

CIS31041 - Practical for Data Mining

Reg.No: SEU.IS.20.ICT.084

Lab Sheet 01

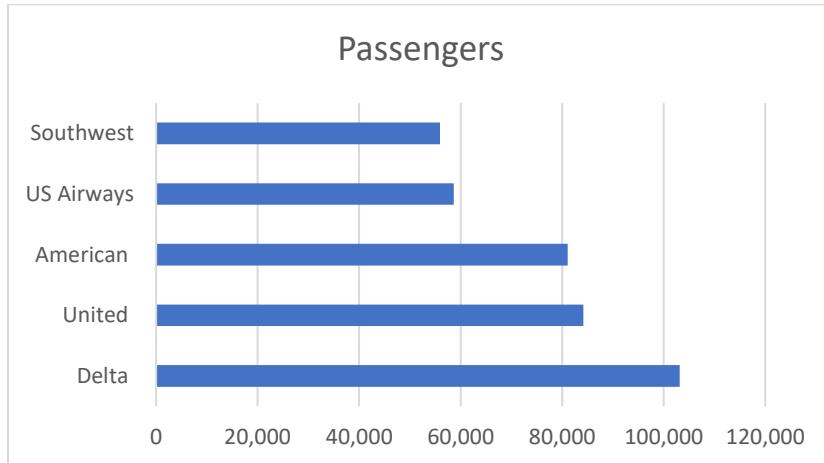
Question 01

The data set shown below the number of passengers flown by top five Airlines.

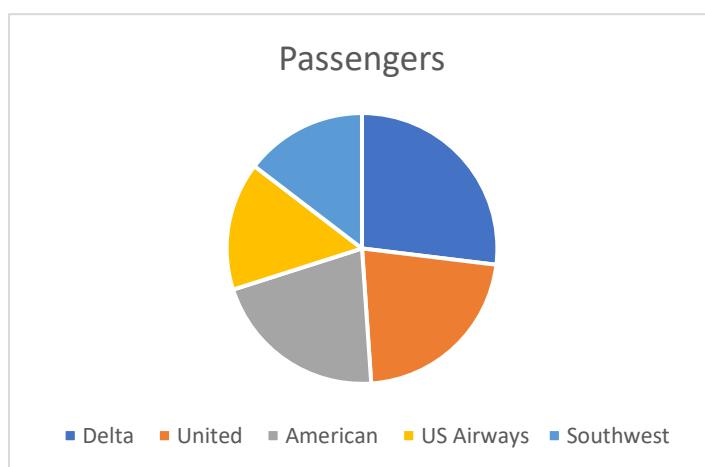
Construct Bar Chart, Line Chart and Pie Chart to depict the following information.

Airlines	Delta	United	American	US Airways	Southwest
Passengers	103,133	84,203	81,083	58,659	55,946

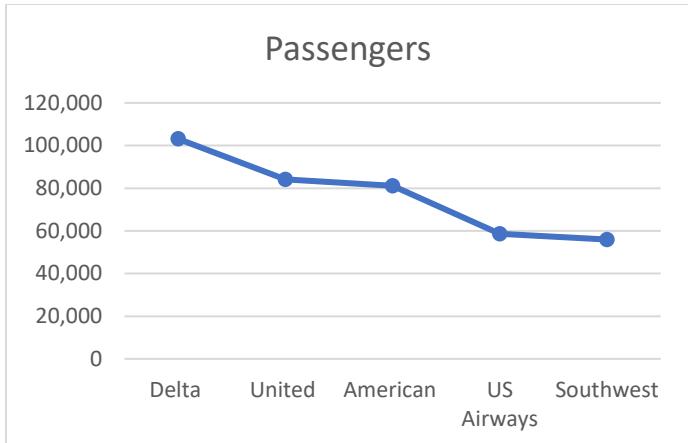
Construct Bar



Pie Chart

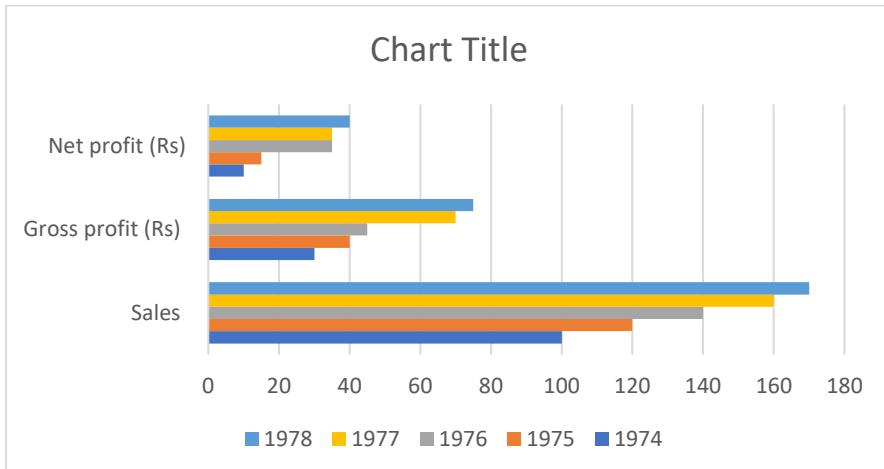


Line Chart

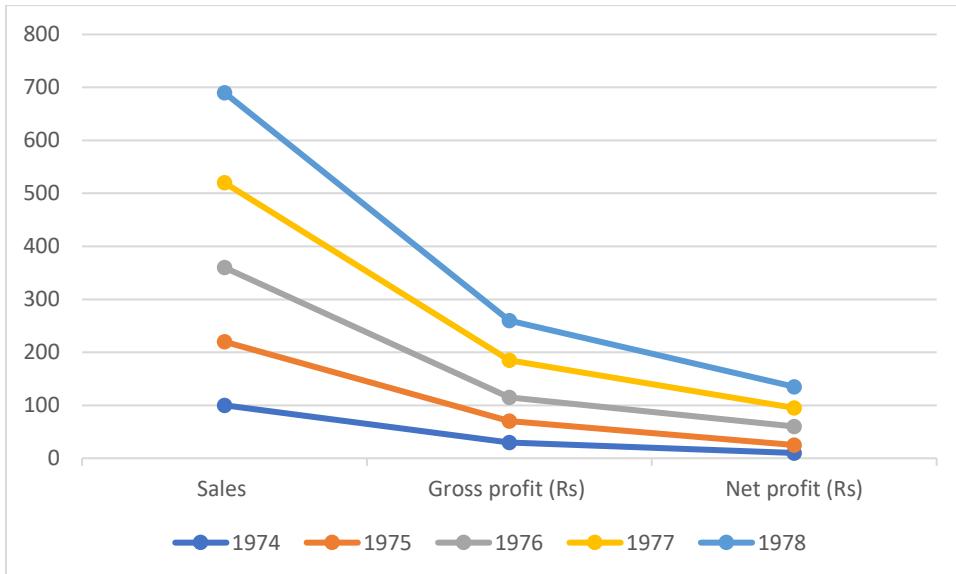


Year	Sales	Gross profit (Rs)	Net profit (Rs)
1974	100	30	10
1975	120	40	15
1976	140	45	35
1977	160	70	35
1978	170	75	40

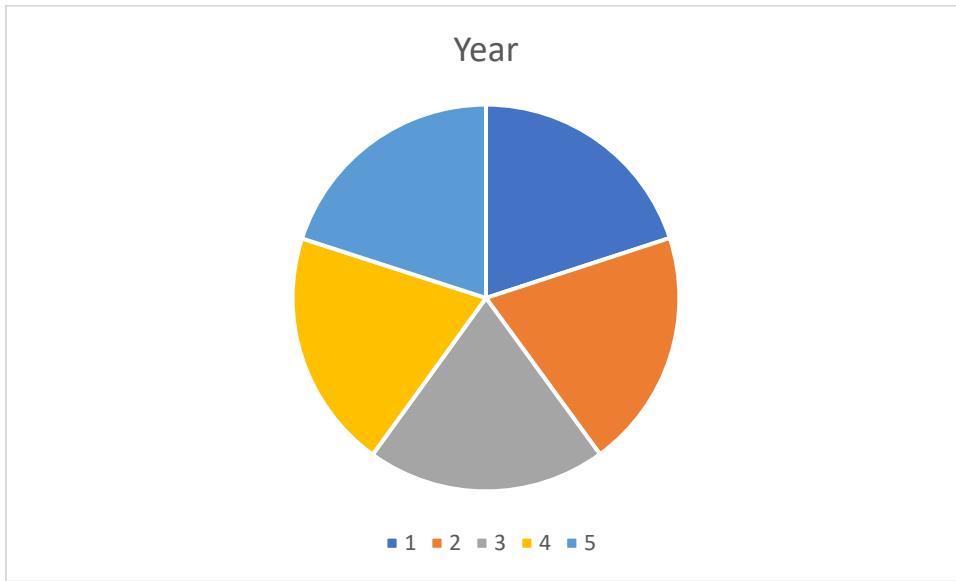
Construct Bar Chart



Line Chart



Pie Chart



Question 02

Calculate Mean, Median, Mode, Quartiles, Inter Quartile Range (IQR)

- i. 36, 48, 90, 67, 81, 41, 28, 34, 90, 77, 90

36	28			Mean	$(28+34+36+41+48+67+77+81+90+90+90)/11$	62
48	34			Median	$(n+1)/2 = (11+1)/2 = 6\text{th}$	67
90	36			Mode		90
67	41			Quartiles 1	$(n+1)/4 = (11+1)/4 = 3\text{th}$	36
81	48			Quartiles 2	$(n+1)/2 = (11+1)/2 = 6\text{th}$	67
41	67			Quartiles 3	$3(n+1)/4 = 3(11+1)/4 = 9\text{th}$	90
28	77			Inter Quartile Range (IQR)	$Q_3 - Q_1 = 90 - 36$	54
34	81					
90	90					
77	90					
90	90					

- ii. 6, 7, 2, 5, 6, 0, 3, 6, 7, 4, 5, 4, 5, 7, 7

6	0			Mean	$(0+2+3+4+4+5+5+5+6+6+6+7+7+7+7)/15$	4.933333
7	2			Median	$(n+1)/2 = (15+1)/2 = 8\text{th}$	5
2	3			Mode		7
5	4			Quartiles 1	$(n+1)/4 = (15+1)/4 = 4\text{th}$	4
6	4			Quartiles 2	$(n+1)/2 = (15+1)/2 = 8\text{th}$	5
0	5			Quartiles 3	$3(n+1)/4 = 3(15+1)/4 = 12\text{th}$	7
3	5			Inter Quartile Range (IQR)	$Q_3 - Q_1 = 7 - 4$	3
6	5					
7	6					
4	6					
5	6					
4	7					
5	7					
7	7					
7	7					

Question 03

- i. Estimate the following for the below grouped frequency distribution table
 ii. Measures the mean, Var, Std, CV for the following data.

Class	1.2 - 1.6	1.6 - 2.0	2.0 - 2.4	2.4 - 2.8	2.8 - 3.2
Frequency	220	150	90	110	280

Class	Frequency	x_i	$f \cdot x_i$	Mean Value	$\sum f_i \cdot x_i / \sum f_i$	$\sum f_i \cdot x_i^2 / \sum f_i - ((\sum f_i \cdot x_i) / \sum f_i)^2$	$\sqrt{\text{Variance}}$	$(\text{standard Deviation} / \text{Mean})$	$1902 / 850$	$[4616.4 - ((1902 / 850) * (1902 / 850))] / 850$	$\sqrt{5.425168}$	$2.237647059 / 2.329199038$	2.237647059
1.2 - 1.6	220	1.4	308										5.425168
1.6 - 2.0	150	1.8	270										2.329199038
2.0 - 2.4	90	2.2	198										0.960693793
2.4 - 2.8	110	2.6	286										
2.8 - 3.2	280	3	840										
Total	850		1902										

Mean = $\frac{\sum f_i x_i}{\sum f_i}$ variance = $\frac{\sum f_i (x_i - \bar{x})^2}{\sum f_i}$ standard deviation = $\sqrt{\sigma}$ coefficient of variation (CV) = $\left(\frac{\sigma}{\bar{x}} \right) \times 100\%$

Class	1.2 - 1.6	1.6 - 2.0	2.0 - 2.4	2.4 - 2.8	2.8 - 3.2	3.2 - 3.6	3.6 - 4
Frequency	3	7	40	50	30	15	5

Class	Frequency	Middle(x_i)	$f_i \cdot x_i$	$f_i \cdot x_i^2$	Mean Value	$\sum f_i \cdot x_i / \sum f_i$	$\sum f_i \cdot x_i^2 / \sum f_i - ((\sum f_i \cdot x_i) / \sum f_i)^2$	$\sqrt{\text{Variance}}$	$(\text{standard Deviation} / \text{Mean})$	$394.8 / 150$	$[1075 - ((394.8 / 150) * (394.8 / 150))] / 150$	$\sqrt{35.8864}$	2.6324
1.2-1.6	3	1.4	4.2	5.88									35.8864
1.6-2.0	7	1.8	12.6	22.68									5.990525853
2.0-2.4	40	2.2	88	193.6									2.275689809
2.4-2.8	50	2.6	130	338									
2.8-3.2	30	3	90	270									
3.2-3.6	15	3.4	51	173.4									
3.6-4.0	5	3.8	19	72.2									
Total	150	18.2	394.8	1075.76									

Question 04

Calculate the correlation coefficient of the following data set and find the linear regression model

- i. (1,30), (2,45), (3,51), (4,57), (5,60), (6,65), (7,70), (8,71)

	x_i	y_i	$x_i \cdot y_i$	x_i^2	y_i^2
	1	30	30	1	900
	2	45	90	4	2025
	3	51	153	9	2601
	4	57	228	16	3249
	5	60	300	25	3600
	6	65	390	36	4225
	7	70	490	49	4900
	8	71	568	64	5041
Total	36	449	2249	204	26541

n=8				
m	$[n \cdot \sum x_i \cdot y_i - \sum x_i \cdot \sum y_i] / [n \cdot \sum x_i^2 - (\sum x_i)^2]$	$[8 \cdot 2249 - 36 \cdot 449] / [8 \cdot 204 - (36)^2]$	5.440476	
c	$(\sum y_i - m \cdot \sum x_i) / n$	$(449 - 5 \cdot 440476 \cdot 36) / 8$	31.642858	
Correlation coefficient	$[n \cdot \sum x_i \cdot y_i - \sum x_i \cdot \sum y_i] / \sqrt{[n \cdot \sum x_i^2 - (\sum x_i)^2] \cdot [n \cdot \sum y_i^2 - (\sum y_i)^2]}$	$[8 \cdot 2249 - 36 \cdot 449] / \sqrt{[8 \cdot 204 - (36)^2] \cdot [8 \cdot 26541 - (449)^2]}$	0.962869	

- ii. (100,75), (120,90), (140,115), (160,140), (170,155)

	x_i	y_i	$x_i \cdot y_i$	$x_i \cdot x_i$	$y_i \cdot y_i$
	100	75	7500	10000	5625
	120	90	10800	14400	8100
	140	115	16100	19600	13225
	160	140	22400	25600	19600
	170	155	26350	28900	24025
Total	690	575	83150	98500	70575

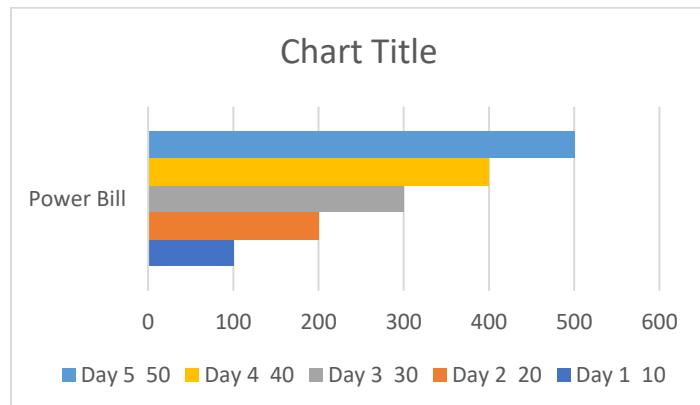
n=8			
m	$[n \cdot \sum x_i \cdot y_i - \sum x_i \cdot \sum y_i] / n \cdot \sum x_i \cdot x_i - (\sum x_i)^2$	$[8 \cdot 2249 - 36 \cdot 449] / 8 \cdot 204 - (36)^2$	5.440476
c	$(\sum y_i - m \cdot \sum x_i) / n$	$(449 - 5.440476 \cdot 36) / 8$	31.642858
Correlation coefficient	$[n \cdot \sum x_i \cdot y_i - \sum x_i \cdot \sum y_i] / \sqrt{(n \cdot \sum x_i \cdot x_i - (\sum x_i)^2) \cdot (n \cdot \sum y_i \cdot y_i - (\sum y_i)^2)}$	$[8 \cdot 2249 - 36 \cdot 449] / \sqrt{8 \cdot 204 - (36)^2} \cdot \sqrt{8 \cdot 26541 - (449)^2}$	0.962869

Exercises:

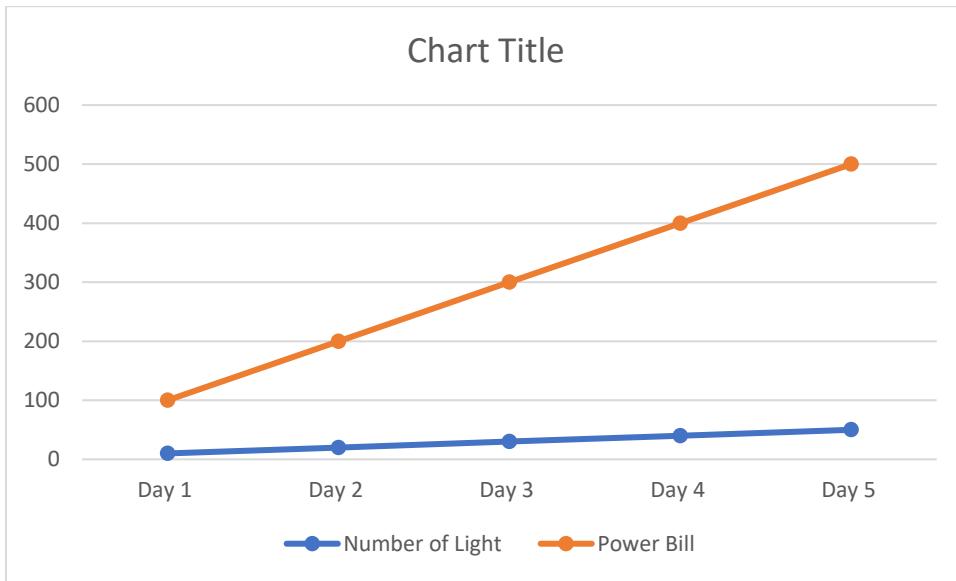
Construct Bar Chart, Line Chart and Pie Chart to depict the following information

Days	Number of Light	Power Bill
Day 1	10	100
Day 2	20	200
Day 3	30	300
Day 4	40	400
Day 5	50	500

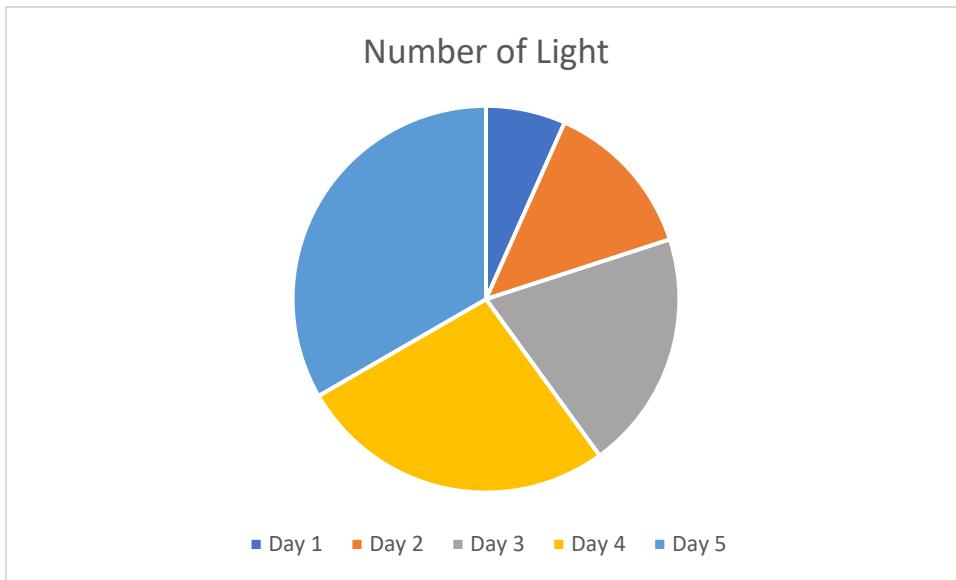
Bar Chart



Line Chart



Pie Chart



Calculate Mean, Median, Mode, Quartiles, Inter Quartile Range (IQR)

- i. 8, 15, 18, 13, 15, 28, 30, 17, 41, 27, 17, 17, 44, 31, 24

8	8	Mean	$(8+13+15+15+17+17+17+18+24+27+28+30+31+41+44)/15$	23
15	13	Median	$(n+1)/2=(15+1)/2=8\text{th}$	18
18	15	Mode		17
13	15	Quartiles 1	$(n+1)/4=(15+1)/4=4\text{th}$	15
15	17	Quartiles 2	$(n+1)/2=(15+1)/4=8\text{th}$	17
28	17	Quartiles 3	$3(n+1)/4=(15+1)/4=12\text{th}$	30
30	17	Inter Quartile Range (IQR)	$Q_3 - Q_1 = 30 - 15$	15
17	18			
41	24			
27	27			
17	28			
17	30			
44	31			
31	41			
24	44			

Calculate the correlation coefficient of the following data set and find the linear regression model

- i. (32,30), (36,30), (39,31), (40,39), (41,37), (41,34), (44,37), (45,41)

	xi	yi	xi*yi	xi*x _i	yi*y _i
	32	30	960	1024	900
	36	30	1080	1296	900
	39	31	1209	1521	961
	40	39	1560	1600	1521
	41	37	1517	1681	1369
	41	34	1394	1681	1156
	44	37	1628	1936	1369
	45	41	1845	2025	1681
Total	318	279	11193	12764	9857

n=8			
m	$[n \cdot \sum xi \cdot yi - \sum xi \cdot \sum yi] / n \cdot \sum xi^2 - (\sum xi)^2$	$(8 \cdot 11193 - 318 \cdot 279) / (8 \cdot 12764 - 318^2)$	0.831984
c	$(\sum yi - m \cdot \sum xi) / n$	$(279 - (0.831 \cdot 318)) / 8$	1.84275
Correlation coefficient	$[n \cdot \sum xi \cdot yi - \sum xi \cdot \sum yi] / \sqrt{[(n \cdot \sum xi^2 - (\sum xi)^2) \cdot (n \cdot \sum yi^2 - (\sum yi)^2)]}$	$(8 \cdot 11193 - 318 \cdot 279) / \sqrt{(8 \cdot 12764 - 318^2) \cdot (8 \cdot 9857 - 279^2)}$	0.820843

Discussion:

