

Data Science Certification Program

1

Course : Machine Learning

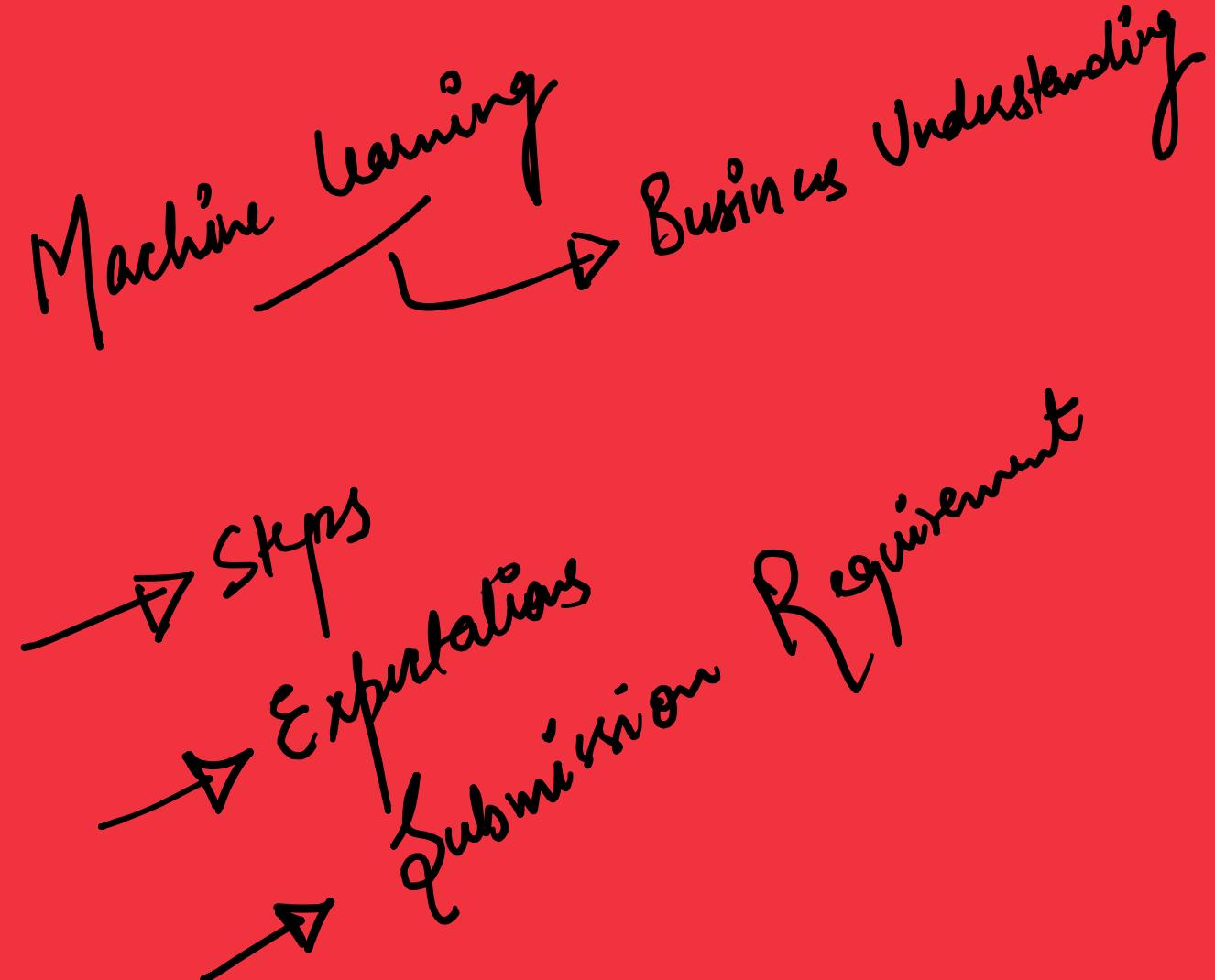
Lecture On : Lead Scoring
Assignment

Instructor : Shivam Garg



Today's Agenda

- 1 Problem Statement
- 2 Assignment Walkthrough
- 3 Doubt Session



Data Science Certification Program

Assignment: Problem Statement

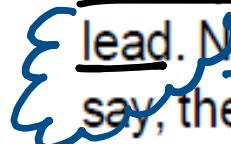
[Ed-tech Industry]

upGrad

Classification
problem

→ upgrad

An education company named X Education sells online courses to industry professionals. The company markets its courses on several websites and search engines like Google. Once these people land on the website, they might browse the courses or fill up a form for the course or watch some videos. When these people fill up a form providing their email address or phone number, they are classified to be a lead. Now, although X Education gets a lot of leads, its lead conversion rate is very poor. For example, if, say, they acquire 100 leads in a day, only about 30 of them are converted.

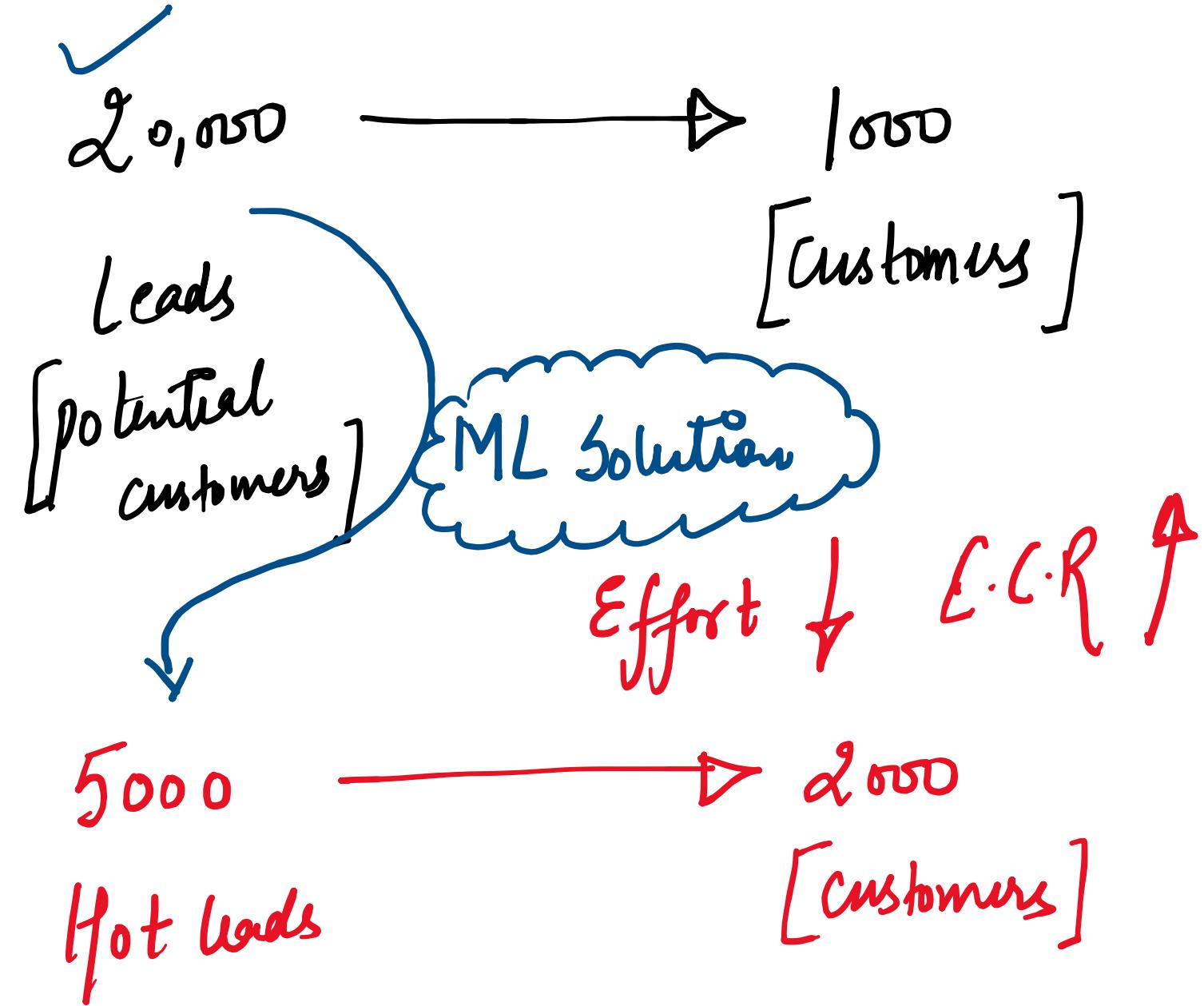


What you need to do?

→ hot leads

- X Education has appointed you to help them select the most promising leads, i.e. the leads that are most likely to convert into paying customers.
- The company requires you to build a model wherein you need to assign a lead score to each of the leads such that the customers with higher lead score have a higher conversion chance and the customers with lower lead score have a lower conversion chance.

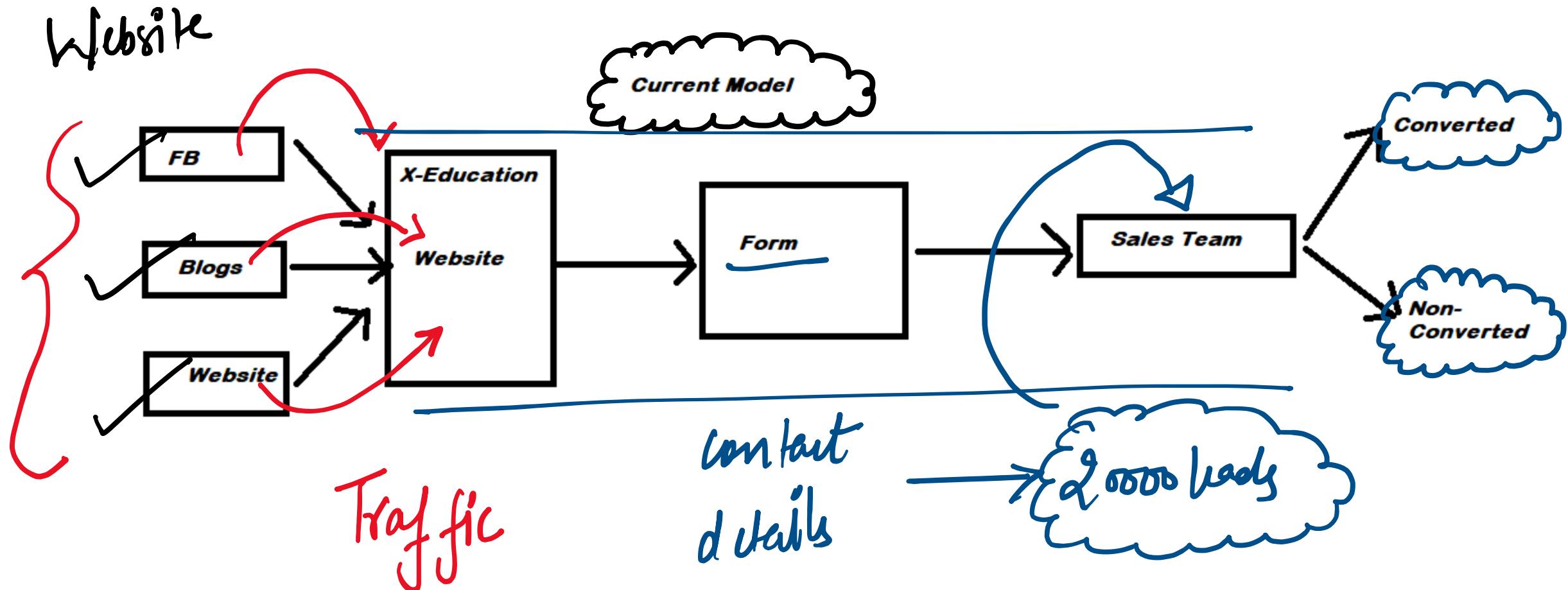
$$\text{Lead Score} = 100 \times \frac{\text{nob.}}{\text{total}}$$



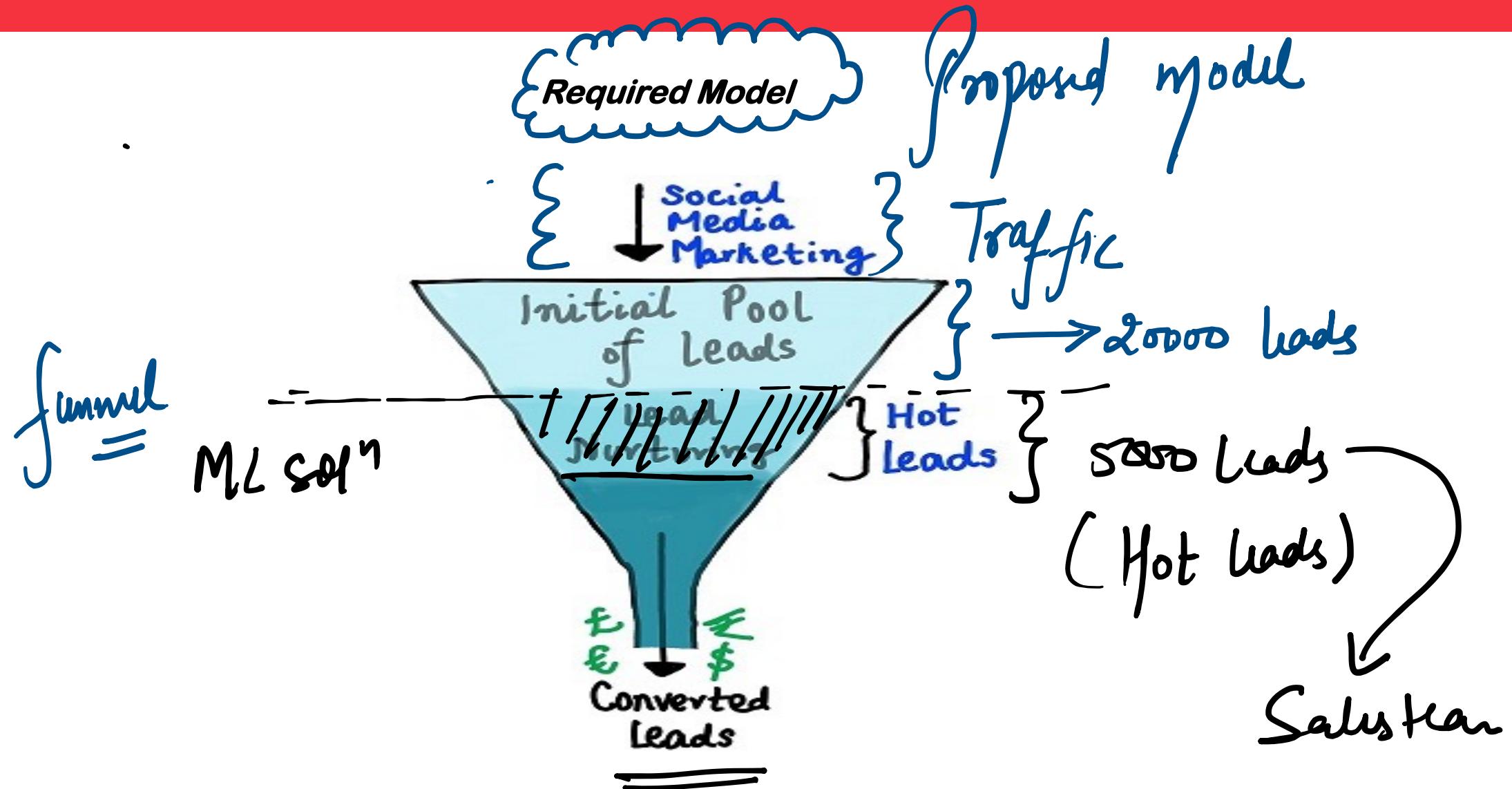
Lead conversion
 rate = $\frac{1000}{20000} \times 100$
 Ratio = 5 %

L.C.R = $\frac{2000}{5000} \times 100$
 = 40 %

Assignment: Problem Statement



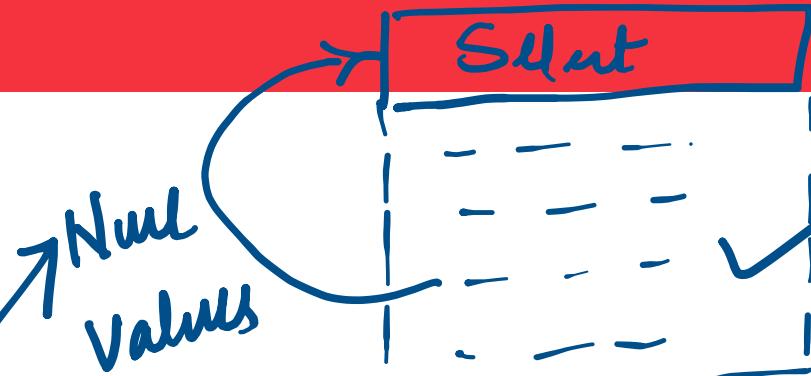
Assignment: Problem Statement



→ Assignment Steps:-

- 1- Data Cleaning / Data preprocessing. {
 - 2- Data preparation } →
 - 3- Model Building & Evaluation }.
-
- ```
graph TD; A["1- Data Cleaning / Data preprocessing."]; B["2- Data preparation"]; C["3- Model Building & Evaluation"]; A --- B --- C; B --- D["{ }"]; A --- D; C --- E["{ }"]; D --- E;
```

## Data Cleaning:



1. Handle the “Select” level that is present in many of the categorical variables.
2. Drop columns that are having high percentage of missing values. [Check all the columns before dropping them.]
3. Check the number of unique categories in each categorical column. Here you may need to do something.
4. For the columns with less percentage of missing, use some imputation technique.
5. Finally check the percentage of rows retained in data cleaning process.
6. Outlier Treatment. & Columns

5 Data Sanity Checks

A

% of missing values

15% X

→ 75% ✓

Select → np.nan } Null

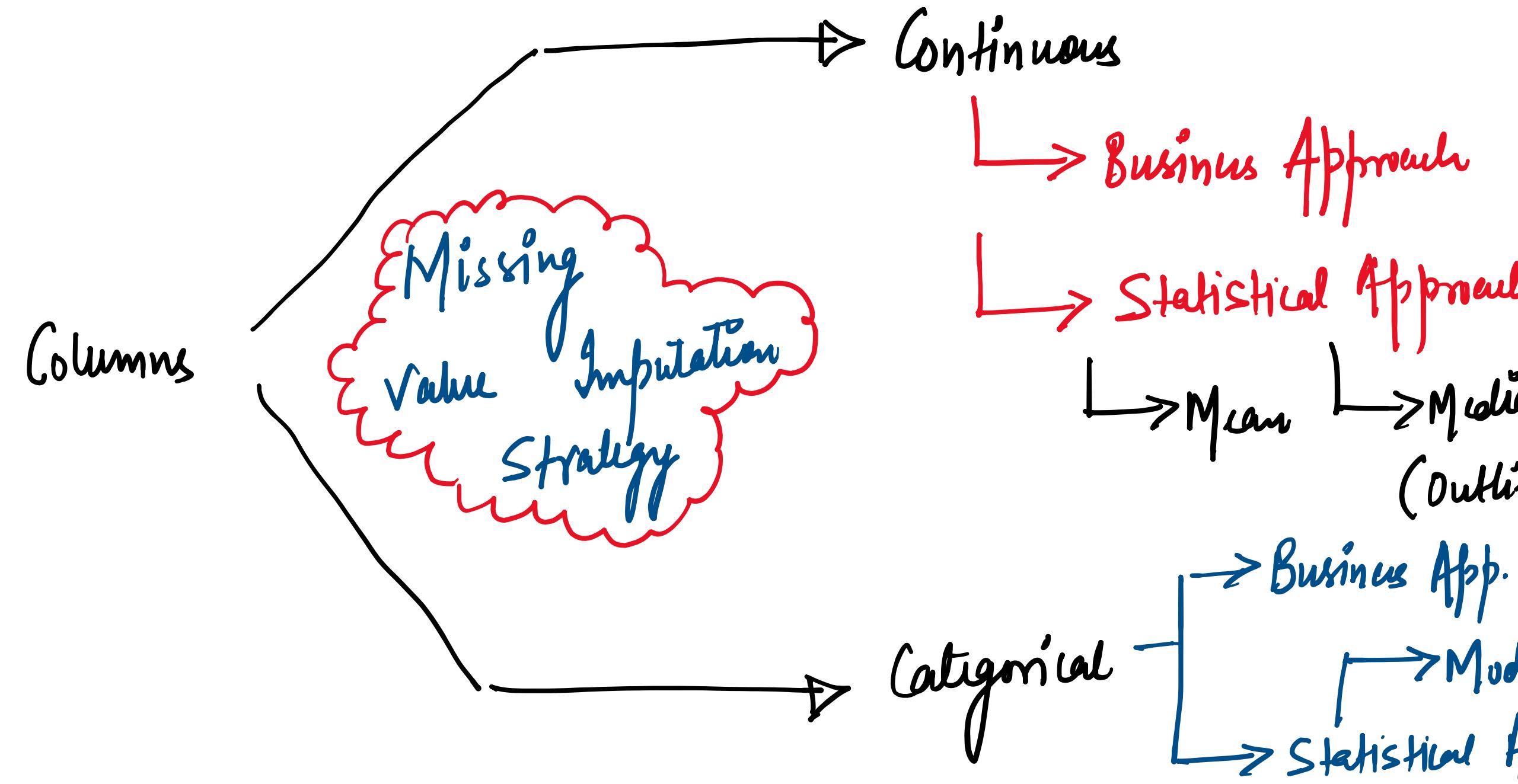
% of missing values → Threshold → Drop columns  
(30-50%) ( $>$  threshold)

Country (100 countries)

|           |       |
|-----------|-------|
| India     | 75%   |
| U.S       | 15%   |
| U.K       | 5%    |
| Australia | 2%    |
| Pak       | 0.1%  |
| Srilank   | 0.01% |
| :         | :     |
|           |       |

Country (5 categories)

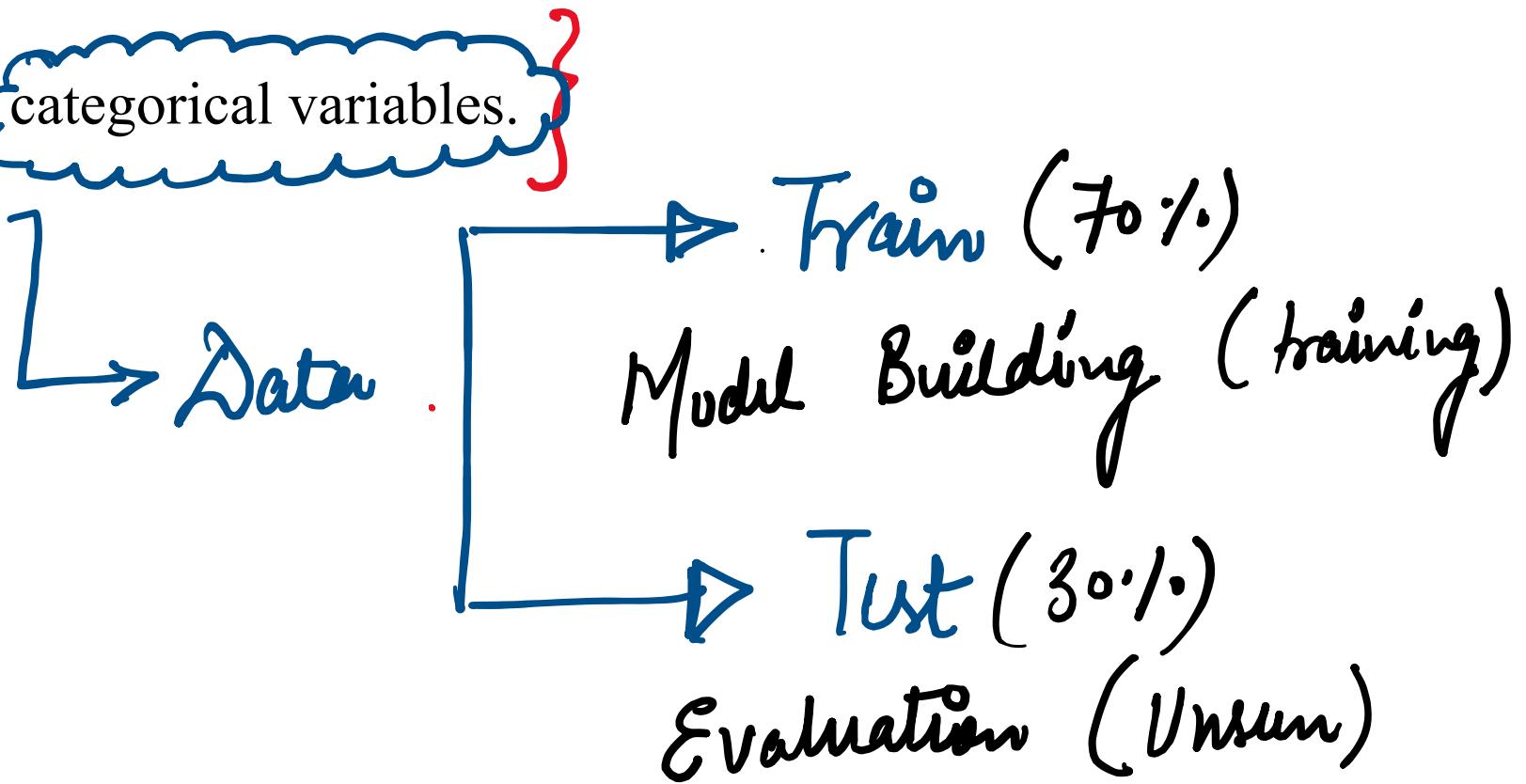
|        |     |
|--------|-----|
| India  | 75% |
| U.S.   | 15% |
| U.K    | 5%  |
| Aus    | 2%  |
| Others | 3%  |



## Data Preparation:

- ✓ 1. Create dummies for all categorical variables.
- ✓ 2. Perform train-test split.
- ✓ 3. Perform scaling.

↓  
[Mandatory]



Categorical column

A

B

A

G

B

A



|   | A | B | C |
|---|---|---|---|
|   | 1 | 0 | 0 |
| A | 0 | 1 | 0 |
| B | 1 | 0 | 0 |
| A | 0 | 0 | 1 |
| G | 0 | 1 | 0 |
|   | 1 | 0 | 0 |

No. of categories(n) = 3

No. of dummies = n - 2

`pd.get_dummies(df[categorical], drop_first=True)`

↓  
first will get drop

{Assignment Steps : }  $P\text{-value} \leq 0.05$   $VIF \leq 5$

upGrad

$$30 \xrightarrow{\text{RFE}} 15 \xrightarrow{\text{P-value}} 12 \xrightarrow{\text{VIF}} 10$$

### Modelling:

1. Use techniques like RFE to perform variable selection.

2. Build a Logistic Regression model with good sensitivity.

3. Check p-value and VIF. *Iterative*

4. Find the optimal probability cutoff. *[ROC-Curve]*  $\rightarrow$  Recall

5. Check the model performance over the test data. (Confusion Matrix, Sensitivity, F1-score etc.)

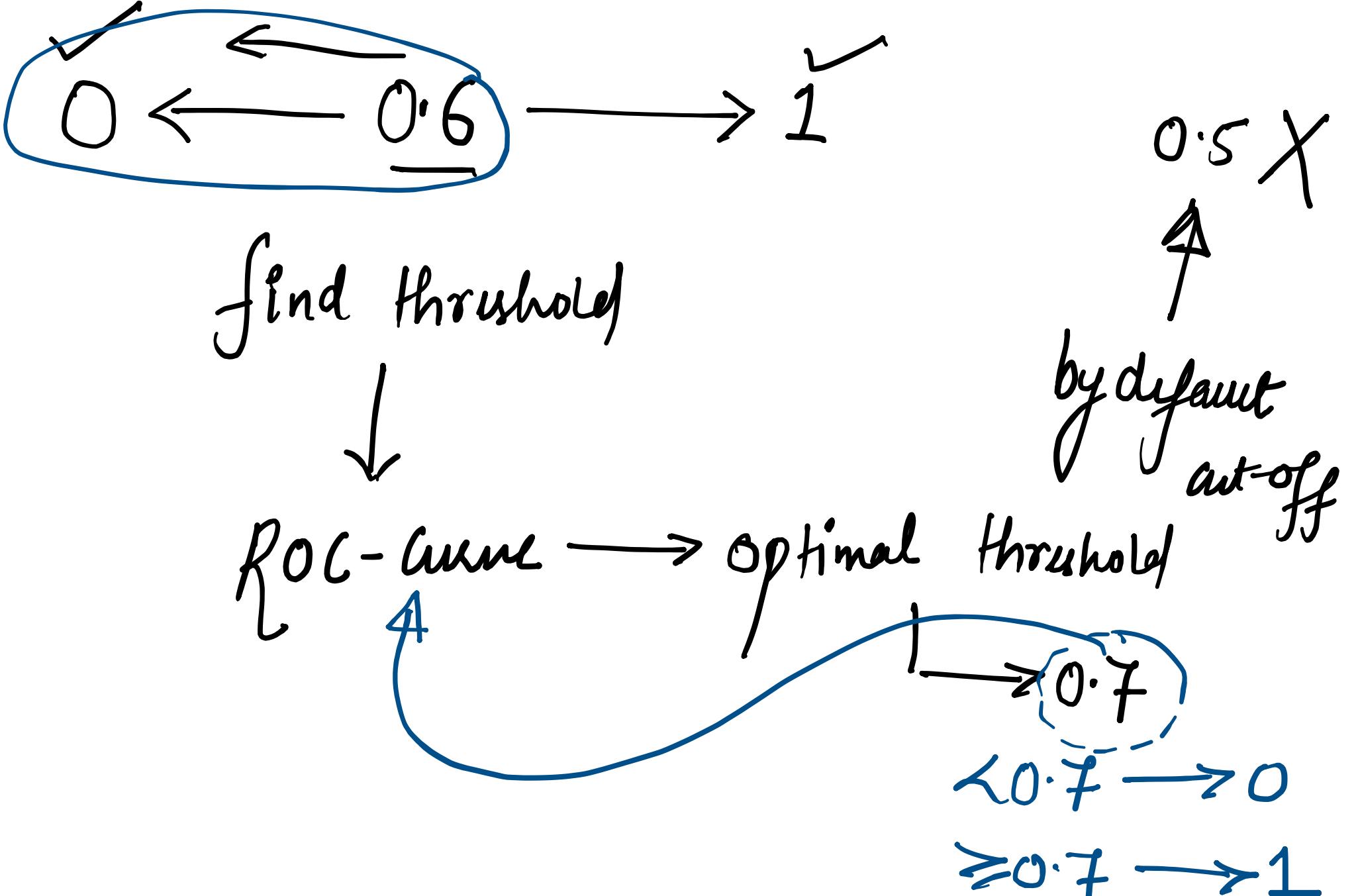
6. Generate the score variable.

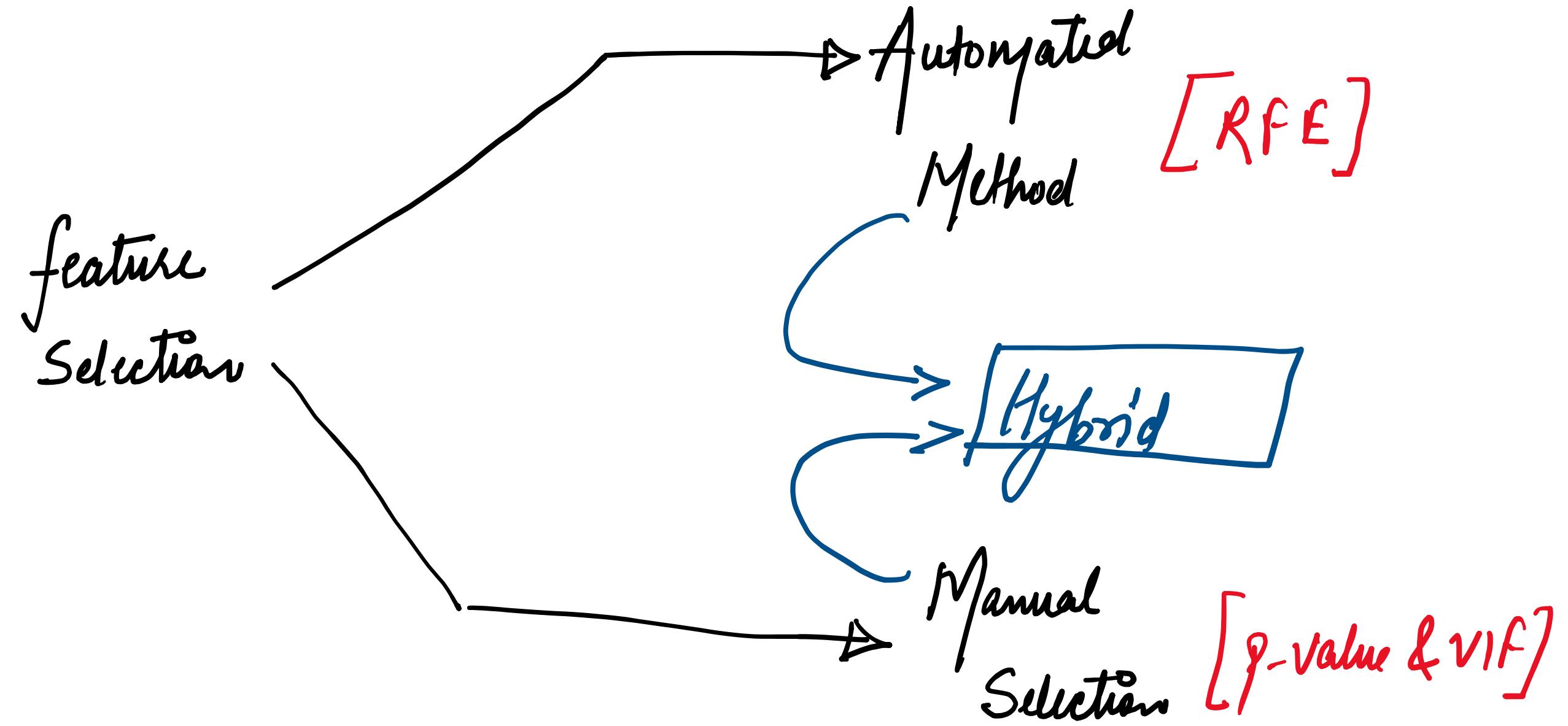
$$\text{Score} = \text{prob.} \times 100$$

$\swarrow$  final features

$\nearrow$  final model

|       |                   |           |
|-------|-------------------|-----------|
| prob. | $0.6 \rightarrow$ | Score     |
|       |                   | <u>60</u> |





Assignment Steps : *(EDA + Data preprocessing + Data prep + modeling + evaluation + ROC + Score)*

upGrad

## Submission:

- ✓ 1. Jupyter Notebook: A well-commented Jupyter notebook with at least Logistic Regression model, the conversion predictions and evaluation metrics.
- ✓ 2. Subjective Answers: The word document filled with solutions to all the problems.
- ✓ 3. The overall approach of the analysis in the presentation:
  - ✓ Mention the problem statement and the analysis approach briefly.
  - ✓ Explain the results in business terms.
  - ✓ Include visualizations and summarize the most important results in the presentation.
- ✓ 4. Summary Report : A brief summary report in 500 words explaining how you proceed with the assignment and the learning that you gathered.

- Jupyter Notebook → .ipynb format
  - Subjective Quest<sup>n</sup>/ps → Word (PDF)
  - Presentation → PPT → (PDF)
  - Summary Report (Word) → PDF
- } Zip file  
↓  
Upgraded portal

PPT (12-15 Slides)

→ Problem Statement + Assumptions (if any)

→ Approach

→ Graphs (EDA Results)

→ Conclusion or Result (Model performance, Conf matrix)

→ Recommendation (Insights, Impact on business,  
Recommendations)

## Points to keep in mind

- 1- Do not mention code in your ppt.
- 2- Only put relevant results (Important)
- 3- Graphs should be readable (not blur) & well-labelled
  - x
  - y
  - title
- 4- Do not write too much theory in your ppt. Always work in bullets.

Short  
Crisp  
Clear

## What to keep in mind:

1. Add comments after every cell of code. So that we can understand your approach and method.
2. Describe the results.
3. Use Stack Overflow for dealing with syntax error, rather than being stuck at one place or waiting for someone to resolve your doubts, act and use the resources available on the internet to save time.
4. Post on the discussion forums for resolving any doubts you have.
5. Finally, write code manually instead of copy pasting from the in-content notebooks provided. Builds a habit of writing code. Its okay to look and write, but do not just copy-paste under any circumstance. Because of just copy pasting a lot of students have faced difficulties in the past when they had to write same code in their interview.