

Q.01 Dataset to draw - tabularly as shown.

CustomerID	LoanType	Gender	MaritalStatus	AccommodationType	YearsInAddress	YearsInJob	Salary	Balance	Target
1	Home / Home Improvement	M/F	Family / Single Others	Family / Other / Company / Own / Rent	0-100 n, numeric	0-100 numeric	€R Any numbers	€R Realno. numeric	??
2	Binary / Categorical with 2 levels	Binary / 2-level Category	Categorical 3-level	S-level					
3	Unique			Categorical					

Summary: 1) CustomerID \rightarrow primary key / unique \rightarrow CustomerID
 (Ignored for modelling)

2) LoanType \rightarrow Loan (0/1) \rightarrow 0 for Home, 1 for Improvement

3) Gender = 0/1, 1 = male, 0 = female

4) Marital Status = Family, Single, Other = Status

5) AccommodationType = Acc = Family, Own / ... etc.

6) YearsInAddress = numeric = YIA

7) YearsInJob = YIJ = numeric

8) Salary numeric & 9) Balance numeric

10) Charts for cat. / bin / num. attr.

Categorical \rightarrow Salary vs. YIA or YIJ or YIA distribution

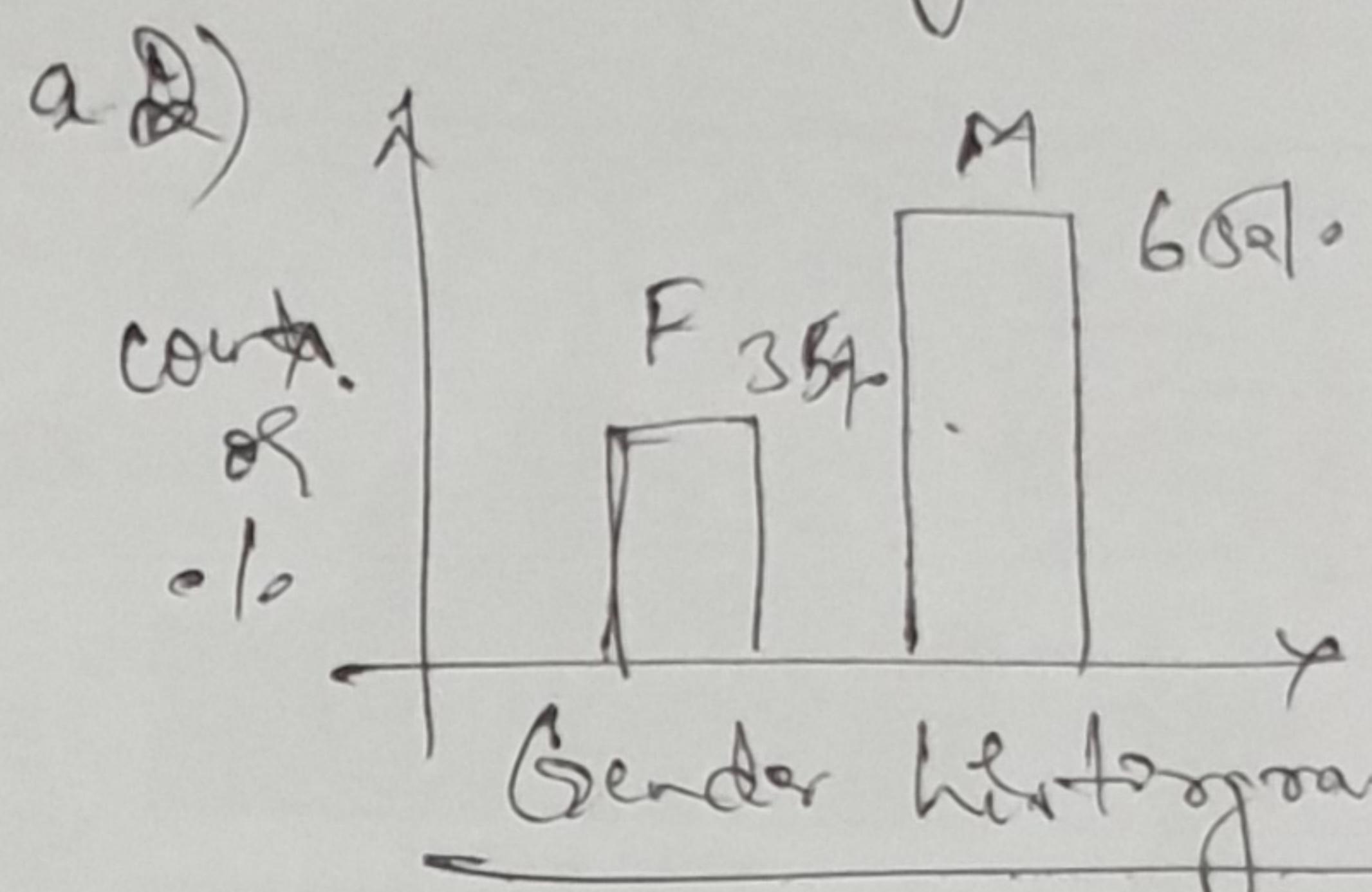
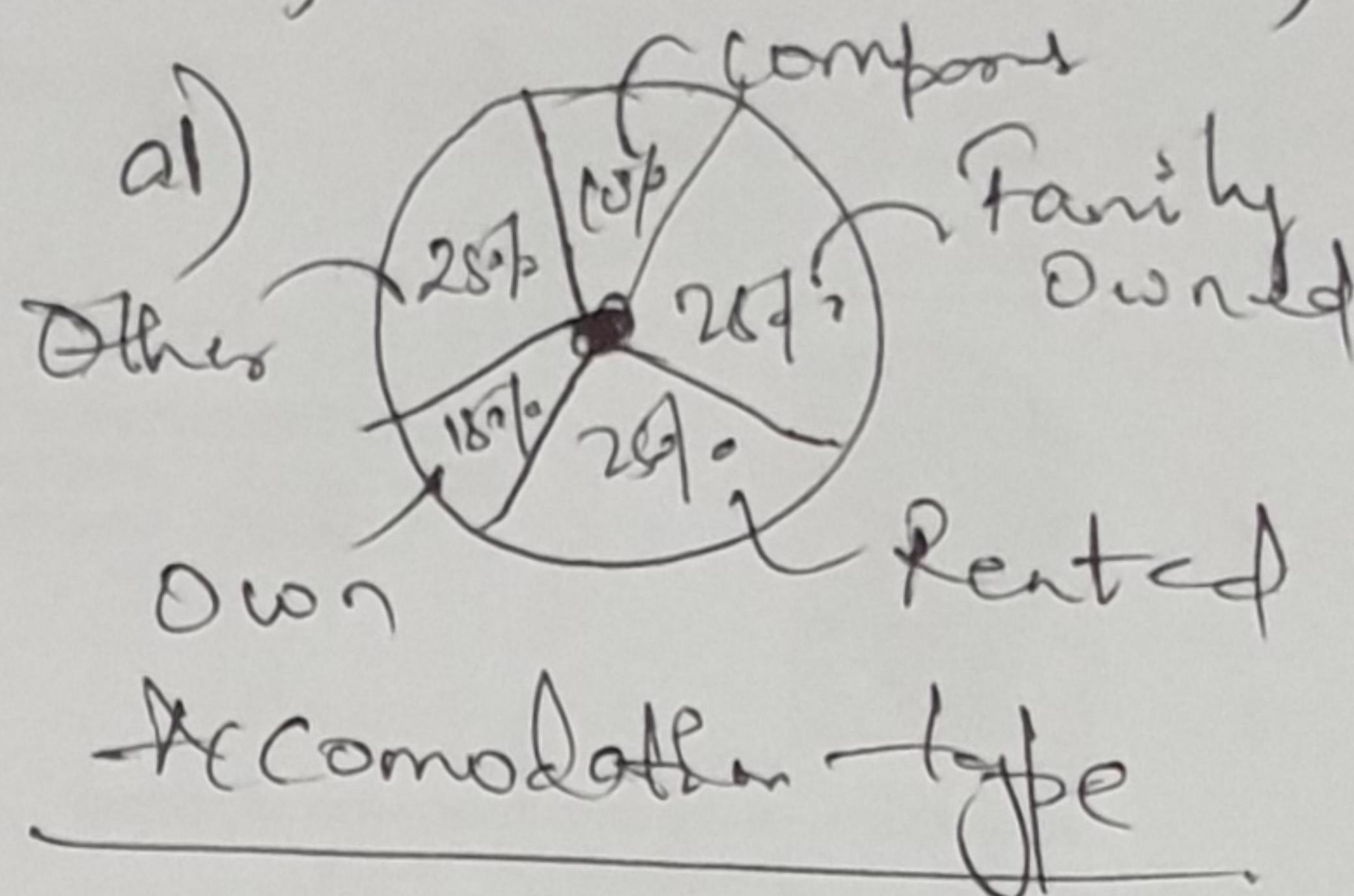
Binary \rightarrow Gender vs. AccommodationType, or Gender histogram

Numeric \rightarrow Salary ranges of salary vs. balance of gender

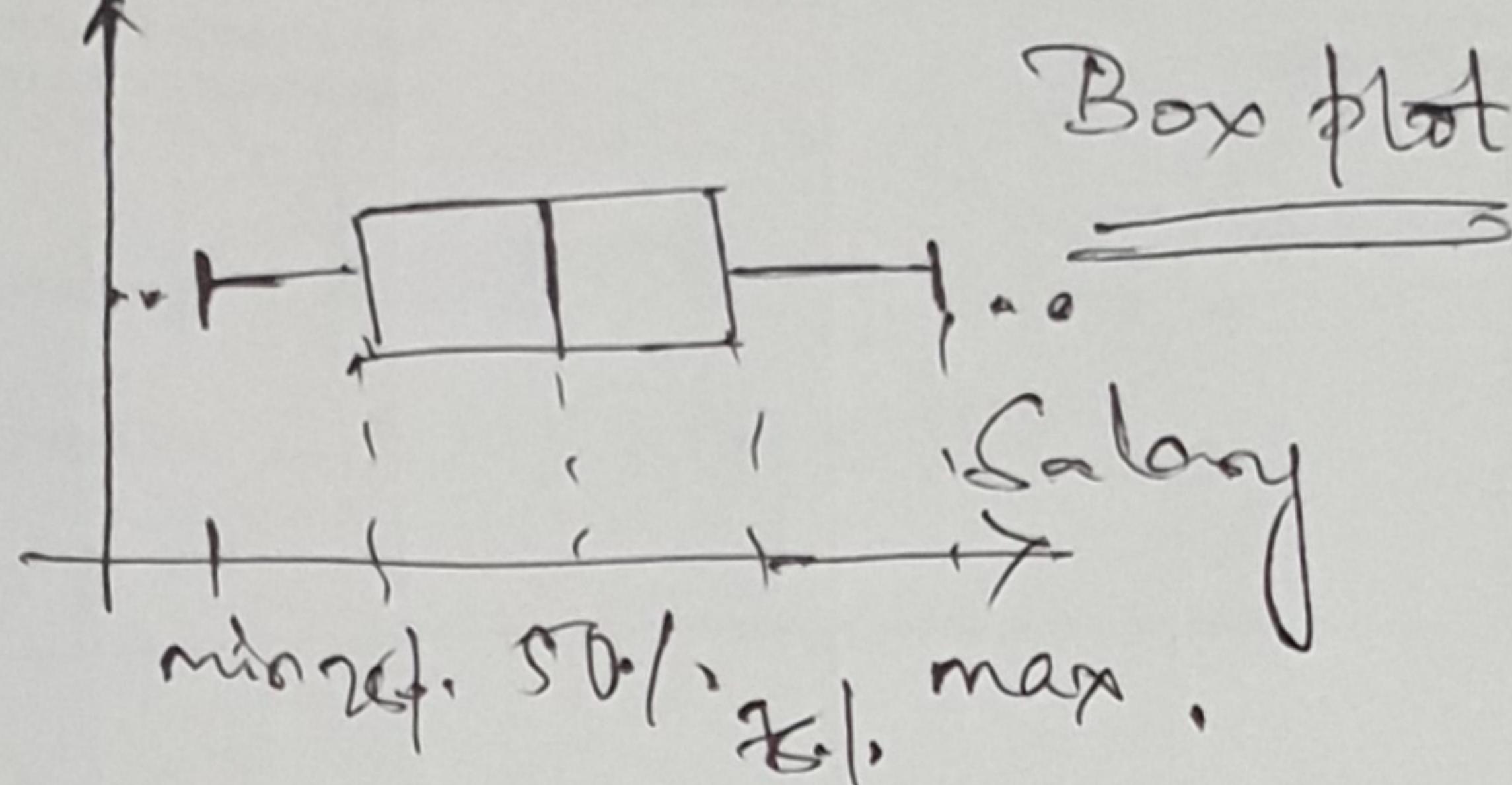
swathanam, 2019-AIRL618, PCAM 7211 Predicting. EC-3R CR

Q1 Ques. a)

PIE



a3)



NOTE: There are many other

ways to slice & dice data
as added on the previous page.

- a1) Check % distribution of diff. accommodation types / any category
- a2) Check freq. distribution of gender / any binary var.
- a3) Distribution plot showing median, 25%, range & outliers (box plot) for salary / or any numeric attribute.

b) Multiple Regression : [Loan sanction amount] (numeric)

Logistic Regression : [Loan Default] or [Loan Approval] (categorical / binary).

Decision Tree: This can predict test values as well as classify categories, hence all of the above target types listed can be used, e.g.,

[Loan sanction amount or Loan Default or Loan Approval]

c) Multiple regression features : ~~ID → dropped.~~
Acc. & Status have to be encoded

or binarized

Ex. Status = F/S/Other.

Loan Type	Gender	Y1	Y2	Salary	Balance
0	0	numeric	numeric	numeric	numeric
1	1	0	0	0	0

Family	0	0	1
Single	0	1	0
Other	1	0	0

Catt. c) So the total columns will be.

After one indicated there are 3 columns for status of 5 for acc.

Net (columns) = Loantype, Gender, YIT, YIA, Salary, Balance, Status_F, Status_O,
Status_A, acc-f, acc-oth, acc-co, acc-o, acc-r

Target = Loan Amount sanctioned

Status_F = 001 = family

Status_O = 010 = single

Status_A = 100 = Other

acc-f = 00001 = family owned

acc-oth = 00010 = other

acc-co = 00100 = company provided

acc-o = 01000 = own, acc-r = 10000 created

d) Numeric attr are YIA, YIT, Salary & Balance.

Issues than can arise.

→ YIA, YIT & Salary cannot be < 0 (invalid/incorrect data)
→ need correction

→ salary can go up to a large number (need normalization)

→ salary can be skewed (detect outliers)

→ Normalization can be done with Z-score $\left[\frac{x_i - \bar{x}}{\sigma_i} \right]$ or IQR.

x_i = value, \bar{x} = mean, σ_i = std. dev.

→ From box-plot, IQR can be values $< 1.5 \times 25\%$
or $> 1.5 \times 75\%$

E-Score : I CANNOT SEE THIS QUESTION.

PROCTOR HAS GIVEN THE BELOW TABLE WITH 5 ENTRIES

F1	F2	CLASS
1	5	C ₁
2	6	C ₁
3	7	C ₂
4	1	C ₂
5	2	C ₂

$$F = \frac{\sum_{i=1}^k p_i (M_i - \mu)^2}{\sum_{i=1}^k p_i \sigma_i^2}$$

$$\text{Mean } M_1 \text{ for class } C_1 = \frac{1+5+2+6}{4} = \frac{14}{4} = 3.5$$

$$\text{Mean } M_2 \text{ for class } C_2 = \frac{3+7+4+1+5+2}{6} = \frac{22}{6} = 3.67$$

$$\text{Global mean } \mu = \frac{1+3+4+5+5+6+7+1+2}{10} = \frac{36}{10} = 3.6$$

C ₁	1, 5, 2, 6
C ₂	3, 7, 4, 1, 5, 2

$$, \quad p_1 = 4/10, \quad p_2 = 6/10$$

$$\Rightarrow p_1 = 0.4, \quad p_2 = 0.6$$

$$\text{Variance } \sigma_1^2 = \frac{1}{n} \sum (f_{C_1} - \bar{f})^2 = \frac{(1-3.6)^2 + (5-3.6)^2 + (2-3.6)^2 + (6-3.6)^2}{4}$$

$$\Rightarrow \sigma_1^2 = \frac{1}{4} \cdot (25/4 + 9/4 + 9/4 + 25/4) = 68/16 = 17/4$$

$$\sigma_2^2 = \frac{1}{n} \sum (f_{C_2} - \bar{f})^2 = \frac{1}{6} \left[(3-3.6)^2 + (7-3.6)^2 + (4-3.6)^2 + (1-3.6)^2 + (5-3.6)^2 + (2-3.6)^2 \right]$$

$$\Rightarrow \sigma_2^2 = \frac{1}{6} (4/9 + 100/9 + 1/9 + 64/9 + 16/9 + 28/9) = \frac{210}{9 \times 6} = \frac{35}{9}$$

$$F = \frac{p_1 (\mu_1 - \mu)^2 + p_2 (\mu_2 - \mu)^2}{p_1 \sigma_1^2 + p_2 \sigma_2^2} = \frac{\frac{4}{10} \times \left(\frac{7}{2} - \frac{18}{5}\right)^2 + \frac{6}{10} \times \left(\frac{4}{3} - \frac{18}{5}\right)^2}{\frac{4}{10} \times \frac{17}{4} + \frac{6}{10} \times \frac{35}{9}}$$

$$= \frac{\frac{4}{10} \times \frac{1}{100} + \frac{6}{10} \times \frac{1}{225}}{\frac{4}{10} \times \frac{17}{4} + \frac{6}{10} \times \frac{35}{9}} = \frac{\frac{1500}{100 \times 225}}{\frac{68 \times 3 + 70 \times 4}{3 \times 4}}$$

$$= \frac{\frac{1}{10}}{\frac{482}{120}} = \frac{\frac{3}{10}}{\frac{482}{120}} = \frac{1}{1605} = 0.001653$$

Q.03. Total deficit = 1800 (All acc. in crores)

Prop. tax deficit = 600 (target = 1800)

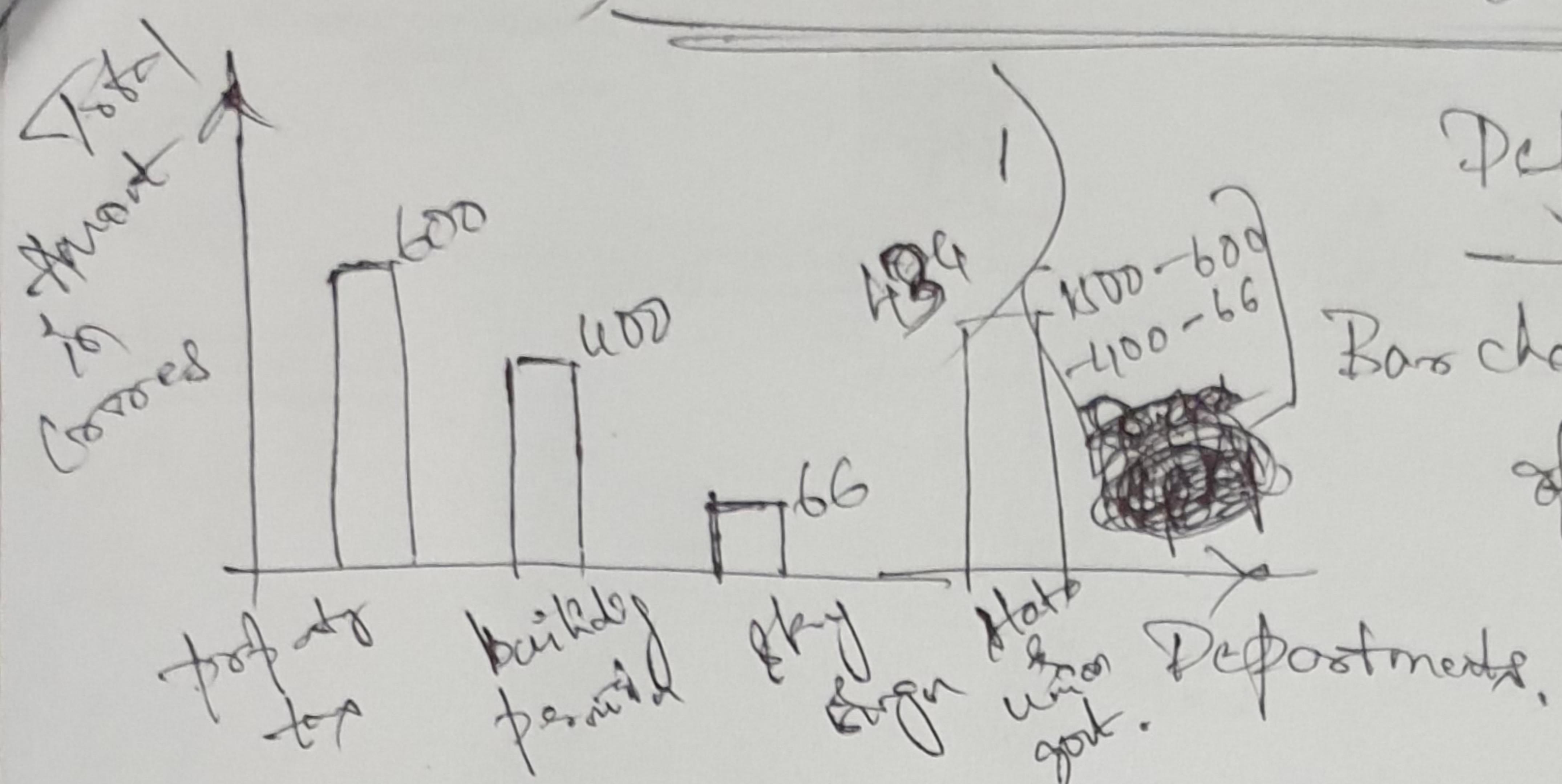
Other taxes & sky sign deficit = 66 (target = 90)

Other taxes & sky sign deficit = 66 (target = 90)
building promotion - [Per. charges] deficit
perm. fees] = 400
(target = 800)

Infund from govt's
State & union = 70% of total revenue.

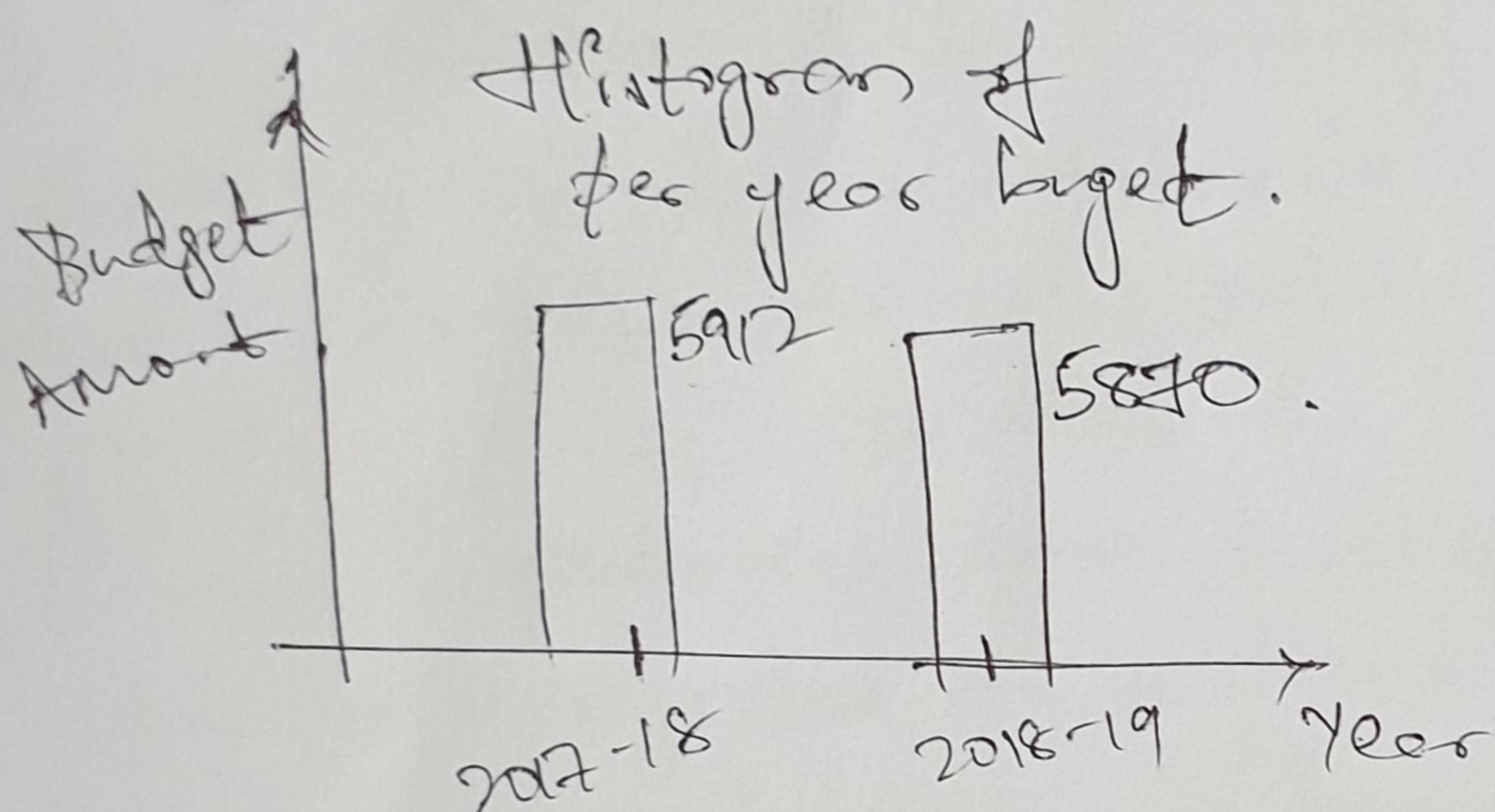
- i) Dominant reason for revenue deficit
 - PMC is dependent on grants by State & union govt. which have given only 70% of revenue & not yielding results on follow up
- ii) LBT (Local body tax has been abolished)
- iii) Individual departments have deficits and are not generating any profits (like property tax, other taxes including sky sign, etc.)
- iv) Story points are already described above.
plots in next page for data visualization.

Ado3. (contd. II) Plot of $\frac{1}{\text{f}} \text{ vs } \text{f}_{\text{Ku}}$ for visualization

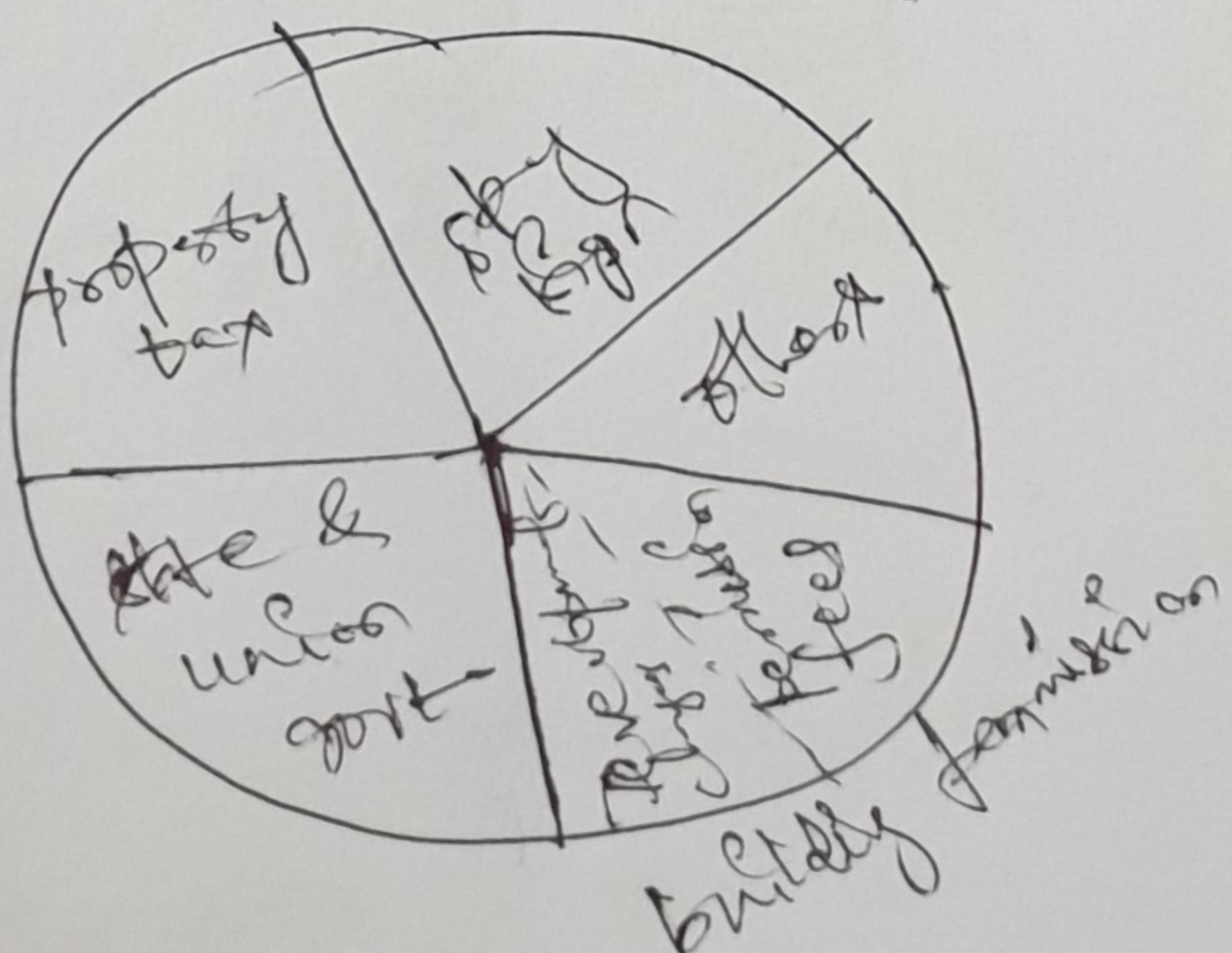


Perfect for department

~~Bank charged for distribution
of fee depositment
deficit.~~



PIE for department distribution %.



[Revenue generating
departments]

A) Hypothesis Testing

i) Hypothesis testing is to test for differences in samples via some statistic (mean, variance, etc.) or between models & model parameters.

ii) Hypothesis testing is usually done with a null hyp. H_0 which states "There is no / no significance". This is typically tested through a p-value which shows the probability of finding the value as extreme (at least as extreme) as observed just by random chance. If ~~keeps~~ the p-value is greater than some threshold then H_0 is rejected and a alternate hyp. H_a which typically has "there is some significance" is selected

iii) Hyp. is post-hoc. which is after finding relevant models / stats / parameters

Ex. In linear regression, (univariate) after finding slope of the line,
 $H_0 \Rightarrow \text{slope} = 0$ | we test if using t-stat / dev. tabulation
 $H_a \Rightarrow \text{slope} \neq 0$ |
 $\text{obs. value} = \text{observed value}$ | $t = \frac{\text{obs. value} - 0}{\text{std. error(obs. value)}}$

Exploratory Data Analysis (EDA)

i) EDA is analyzing & visualizing data to be done to identify patterns / usages for the data to build a suitable model

ii) EDA is done through plots / charts & graphs to get more insights into the data & the process that generates the data (underlying)

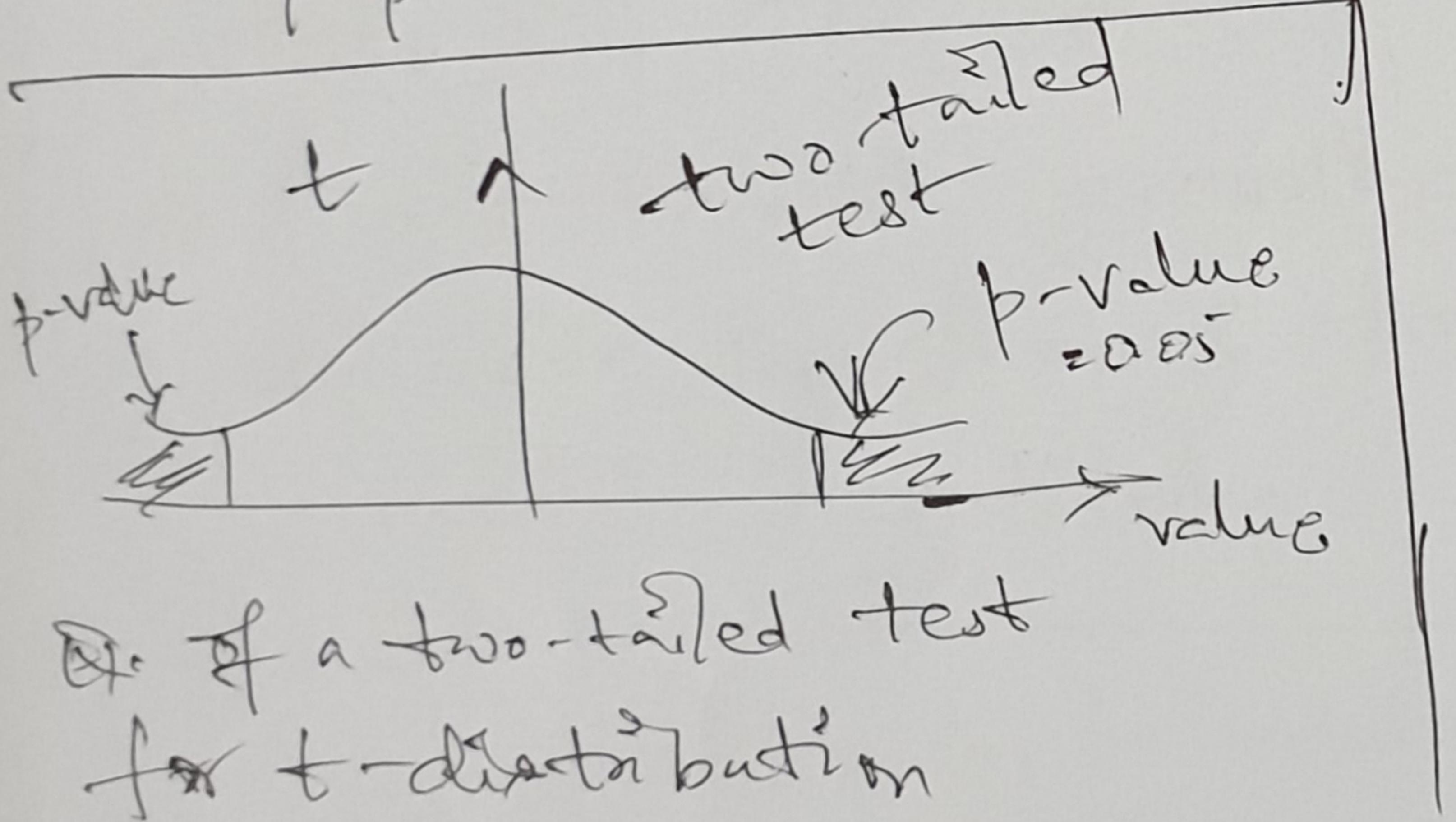
iii) EDA is prior to finding stats / models / posars & involves finding statistics & estimating parameters

Ex. Finding data summary (mean, variance, distribution), plotting them, finding null / missing values

4 - Contd. Hyp.

(ii) t-stat says, how far the slope is from '0' & the probability of finding such a value by chance if H_0 is true
 Note: For multiple reg. this will be F-stat.

Other tests are between samples for $\mu = 0$ & $\sigma^2 = 0$
 where H_0 = samples are from same population.
 H_a = samples are not from same population



Ex. of a two-tailed test
 for t-distribution
 p-values are typically 0.05 or 0.01

EDA

(iv) EDA gives outputs in tables, also in graphs and charts which can be univariate/multi-variate
Ex. box, histogram (cum^2)
 scatter contour (multi)
 b)

- Q.Q. Contd. B) i) Variables with 90% missing values.
→ this variable will be dropped as there is only 10% information. We cannot assume the distribution from which the data was generated or the probability of it or the likelihood of data.
→ as population is not known, we cannot estimate any statistic (mean, dispersion, etc.) of the population parameters & hence will discard it.
- ii) Variables strongly correlated carry the same information & hence will be dropped as well. They will face difficulty in optimization.
→ For probabilistic generative methods, this will overestimate the prior & probabilities affecting the model negatively.

- C) NMR data can be imputed by these 2 ways -
- Add values at the end of distribution
 - Any random value (which is not the mean / median)
 - Also add a undergone indicator column & check model performance.
- ~~Male~~
- In given data set →
- $\Rightarrow \text{Male} = 57$
 - $\Rightarrow \text{Avg. of highest or least 2 values, etc.}$
 - \Rightarrow Additional column.
- ~~Female = 18~~
- ~~57 tail / end~~
~~# distribution~~

