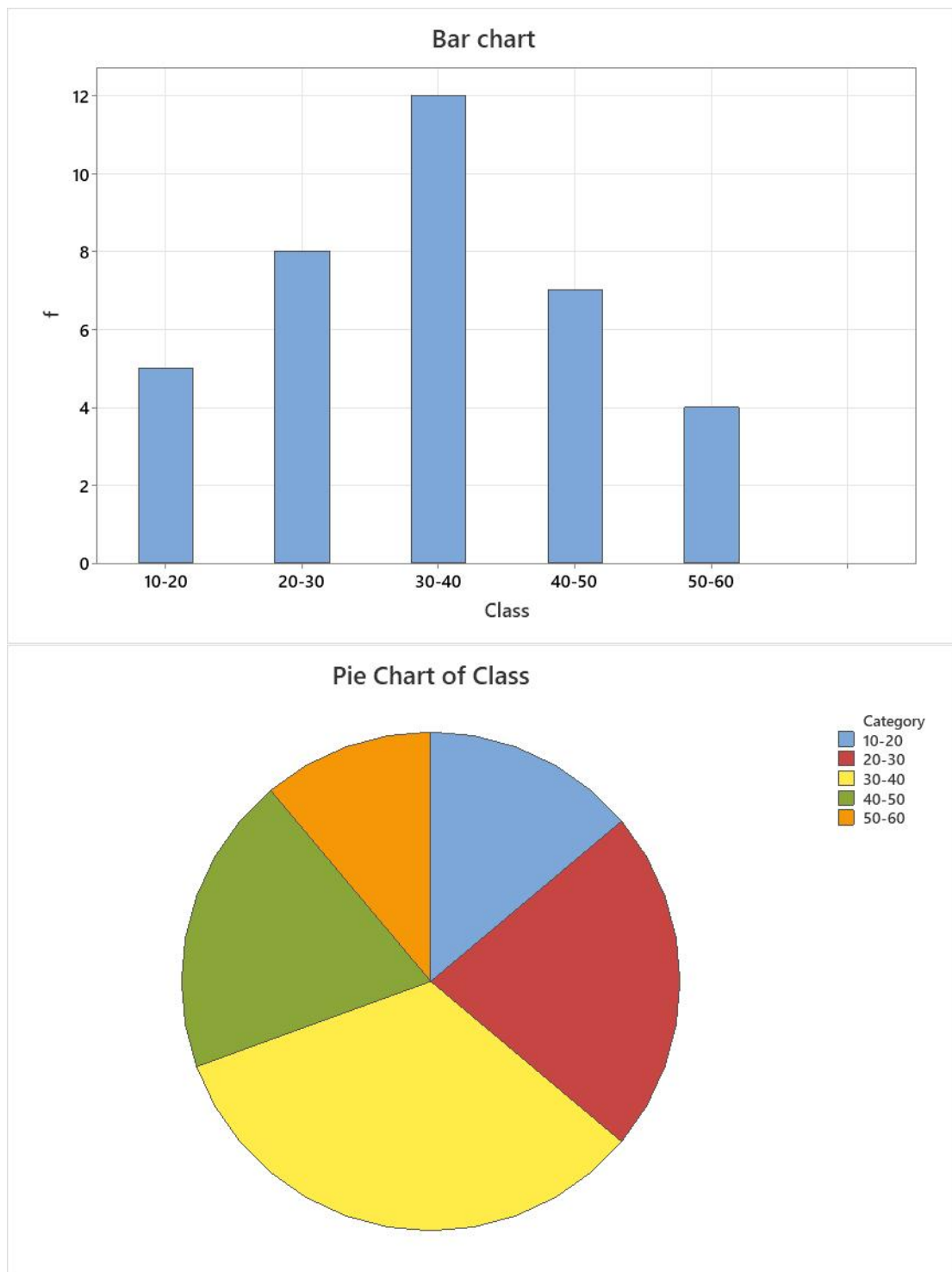
$$\begin{aligned} \text{Range coefficient} &= \frac{\text{Mid-points of the Highest Class} - \text{Mid-points of the Lowest Class}}{\text{Mid-points of the Highest Class} + \text{Mid-points of the Lowest Class}} \\ &= \frac{55 - 15}{55 + 15} \\ &= 0.57 \end{aligned}$$

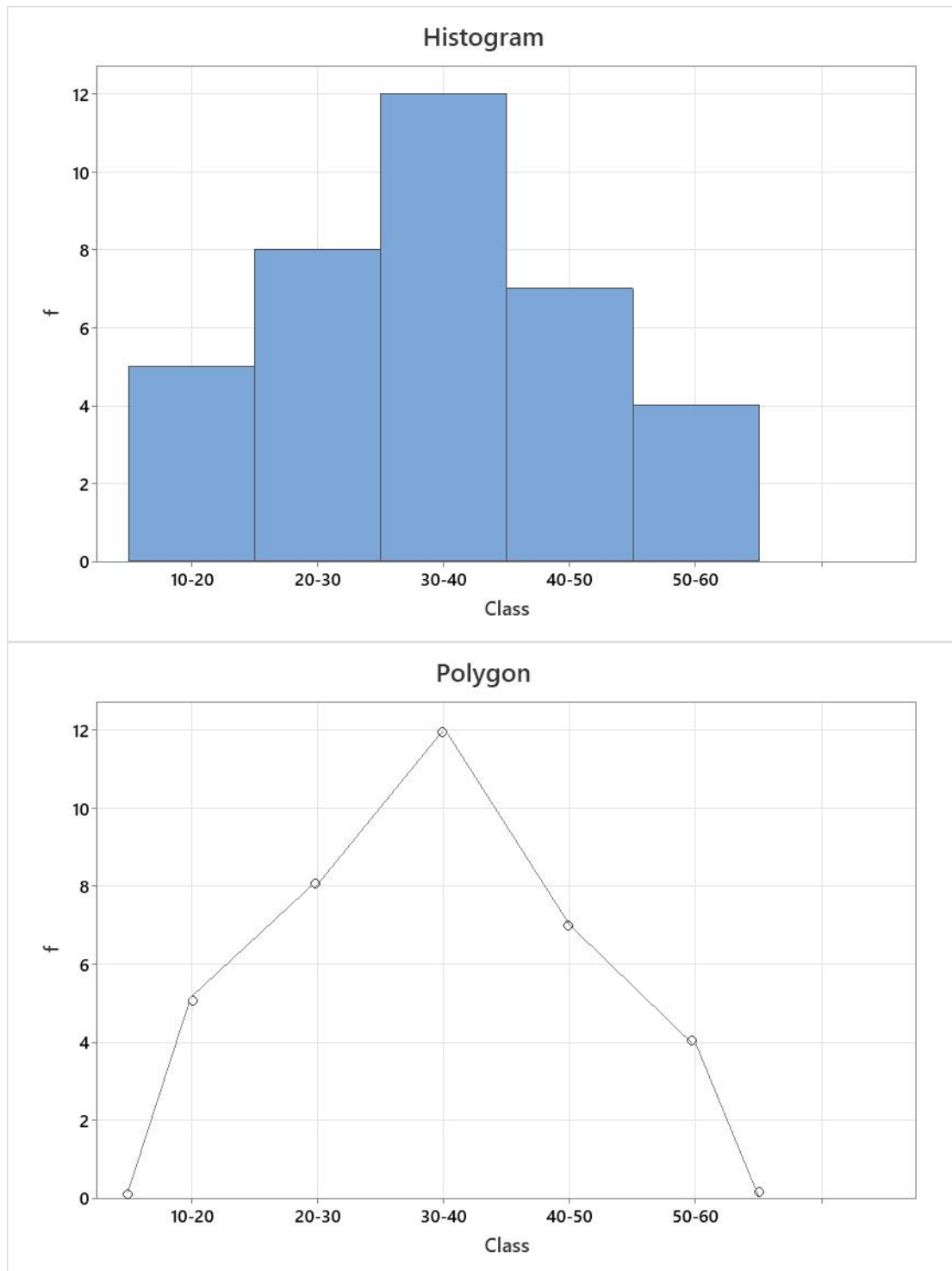
$$\begin{aligned} \text{Quartile deviation} &= \frac{Q_3 - Q_1}{Q_3 + Q_1} \\ &= \frac{42.86 - 25}{42.86 + 25} \\ &= 0.26 \end{aligned}$$

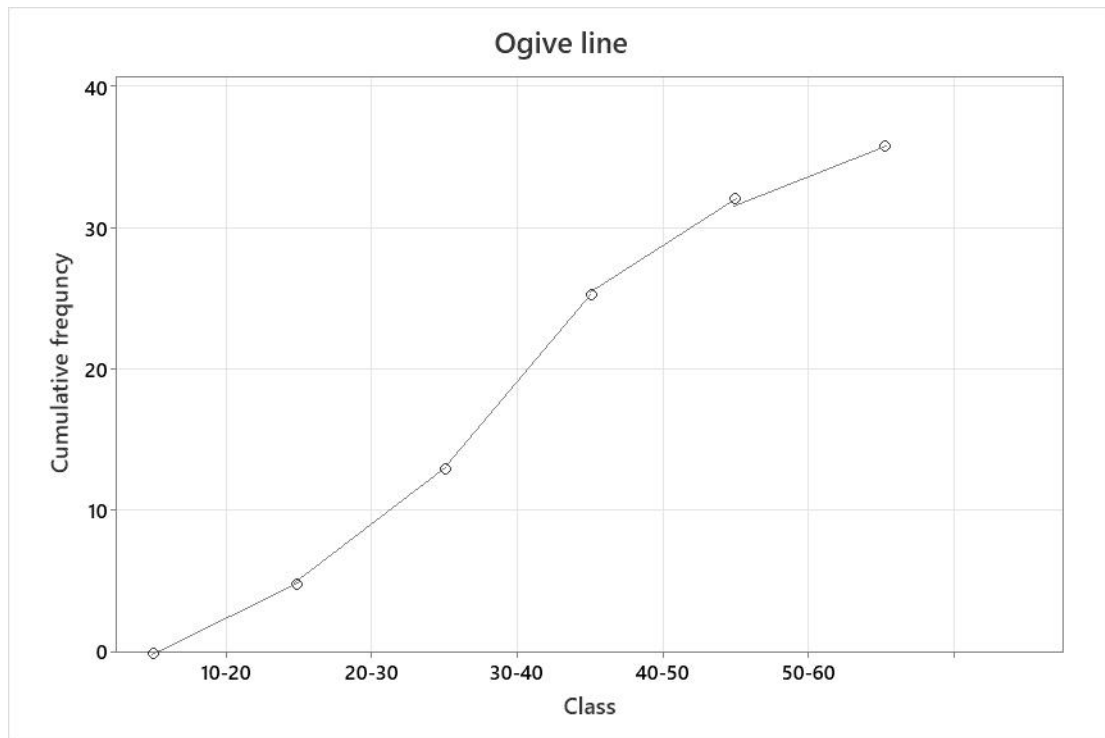
$$\begin{aligned}\text{variance} &= (\sum f(x - \text{mean})^2)/n \\ &= (5075)/36 \\ &= 140.97\end{aligned}$$

$$\begin{aligned}\text{standard deviation} &= \sqrt{\text{variance}} \\ &= \sqrt{140.97} \\ &= 11.87\end{aligned}$$

f) bar diagram, pie chart, histogram, and polygon, ogive line







7.

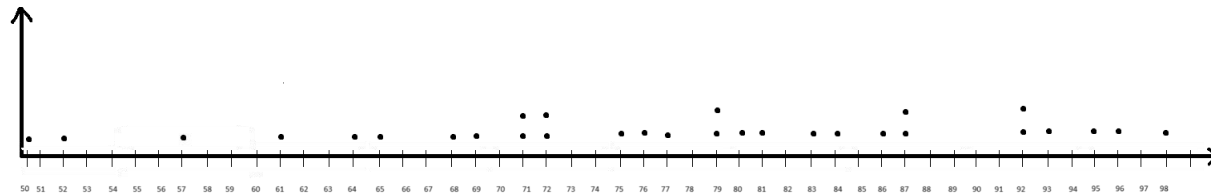
a)

Stem-and leaf display:

5	2 0 7
6	5 9 1 8 4
7	5 9 1 2 6 9 7 1 2
8	0 7 1 6 3 4 7
9	6 3 5 2 2 8

b)

Dot plot:



c)

The dot plot is more efficient for drawing conclusion about scattering, as we can easily visualize outliers, modes, minimum and maximum etc.

8.

Kurtosis is the degree of peakedness of a distribution, usually taken in relation to a normal distribution. It quantifies how much the tails of a distribution differ from those of a normal distribution.

Here's how kurtosis relates to understanding a frequency distribution:

i) High kurtosis indicates that the distribution has heavy tails and a sharp peak, suggesting that the data has more outliers compared to a normal distribution.

ii) Low kurtosis indicates that the distribution has lighter tails and is more flat-topped, suggesting that the data has fewer outliers compared to a normal distribution.

There are different measures of kurtosis, but the most commonly used is Pearson's moment coefficient of kurtosis, denoted by β_2 .

a)

Given the first four moments of distribution,

$$\mu'_1 = -1.5$$

$$\mu'_2 = 17$$

$$\mu'_3 = -30$$

$$\mu'_4 = 108$$

So,

central moments,

$$\mu_1 = 0$$

$$\begin{aligned}\mu_2 &= \mu'_2 - \mu'^2_1 = 17 - (-1.5)^2 \\ &= 14.75\end{aligned}$$

$$\begin{aligned}\mu_3 &= \mu'_3 - 3\mu'_1\mu'_2 + 2\mu'^3_1 \\ &= -30 - 3(-1.5)(17) + 2(-1.5)^3 \\ &= 39.75\end{aligned}$$

$$\begin{aligned}\mu_4 &= \mu'_4 - 4\mu'_1\mu'_3 + 6\mu'^2_1\mu'_2 + 3\mu'^4_1 \\ &= 108 - 4(-1.5)(-30) + 6(-1.5)^2(17) + 3(-1.5)^4 \\ &= 172.6875\end{aligned}$$

b)

$$\begin{aligned}
 \beta_1 &= \mu_3^2 / \mu_2^3 \\
 &= (39.75)^2 / (14.75)^3 \\
 &= 0.4924 > 0, \text{ so, positively skewed}
 \end{aligned}$$

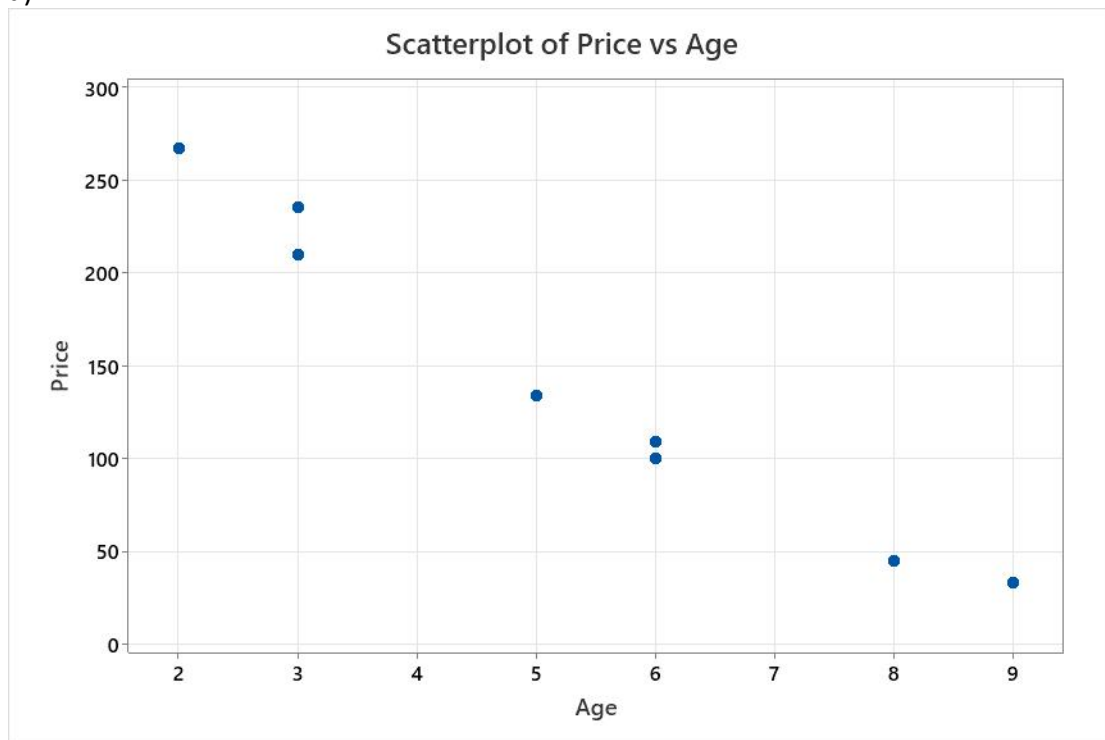
$$\begin{aligned}
 \beta_2 &= \mu_4 / \mu_2^2 \\
 &= (172.6875) / (14.75)^2 \\
 &= 0.7937 < 3, \text{ so, platykurtic}
 \end{aligned}$$

c)

While skewness tells us about the asymmetry of the distribution kurtosis provides additional information about the shape of the distribution's tails. In this case, both measures are valuable for understanding the distribution's characteristics with kurtosis indicating a distribution with heavier tails and positive skewness indicating that the mode > median > mean.

9.

a)



by looking at the scatter diagram , we can observe that there exists a strong linear relationship between car age and price. If a straight line is drawn through the points, the points will be scattered closely around the line.

b)

Age x	Price y	xy	x ²
8	45	360	64
3	210	630	9
6	100	600	36
9	33	297	81
2	267	534	4
5	134	670	25
6	109	654	36
3	235	705	9
Σx = 42	Σy = 1133	Σxy = 4450	Σx^2 = 264

So,

$$\bar{x} = \Sigma x / n = 42 / 8 = 5.25$$

$$\bar{y} = \Sigma y / n = 1133 / 8 = 141.625$$

Now,

$$\begin{aligned} ss_{xy} &= \Sigma xy - (\Sigma x)(\Sigma y) / n \\ &= 4450 - (42)(1133) / 8 \\ &= -1498.25 \end{aligned}$$

$$\begin{aligned} ss_{xx} &= \Sigma x^2 - (\Sigma x)^2 / n \\ &= 264 - (42)^2 / 8 \\ &= 43.5 \end{aligned}$$

So,

$$b = \frac{\Sigma xy - \frac{\Sigma x \Sigma y}{n}}{\Sigma x^2 - \frac{(\Sigma x)^2}{n}}$$

$$= \frac{ss_{xy}}{ss_{xx}}$$

$$= \frac{-1498.25}{43.5}$$

$$= -34.4425287356$$

$$\begin{aligned} a &= \bar{y} - b\bar{x} \\ &= 141.625 - (-34.44)(5.25) \\ &= 322.435 \end{aligned}$$

Thus the regression line estimated by employing least-square method is,

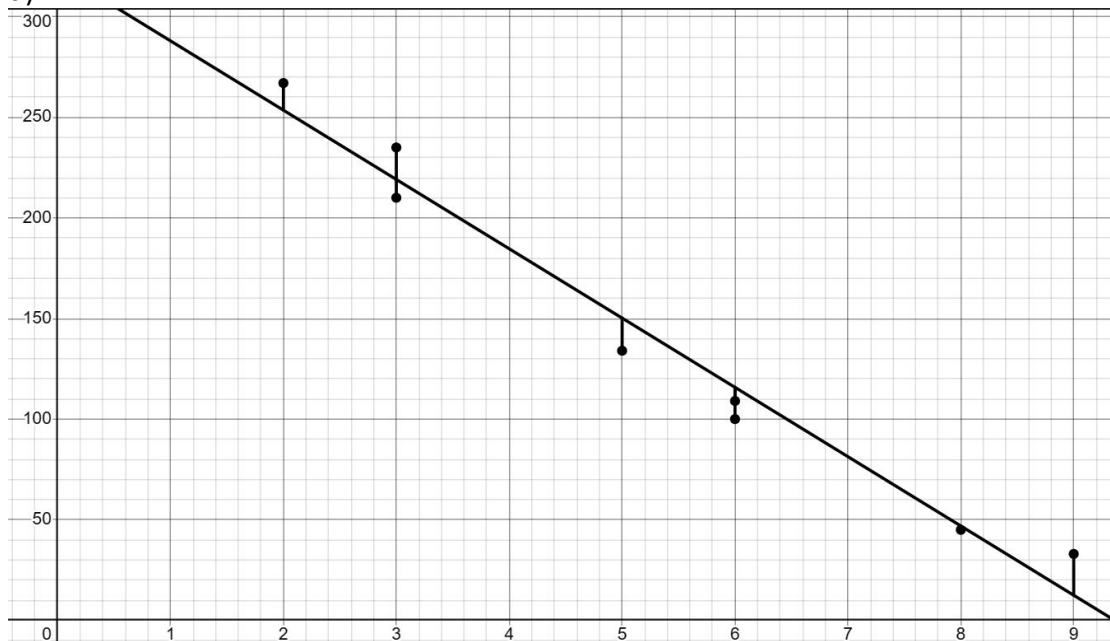
$$\hat{y} = 322.435 - 34.44x$$

c)

'b' represents an estimate of the average change in the value of the dependant variable y for each unit change in the independent variable x. In this particular case, the slope 'b' is negative which implies that as the age of cars (x) increases, the price decreases. So, the value $b = -34.44$ means that for an average increase of one year age of a car, the price would decrease on average by 34.44 hundred dollars.

The estimate of 'a' locates the regression line at point when $x=0$. Thus, if the age of a car is 0 years aka brand new, the average increase in price will be almost 322.435 hundred dollar.

d)



Graph of scatterplot and regression line and error shown with vertical lines.

Table for correlation coefficient,

Age x	Price y	xy	x^2	y^2
8	45	360	64	2025
3	210	630	9	44100
6	100	600	36	10000
9	33	297	81	1089
2	267	534	4	71289
5	134	670	25	17956
6	109	654	36	11881
3	235	705	9	55225
Σx = 42	Σy = 1133	Σxy = 4450	Σx^2 = 264	Σy^2 = 213565

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

$$= \frac{4450 - \frac{(42)(1133)}{8}}{\sqrt{264 - \frac{(42^2)}{8}} \sqrt{213565 - \frac{(1133^2)}{8}}}$$

= -0.986, value close to -1 indicates strong linear relationship with negative slope.

Since the regression model, $\hat{y} = 322.435 - 34.44x$

Regression coefficient = -34.44

Which indicates negative slope and also means that for an average increase of one unit in the independent variable, the value of dependant variable would decrease on average by 34.44 unit.

Therefore, data is not perfectly scattered.

e)

The best prediction of the price of a 7-year-old car of this model is,

$$\hat{y}_{(7)} = 322.435 - 34.44(7)$$

$$= 81.355$$

So, 81.355 hundreds of dollars or \$8135.5

f)

The best prediction of the price of a 18-year-old car of this model is,

$$\hat{y}_{(18)} = 322.435 - 34.44(18)$$

$$= -297.485$$

But, price can't be negative.

10.

a)

Regression Equation

$$\hat{y} = 15.07 + 0.1669 x_1 - 0.1323 x_2$$

Figure 1: obtained from MINITAB

b)

The regression coefficients are 15.07, 0.1669 and – 0.1323

The value of $a = 15.07$ gives the value for \hat{y} when $x_1 = 0$ and $x_2 = 0$. However, since $x_1 = 0$ and $x_2 = 0$ do not occur together in the sample data, the estimate is invalid.

The value $b_1 = 0.1669$ gives the change in \hat{y} for a one-unit change in x_1 when x_2 is held constant.

The value $b_2 = -0.1323$ gives the change in \hat{y} for a one-unit change in x_2 when x_1 is held constant.

c)

Model Summary

S	R-sq	R-sq(adj)	R-sq(pred)
1.48799	97.09%	96.37%	93.00%

Figure 2: obtained from MINITAB

standard deviation of errors, $s_e = 1.488$

the coefficient of multiple determination, $R^2 = 0.9709$

the adjusted coefficient of multiple determination, $\bar{R}^2 = 0.9637$

d)

$$\hat{y} = 15.07 + (0.1669)x_1 - (0.1323)x_2 = 15.07 + .1669(87) - .1323(54) = 22.4461$$

e)

$$\hat{y} = 15.07 + (0.1669)x_1 - (0.1323)x_2 = 15.07 + .1669(95) - .1323(49) = 24.4428$$

f)

$$df = n - k - 1 = 11 - 2 - 1 = 8$$

The 99% confidence interval for B_1 is

$$b_1 \pm t_{s_{b_1}} = .167 \pm (3.355)(.034) = .167 \pm .114 = .053 \text{ to } .281$$

g)

$$H_0: B_2 = 0, H_1: B_2 < 0$$

Since σ_{ϵ} is unknown, use the t distribution.

For $\alpha = .01$ with $df = 8$, the critical value of t is -2.896 .

$$t = (b_2 - B_2) / s_{b_2} = -1.919$$

Do not reject H_0 since $-1.919 > -2.896$.

So, B_2 is not negative.