

Categorical Data

Numerical Data

Summary Table for one variable

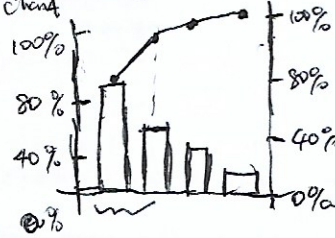
Contingency Table for two variables

Scatter plot

Frequency Distribution and cumulative Distn



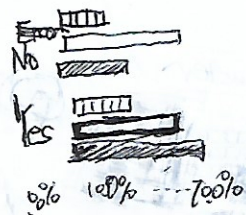
Pareto Chart



The Viteri few

Separate the "vital few" from the "trivial many"

Side by side Chart

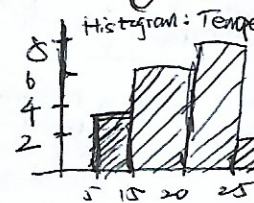


Large Medium Small

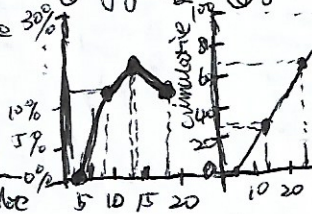
Donut chart

ordered array
stem and leaf display

Histogram



Polygon



The frequency polygon

Cumulative polygon

Central tendency

Mean Arithmetic

Affected by extreme values

Median : Median position: $\frac{n+1}{2}$ → even the average of two middle numbers
less sensitive than the mean to extreme values

Mode : Not affected

Weighted Mean : $\bar{X}_w = \frac{\sum w_i x_i}{\sum w_i}$: gives formal equal weight to all observations

Case	Grade	Grade
I	3	A
II	3	B
III	5	B

$$\frac{(3 \times 4) + (3 \times 3) + (5 \times 3)}{3 + 3 + 5} = \frac{12 + 9 + 15}{11} = \frac{36}{11} \approx 3.27$$

Geometric Mean : measure the rate of change of a variable over time

$$\bar{X}_G = (X_1 \times X_2 \times \dots \times X_n)^{1/n}$$

Geometric mean rate of return:

$$\bar{R}_G = [(1+R_1) \times (1+R_2) \times \dots \times (1+R_n)]^{1/n} - 1$$

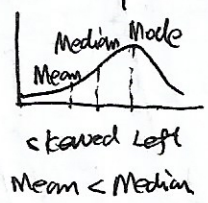
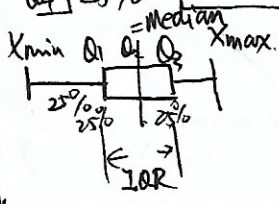
investment from = \$100,000 → 15% decrease → 30000 → 100% increase → \$100,000

$$\bar{R}_G = [(1 + (-0.5)) \times (1 + 1)]^{1/2} - 1 = 0\%$$

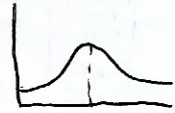
$$\bar{X} = \frac{(-0.5) + 1}{2} = 0.25 = 25\%$$

Quantiles

Quantile	Percentage	Formula
Q1	25%	$Q_1 = \frac{(n+1)}{4}$
Q2	50%	$Q_2 = \frac{(n+1)}{2}$
Q3	75%	$Q_3 = \frac{3(n+1)}{4}$
Q4	100%	



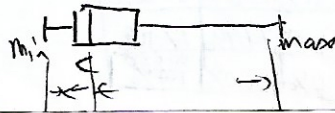
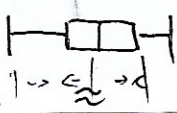
Mean < Median



Mean = Median = Mode



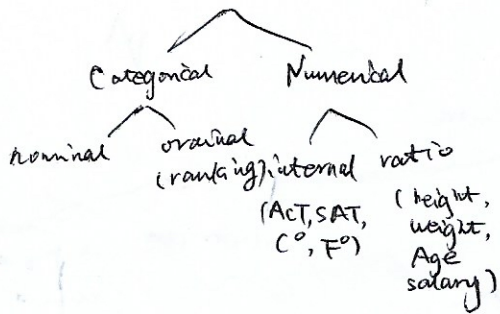
Mean > Median



mark the outlier with a star



DOVA: Define, Collect, Organize, Visualization, Analyse.



when the researcher cannot get a complete list of the members of a population.
 ② a list of subjects so utterly scattered as to be unusable.

vs. random samples:
 ① More cost effective, saving time
 ② less efficient
 ③ units close to each other may be very similar

ensure the representation of individuals across the entire population
 ① Better coverage
 ② convenient

Difficulty in identifying appropriate strata.
 survey errors:
 ① coverage error or selection bias: excluded from frame
 ② Nonresponse error or bias: can be leveraged by survey designers
 ③ Sampling error: random differences from sample to sample
 ④ Measurement error: bad or leading problem

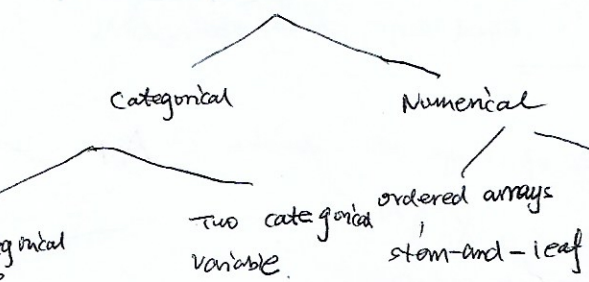
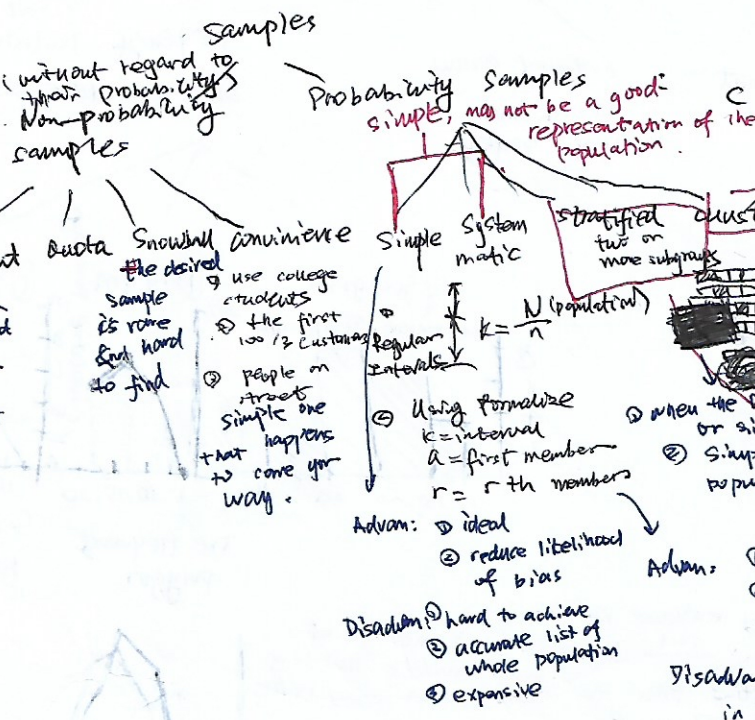
Frequency distributions & cumulative Distributions

Relative & Percent Frequency Distribution

Class	Midpoints	Frequency
$10 < X < 20$	15	3
$20 \leq X < 30$	25	6
$30 \leq X < 40$	35	5
$40 \leq X < 50$	45	4
$50 \leq X < 60$	55	2
Total		20

Relative & Percentage Frequency Distribution	Cumulative Frequency Distribution
---	--------------------------------------

Class	No	Yes	Total
Small	30.71%	30.71%	47.50%
Medium	29.81%	61.30%	35.00%
Large	19.40%	7.69%	17.50%
Total			100.00%



	Percentage
1	34%
2	29%
3	35%

Summary Table

Contingency Table

	No	Yes	Total
Small	170	20	190
Medium	100	40	140
Large	65	5	70
Total	335	65	400

	No	Yes	Total
Small	42.5%	5.0%	47.50%
Medium	25.0%	10.0%	35.00%
Large	16.4%	1.25%	17.50%
Total	83.71%	16.29%	100.00%

Percentage of Row Totals

	No	Yes	Total
Small	89.47%	10.53%	100.0%
Medium	71.43%	28.57%	100.0%
Large	92.86%	7.14%	100.0%

Percentage of Column Totals

	No	Yes	Total
Small	30.71%	30.71%	47.50%
Medium	29.81%	61.30%	35.00%
Large	19.40%	7.69%	17.50%
Total			100.00%