

Week - 2

12/06/21

L.2.1 :- Describing Categorical Data (Frequency Distribution)

* Describing Categorical Data :- Single Variable :-

• Frequency Distribution :-

A Frequency distribution of qualitative data is a listing of the distinct values and their frequency / (count).

Each row of a frequency table lists a category along with the number of cases in this category.

Count : #

Ex :- (i). A, A, B, C, A, D, A, B, D, C

(ii). A, B, B, C, A, D, A, B, D, C, p, B, C, D, A, C, D, D.

* Construct a Frequency Distribution :-

The steps to construct a frequency distribution.

Step 1: List the distinct values of the observations in the data set in the first column of a table.

Step 2: For each observation, place a tally mark in the second column of the table in the row of the appropriate distinct value.

Step 3: Count the tallies for each distinct value and record the totals in the third column of the table.

Category	Tally mark	Frequency
A		4
B		2
C		2
D		2
Total		10

Category	Tally mark	Frequency
A		4
B		3
C		4
D		5
Total		18

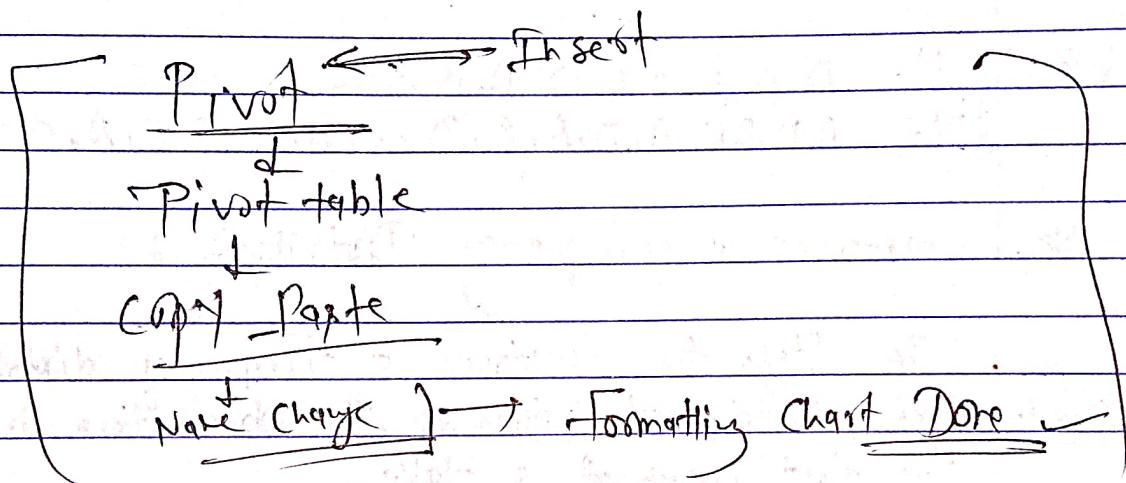
* Frequency table in a google sheet :-

Step 1 Select / Highlight the cells having data you want to visualize.

Step 2 In the formatting bar click on the Data Option.

Step 3 In the data Option go to Pivot Table option and create a new sheet. New - Insert

Step 4 After creating Pivot Table go in Pivot Table Editor and in that first add rows and then values.



* Frequency table gives the count of each variable, each Categorical Variable.

* Relative frequency :-

The Ratio of the frequency to the total Numbers of Observation is called relative frequency.

* The steps to construct a relative frequency distribution.

Step 1 Obtain a frequency distribution of the data.

Step 2. Divide each frequency by the total Numbers of Observations.

(I)				(II)			
Category	TallyMark	Fre	R.Fre	Category	T.M	Fre	R.freqe
A		3	0.4	A		5	0.4
B		2	0.2	B		3	0.2
C		2	0.2	C		3	0.2
D		2	0.2	D		3	0.2
Total		10	1	Total	15	15	1

* Important of Relative Frequency (Why?)

For Comparing two data sets.

Because relative frequencies always fall b/w 0 and 1, they provide a standard for comparison.

* In Google sheet :-

Previous	R.Fre
11	(Function)

* Summary :-

①
②

Constructing a frequency table.

Notion of relative frequency and Constructing a. Relative frequency table.

L.2, 2 - Describing Categorical Data -

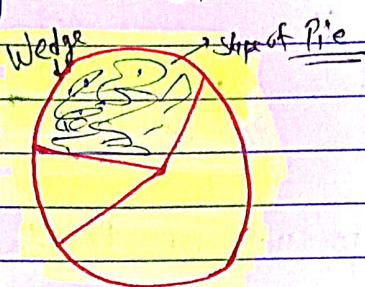
(Charts of Categorical Data)

* Charts for Categorical Data :-

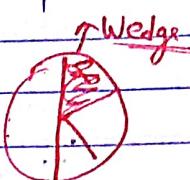
- The two most common displays of a Categorical Variable are a Barr chart and a Pie chart.
- Both describe a Categorical Variable by displaying its frequency.

* Pie charts :-

A pie chart is a circle divided into pieces proportional to the relative frequencies of the qualitative data.



WEDGES



* Construction of a Pie-Chart :-

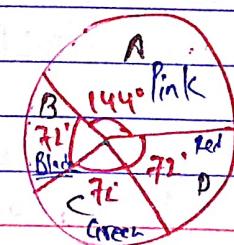
Step 1 Obtain a relative-frequency distribution of the data.

Step 2 Divide a circle into pieces proportional to the relative frequencies. ($R.Frequency \times 360^\circ$)

Step 3 Label the slices with the distinct values and their relative frequencies.

Ex 1 - A, A, B, C, A, D, A, B, D, C.

Category	Tally Mark	Fre	R. Fre	Degree
A		4	0.4	144
B		2	0.2	72
C		2	0.2	72
D		2	0.2	72
Total		10	1	360°



* Pie-chart

* Pie - Chart in a google sheet :-

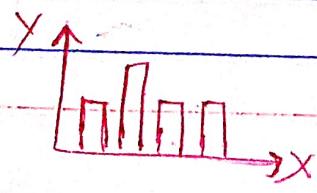
- Step 1 Select / Highlight the cells having data you want to visualize.
- Step 2 Click the Insert chart option in Google sheets toolbar.
- Step 3 Change the Visualization type in Chart editor
- Step 4 Select in Chart editor Chart type to Pie chart.

Summary :-

- ① A Pie chart is used to show the proportions of a Categorical Variable.
- ② A pie chart is a good way to show that one category makes up more than Half of the total.

* Bar Chart :-

A Bar chart Displays the distinct Values of the Qualitative Data on a horizontal axis and the relative frequencies (or frequencies or Percentages) of those Values on a vertical axis. The frequency / relative frequency of each distinct value is represented by a Vertical Bar whose height is equal to the frequency / relative frequency of that value. The bars should be positioned so that they do not touch each other.



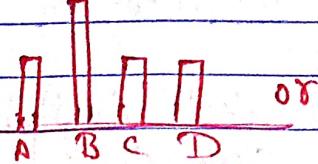
frequency — Vertical Bar
Value — horizontal axis

(q)



* Construction of a Bar chart :-

To Construct a Bar Chart.



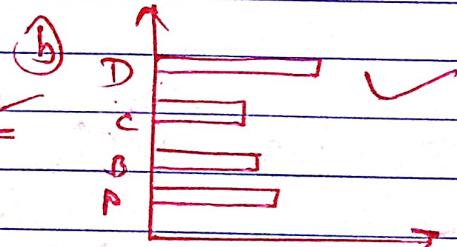
Step 1. Obtain a frequency / relative frequency distribution of the Data.

Step 2. Draw a Horizontal axis on which to place the bars and a Vertical axis on which to display the frequencies / relative frequencies.

Step 3 For each distinct Value, Construct a Vertical Bars Whose height equal the frequency / relative frequency of that Value.

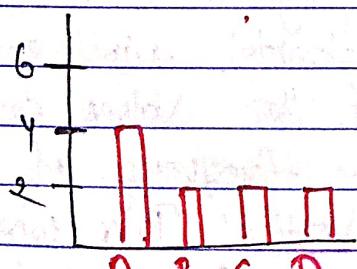
Step 4 Label the bars with the distinct Value, the Horizontal axis with the name of the variable, and the vertical axis with "frequency / relative frequency".

Both (a) and (b) Types are possible.



from

Ex(1)



Copy

* for Ordinal → Maintain Order in Pareto chart
* for Nominal → Name, Value.

* Bar chart in Google sheet :-

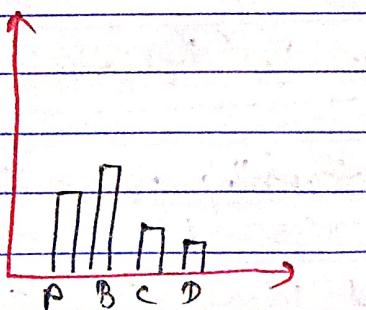
- Step1 Select / highlight the cells having Data you want → visualize.
- Step2 Click the insert Chart option in Google Sheets toolbar.
- Step3 Change the visualization type in Chart editor.
- Step4 Select in Chart editor Chart type to Bar chart.

* Pareto Charts :-

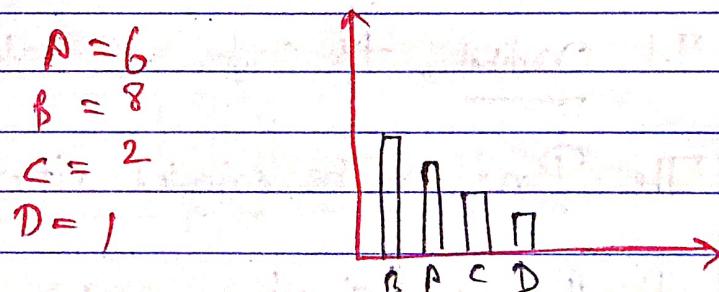
When the Categories in a bar chart are sorted by frequency, the bar chart is sometimes called a Pareto Chart.

Pareto charts are popular in quality control to identify problems in a business process. → (Nominal → Name, Value) | (Ordinal → Order)

- If the Categorical Variable is ordinal, then the bar chart must preserve the ordering.

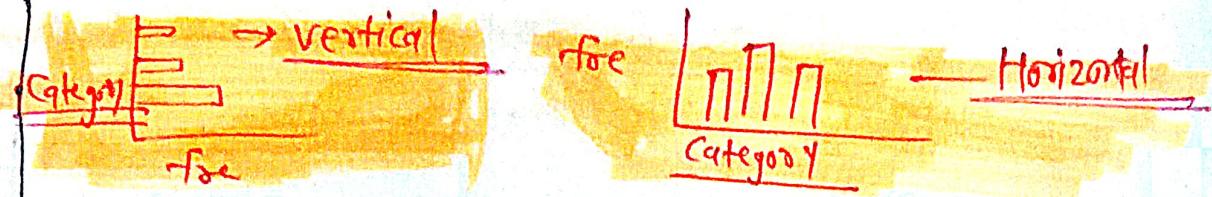


Bar chart



Pareto Chart

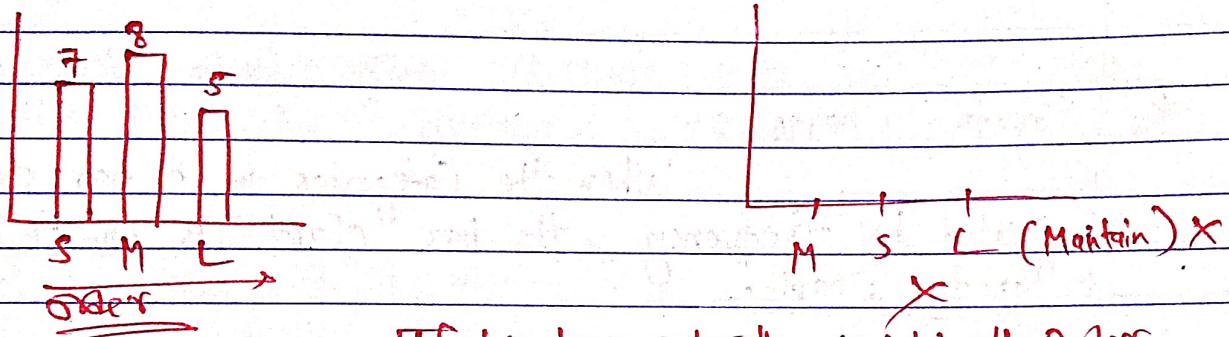
* Ascending Order or
Descending Order.



Ex. Ordinal Data

Size of T-Shirt - L, M, M, S, L, S, S, M, L, M, M, S, S, L, M, S, M, S, L, M.

Size	Tally Marks	freq	relative freq
Small		7	0.35
Medium		8	0.40
Large		5	0.25
Total		20	1



If you have order then maintain the Order.

* Summary

- ① A Bar chart is used to show the frequencies/relative frequencies of a Categorical Variable.
- ② If Ordinal, the Order of Categories is Preserved.
- ③ The Bars can be oriented either horizontally or vertically.
- ④ A Pareto chart is a bar chart where the Categories are sorted by frequency.

L.2.3. Describing Categorical Data (Part)

Practices while graphing Data-1

Note:

Know your purpose,

Have a purpose for every table or graph you create.
Choose the table/graph to serve the purpose.

- * Pie charts are best to use when you are trying to compare parts of a whole.
- * Bar graphs are used to compare things b/w different groups.

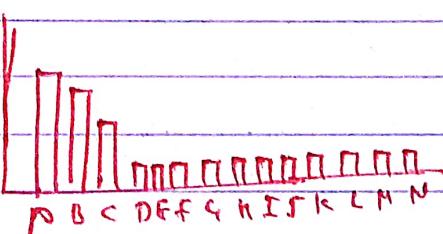
- * Label Your Data : - (Tutorial from sheet)

- Label your chart to show the categories and indicate whether some have been combined or omitted.
- Name the bars in a bar chart.
- Name the slices in a pie chart.
- If you have omitted some of the cases, make sure the label of the plot defines the collection that is summarized.

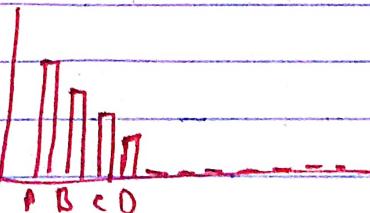
- * Many Categories.

A Bar chart or pie chart with too many categories might conceal the more important categories. In some case, grouping other categories together might be done.

(1)



(2)



L.2.4. D. C. Data

C Best Practices while Graphing Data-2

13/06/2024

* The area Principle :-

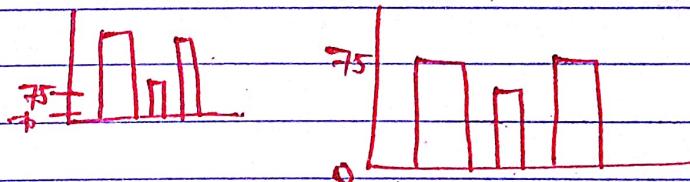
- Display of Data must Obey a fundamental rule called the area principle.
- The area principle says that the area occupied by a part of the graph should correspond to the amount of data it represent.
- Violations of the area principle are a common way to mislead with statistics.

* Misleading Graphs : ① Violating Area Principle.

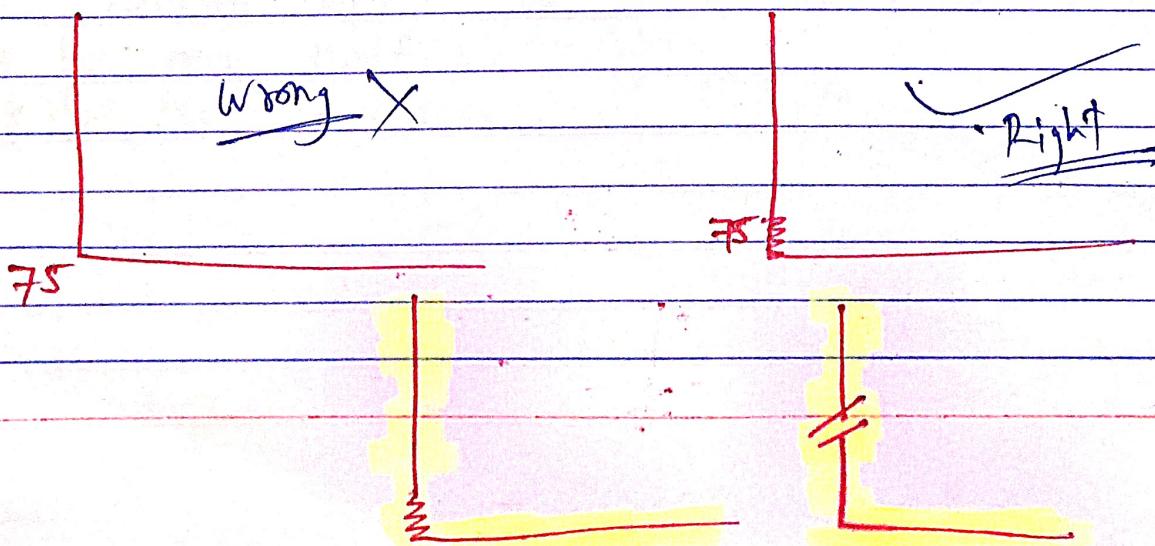
- Decorated Graphics : Charts Decorated to Attract attention often Violate the area area Principle.

② Truncated Graph : → Because of loss of Data

- Another common violation is when the Baseline of a bars chart is not at zero.



③ Indicating a y-axis break : -



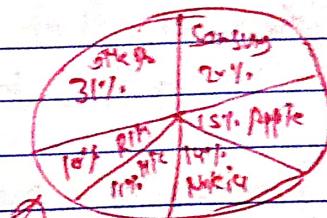
1.9.5 - Describing Categorical Data :

(Mode and Median)

* Round-off errors :-

- Important to check for round-off errors.

		2.5e	
\rightarrow	A	22.5	23
	B	35.5	36
	C	12.5	13
	D	11.20	11
	E	17.50	19
II		100	102

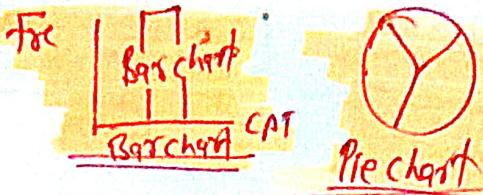


- When table entries are percentages or proportions, the total may sum to a value slightly different from 100% or 1. This might result in a pie chart where the total does not add up.

* Sectional Summary :-

- ① Know your purpose and choose table/graph appropriately
- ② Label your charts
- ③ Handle multiple categories appropriately →
- ④ Respect area principle.
 - 4.1 Avoid Overly decorated Graphs.
 - 4.2 Avoid truncated graphs - Use special symbols to indicate vertical axis has been modified.
- ⑤ Check for round-off errors.

only counting ✓
→ Nominal Ordinal ← order



* Summarizing Categorical Data :-

Table | For / R. For

- Graphical Summaries of Categorical Data : Bar chart and Pie chart.
- Need for a Compact measure.
- Numbers that are used to describe data sets are called Descriptive measures.
- Descriptive measures that indicate where the centre or most typical value of a data set lies are called measures of central tendency.

1. Mode :-

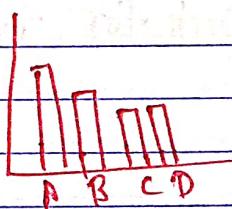
The Mode of a Categorical Variable is the most Common category, the category with the highest frequency.

- The Mode labels :-
- The longest bar in a bar chart.
- The widest slice in a pie chart.
- In a Pareto chart, the mode is the first category shown.

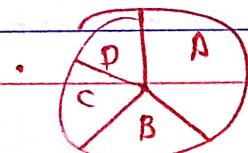
Ex ①. A, B, C, D, A, A, B, A, C, D, A

①

Mode = A (highest frequency).



②



- * We can't define median for a Nominal Categorical Data.
- * So for median we have ordinal Data we have median.
- * Median is Value that divides data set in 2 halves.

* Bimodal and Multimodal Data :-

- If Two or more categories tie for the highest frequency, the data are said to be bimodal (in case of two) or multimodal (more than two).

* Bimodal ————— 2 Modes

* Multimodal ————— More than 2 Modes.

* Median :-

Ordinal data offers another summary, the median, that is not available unless the data can be put into order.

* Definition :-

The median of an ordinal variable is the category of the middle observation of the sorted values.

If there are an even number of observations, choose the category on either side of the middle of the sorted list as the median.

If it odd then exactly that value.

Ex ①. P, P, B, A, C, D, P, B, B, D.

[Even]

$$\text{Median} = \underline{\underline{B}}$$

P, P, B, C, D, D.

either B or C. (median) (3rd Observation or 4th Observation)

③ P, A, B, C, D

[Odd]

Median = B — 3rd Observation

- Mode — Nominal + Ordinal
- Median — only for Ordinal Data

* Sectional Summary :-

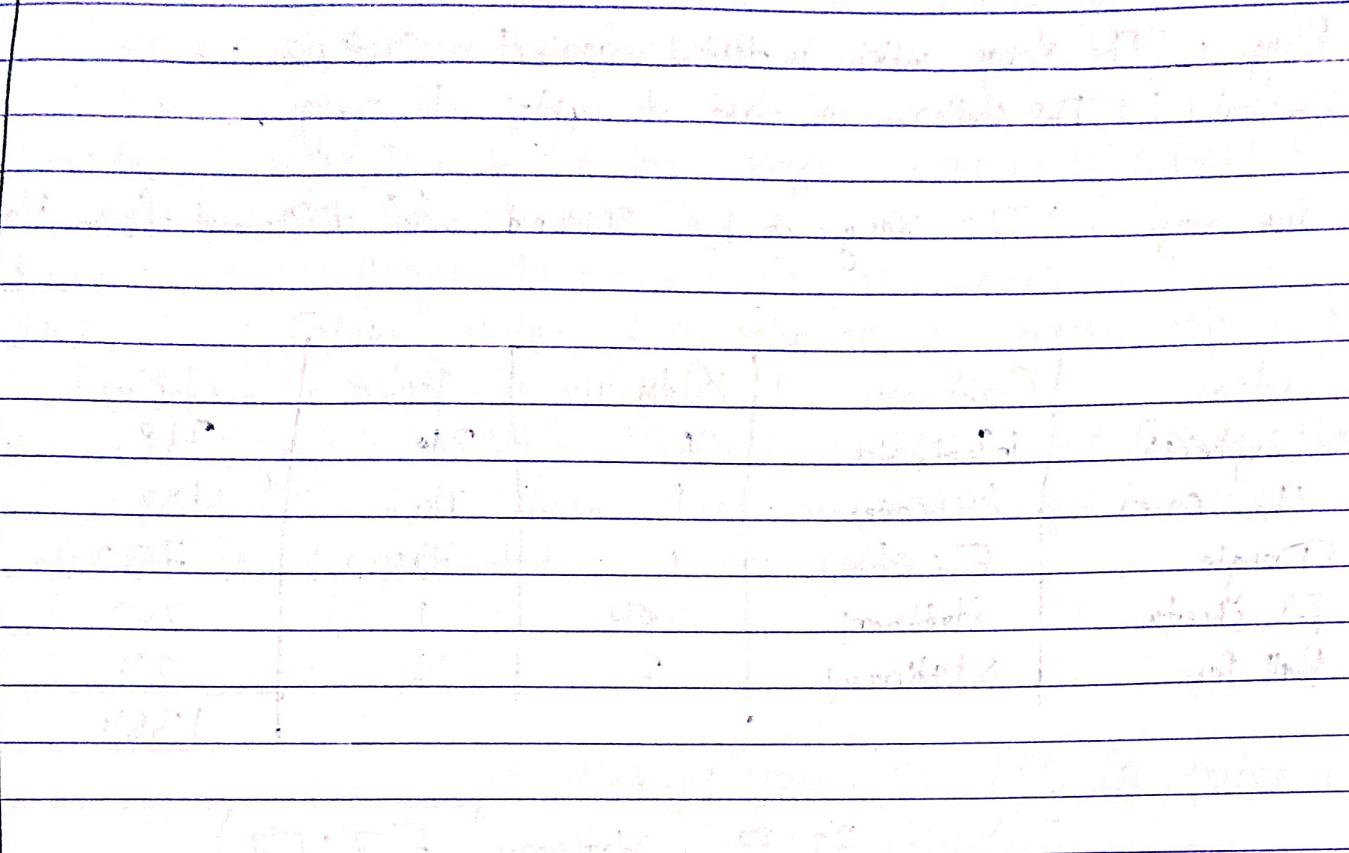
* (Nominal + ^{ordinal})

- The mode of a Categorical Variable is the most Common Category.
- The median of an ordinal variable is the Category of the middle observation of the sorted Value. (~~Nominal + ordinal~~)

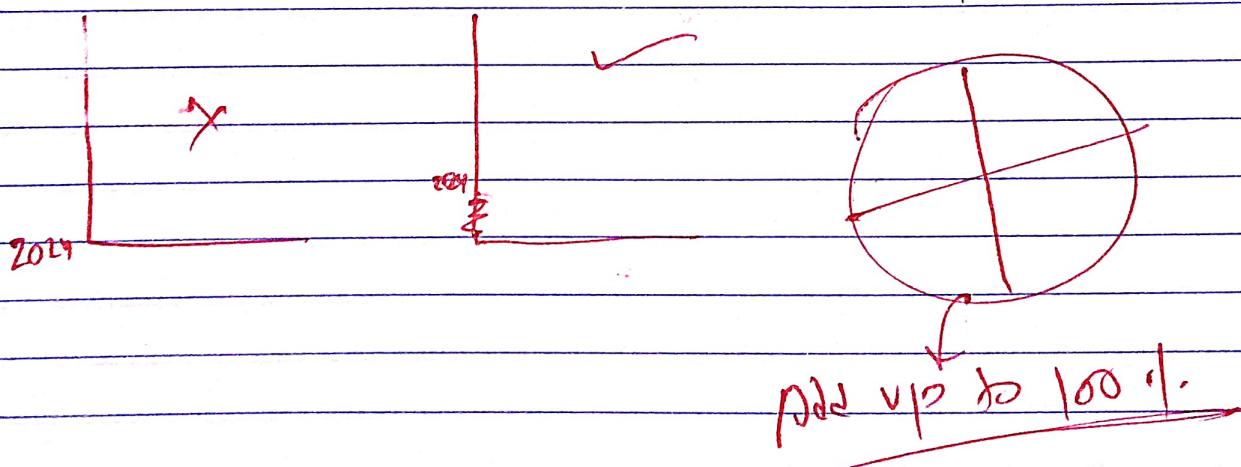
* Categorical Data ####

Week - 2 Tutorial - 1 (Problem chart & table) .

Q. Solving :-



Tutorial - 2 - Problems Misleading Graph :-



Tutorial - 3 (SUMIF in Google Sheets)

* Syntax: `(SUMIF (range, criterion, [sum_range]))`

- Range: The Range which is tested against criterion.
- Criterion: The pattern or test to apply to range.
- Sum_range - The range to be summed , if different from Range.

Item	Category	Qty	Price	Cost
Earphones	Electronics	1	210	210
Phone Cover	Accessories	1	740	740
Dongle	Electronics	1	790	790
PY sheets	Stationary	200	1	200
Ball Pens	Stationary	2	12	24
				1364

= `Sumif(B3:B7, "Stationary", E3:E7)`

Tutorial - 7 - (VLOOKUP in Google sheet)

- Vlookup :- Vertical lookup . Searches down the First Column of a range for a key and returns the Value of a specified cell in the row found.
- Syntax . Vlookup (search_key , range , index , [is_sorted])
- Search key :- The Value to search for .
- Range :- The Range to consider for the search .
- index : The Column index of the value to be returned , where the first column in range is numbered 1 .
- is_sorted : (True by default) Indicates whether the column to be searched (the first column of the specified range) is sorted . False is recommended in most cases .

= Vlookup("place", A1 : A{9}, 2, False)