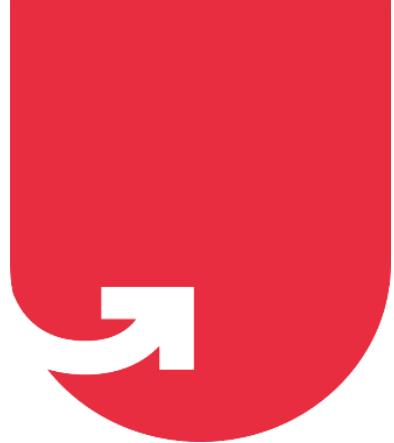


**When tomorrow
is the last date of assignment
submission**



Data Science Certification Program

2

Course : Machine Learning

Lecture On : Linear
Regression Assignment

Instructor : Shivam Garg



Today's Agenda

- 1 Assignment Walkthrough
- 2 Step by Step Approach
- 3 Doubt Session (QnA)



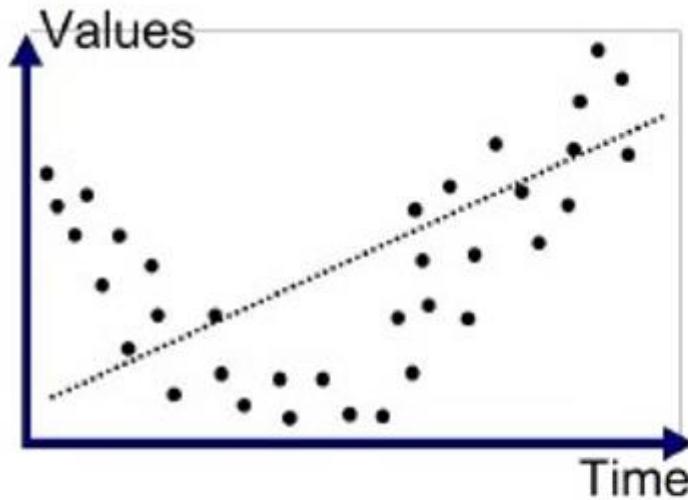
Linear Regression:

Quick Revision

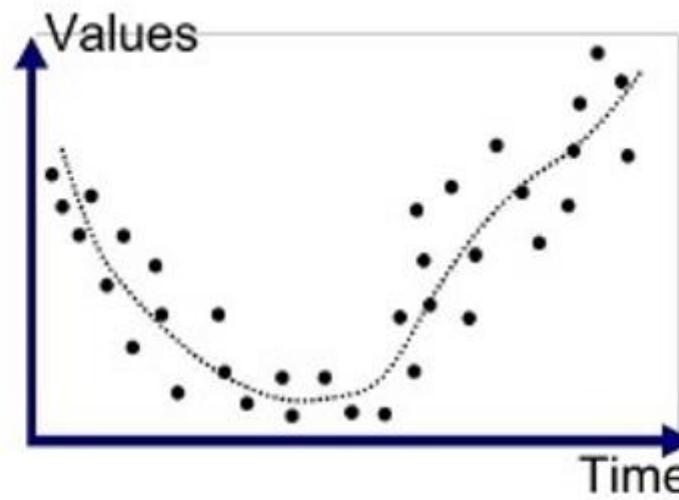
Years of Experience	Salary in 1000\$
2	15
3	28
5	42
13	64
8	50
16	90
11	58
1	8
9	54



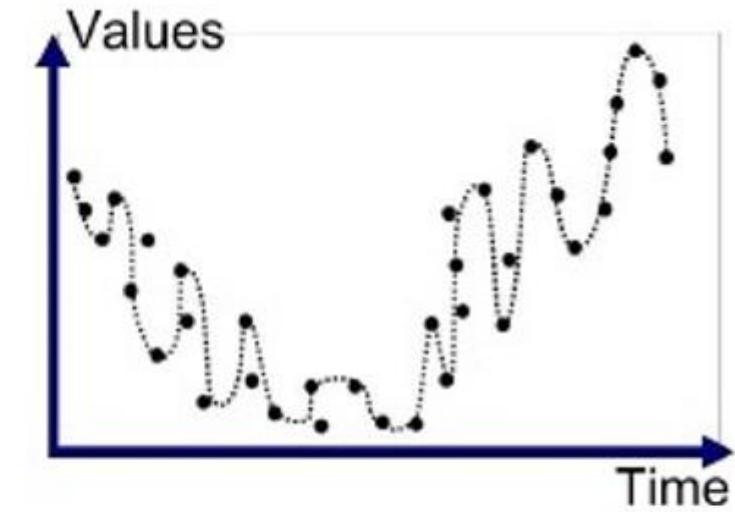
Linear Regression:



Underfitted

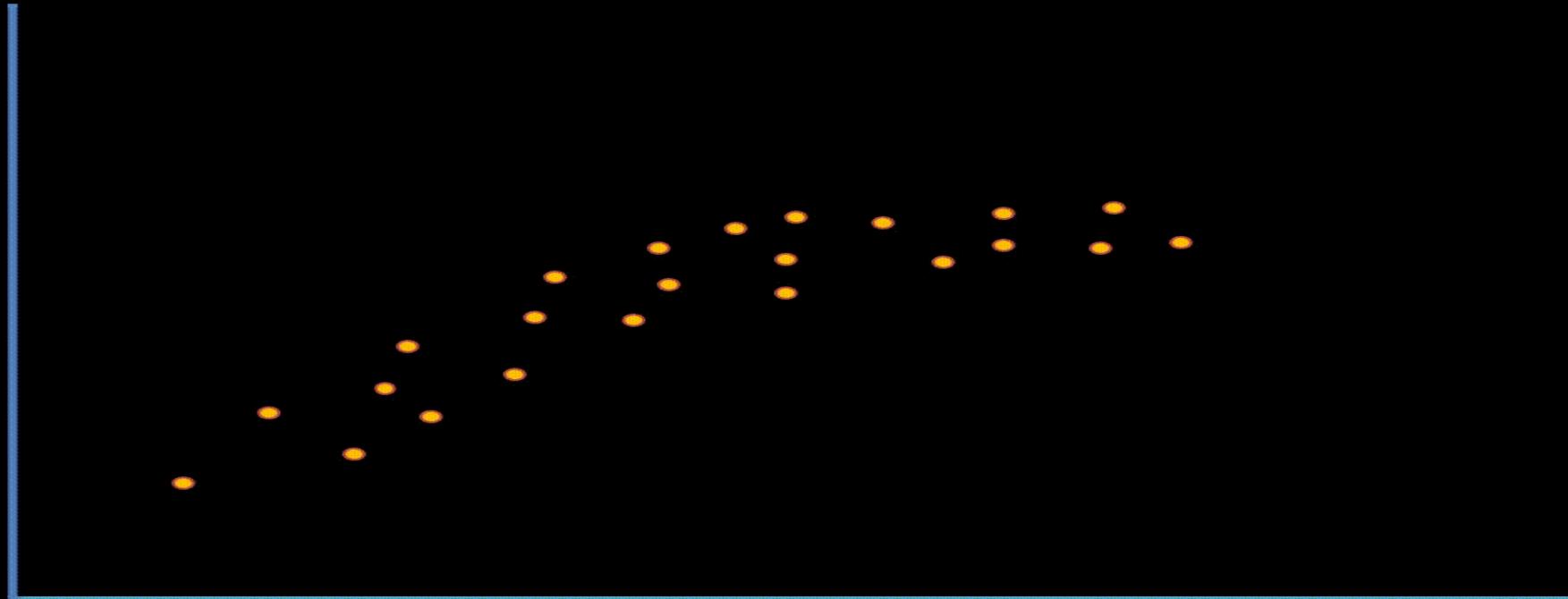


Good Fit/R robust



Overfitted

Overfitting and Underfitting



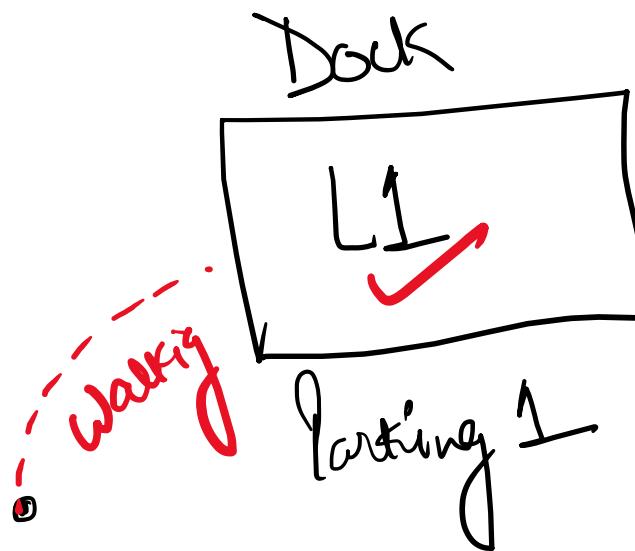


Assignment Problem Statement

Bikesharing

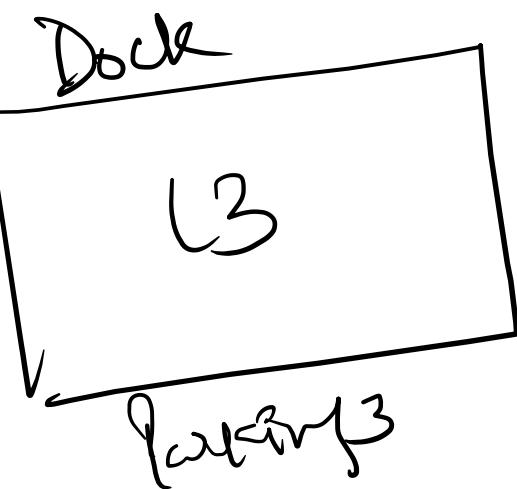
- A bike-sharing system is a service in which bikes are made available for shared use to individuals on a short-term basis for a price or free. Many bike share systems allow people to borrow a bike from a "dock" which is usually computer-controlled wherein the user enters the payment information, and the system unlocks it. This bike can then be returned to another dock belonging to the same system.
- A US bike-sharing provider BoomBikes has recently suffered considerable dips in their revenues due to the ongoing Corona pandemic. The company is finding it very difficult to sustain in the current market scenario. So, it has decided to come up with a mindful business plan to be able to accelerate its revenue as soon as the ongoing lockdown comes to an end, and the economy restores to a healthy state.

Mumbai

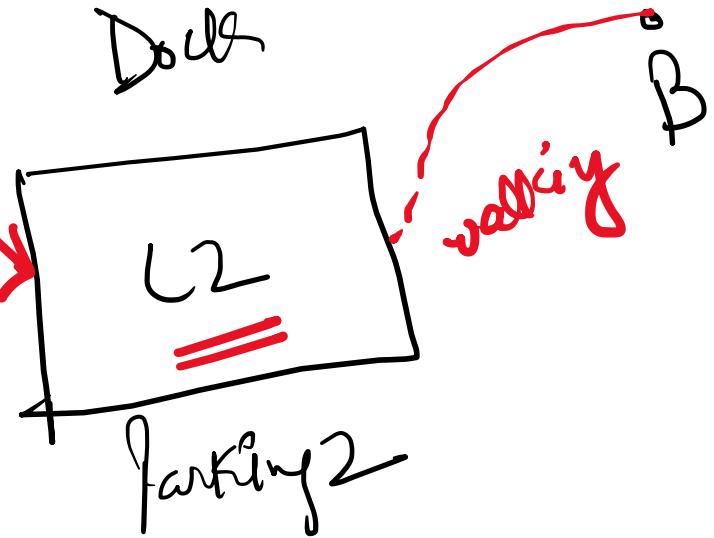


walking

A



Ride
using
bike



Dock

Parking 2

B

walking

X/Y/Z

variables
y-axis

- They have contracted a consulting company to understand the factors on which the demand for these shared bikes depends. Specifically, they want to understand the factors affecting the demand for these shared bikes in the American market. The company wants to know:

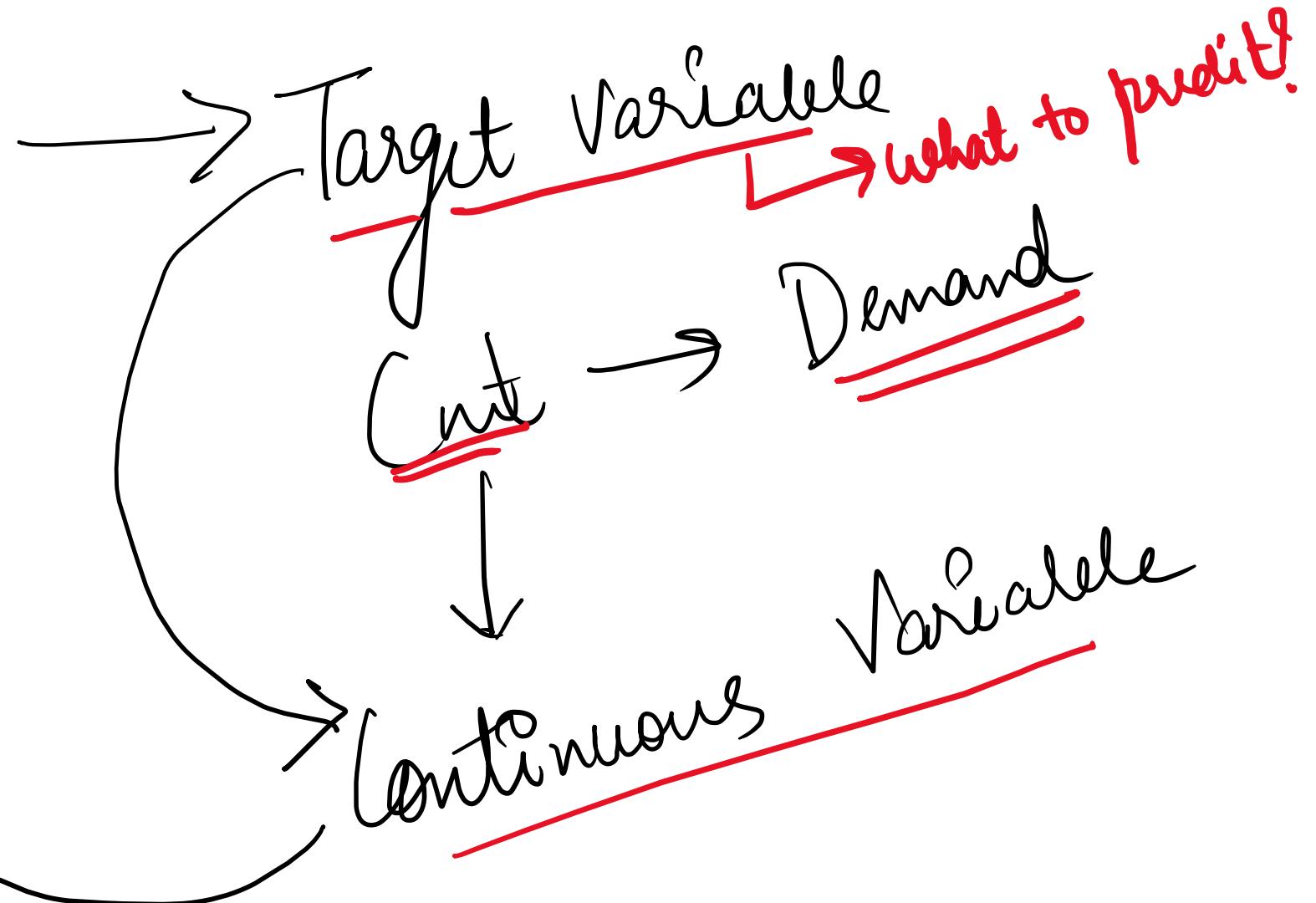
1. Which variables are significant in predicting the demand for shared bikes.
2. How well those variables describe the bike demands

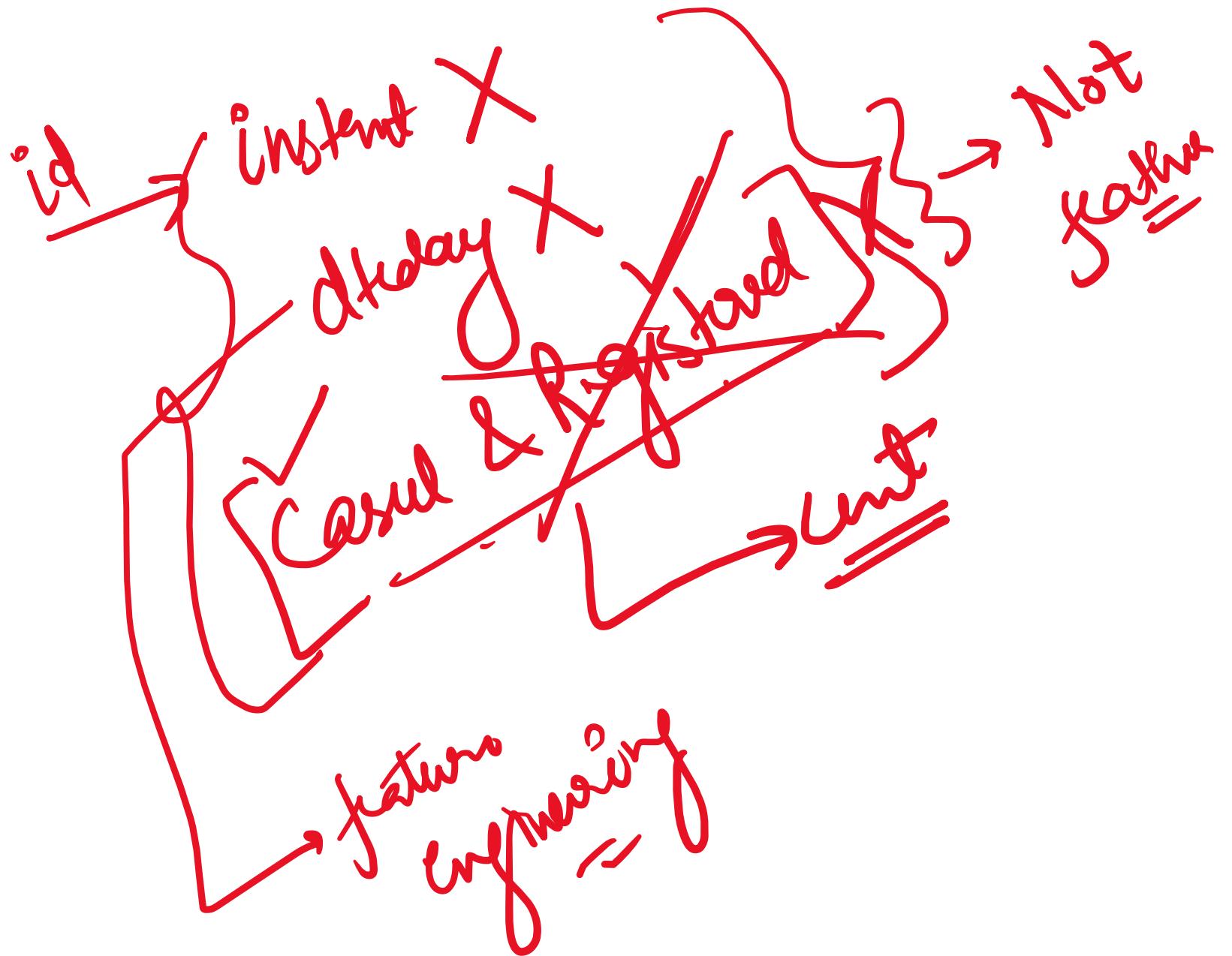
20+

What you need to do?

- Create a linear model that describes the effect of various features on price.
- The model should be interpretable so that the management can understand it.

Regression





① Understand the problem Statement & understand
the data dictionary.

② Load the data



```
graph TD; A[② Load the data] --> B["fd.read_csv()"]
```

③ Sanity Check

- (i) Shape of the data
- (ii) No of columns
- (iii) Missing values
- (iv) Data types ex int

④ Data prepossessing

↳ data driven

↳ Encoding

Weather sit

1
2
3
A

weather

1 → summer
2 → winter
3 → rainy
4 →

⑤ Dummy Creation \rightarrow Only Categorical Variables.

unordered Weathersit

1	1	0	0	0
2	0	1	0	0
3	0	0	1	0
4	0	0	0	1
5	1	0	0	0
6	0	1	0	0
7	0	0	1	0
8	0	0	0	1

No. of unique values = $A = n$ (drop-first = True)

No. of dummies = 3

$n = 8$ (n+1)

w-1 w-2 w-3 w-4

⑥ EDA

Exploratory Data Analysis (Univariate
+
Bivariate)

⑦ Train - Test Split

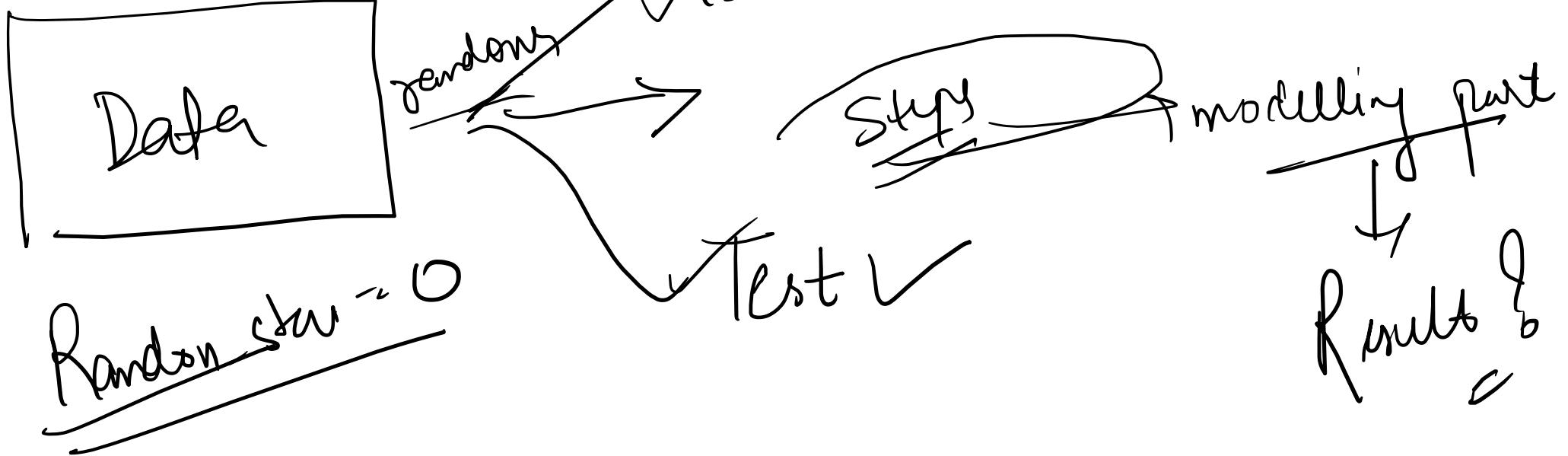
Data
Stratify

Train data → Train my model
Random Split $70:30$, $75:25$, $80:20$
Test data
[most preferred]

from sklearn.preprocessing import train_test_split
~~random_state = 42~~
 $(70, 30)$

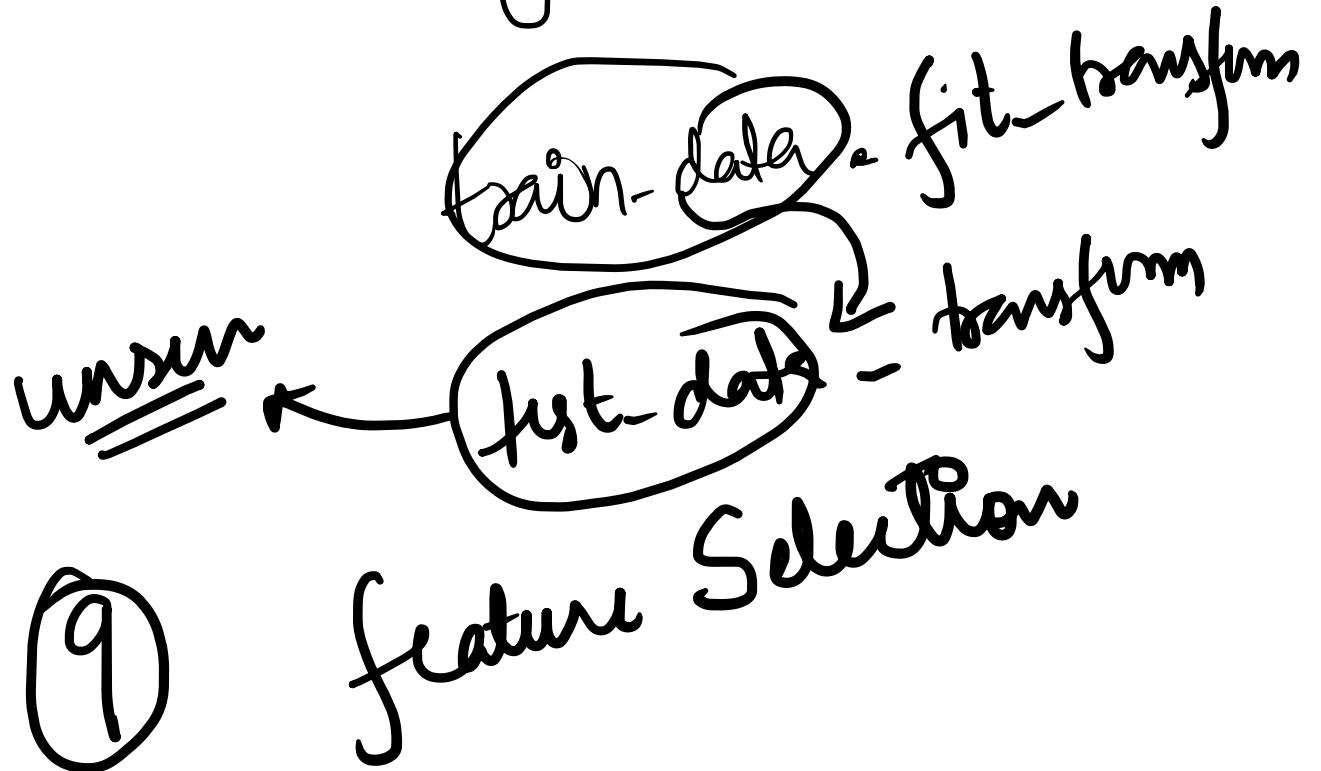
Model
Stratify = ==

Abholt's you off



Reproducibility
→ Randomness

⑧ Scaling → Minmax Scaler / Standard Scaling



⑨ feature Selection

evaluating \bar{x}_1, \bar{y}

$$30 \rightarrow \frac{5}{30}$$

test data

$$\frac{30 - 25}{30} = \frac{5}{30} \quad 2^{\circ} \xrightarrow{\text{Scaling}} 1/6$$

transform

$$\frac{2^{\circ} - 25}{30} = \frac{5}{30} = \frac{1}{6}$$

$$\frac{x - \bar{x}}{s}$$

$$\frac{30 - 30}{30} \approx 0$$

$$\frac{2^{\circ} - 30}{6} = -6$$

$$\frac{30}{1} = 30$$

Product "n"
↳ you will get one
hit

$$\left(\begin{array}{l} w_1 = 0 \\ w_2 = 1 \\ w_3 = 0 \\ w_4 = ? \end{array} \right)$$

$$\left(\begin{array}{l} w_1 = 0 \\ w_2 = 0 \\ w_3 = 0 \\ w_4 = 1 \end{array} \right) \xrightarrow{\text{3}} \text{z}$$

multicollinearity

① feature Selection

- (i) RFE → Automatic way of feature Selection
- (ii) VIF → Manual method
↳ p-value / VIF
- (iii) Hybrid (RFE + VIF)

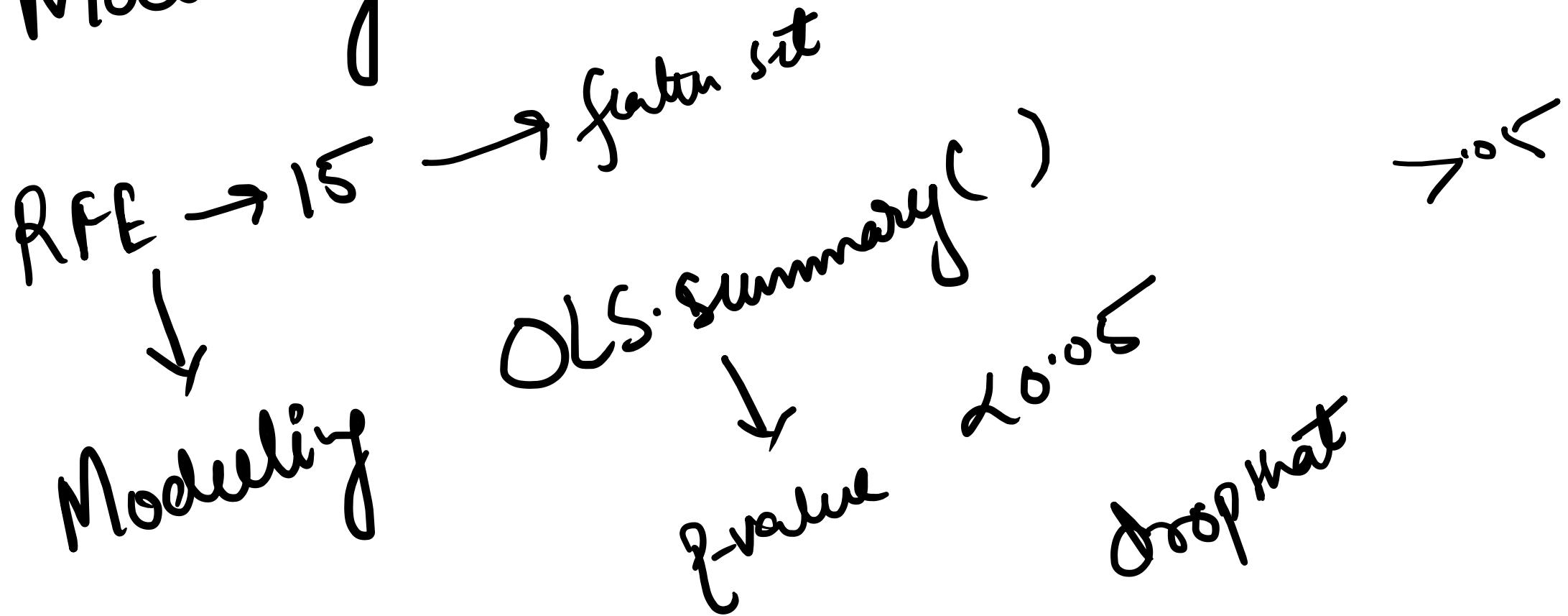
→ 33 features

→ RFE (15) variable → 15
keep 15

→ VIF → 15
L

VIF > 5
p-value > 0.05

⑩ Modelling & Evaluate



Drop one feature

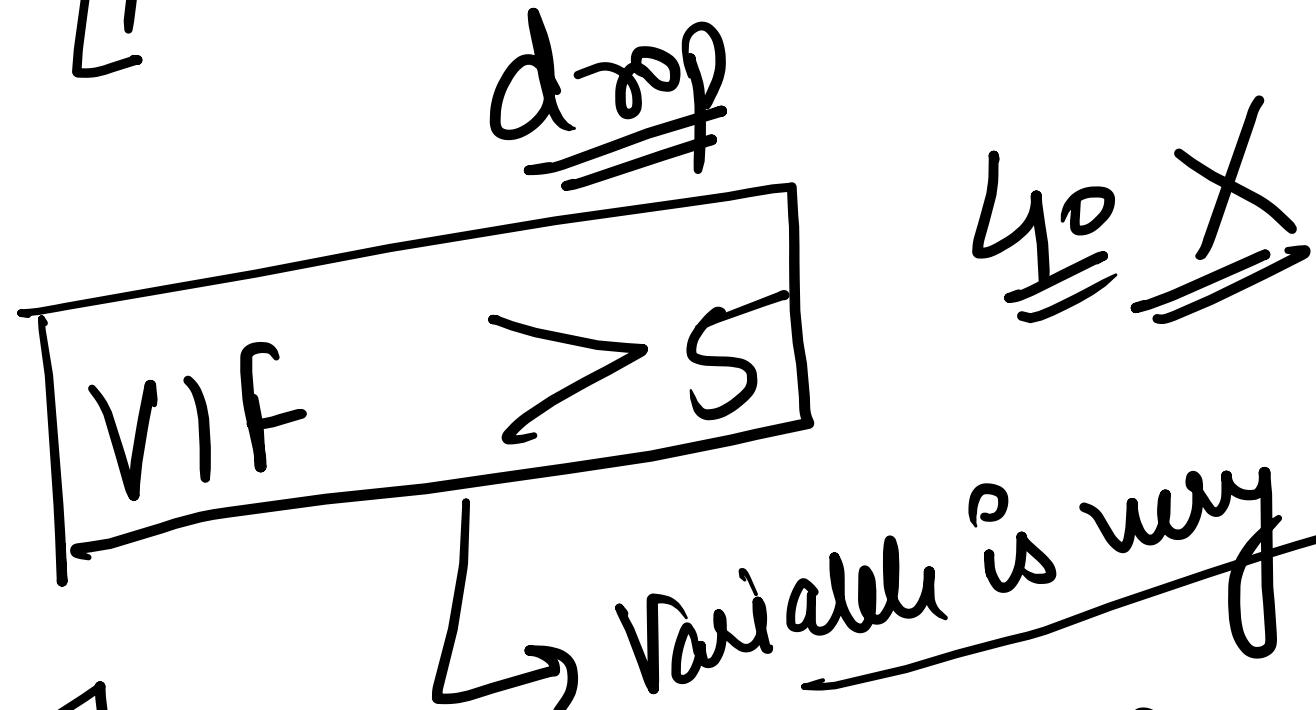
$\checkmark x_1$
 ~~x_2~~
 x_3

$0.06 \} \text{ altogether}$
 0.1

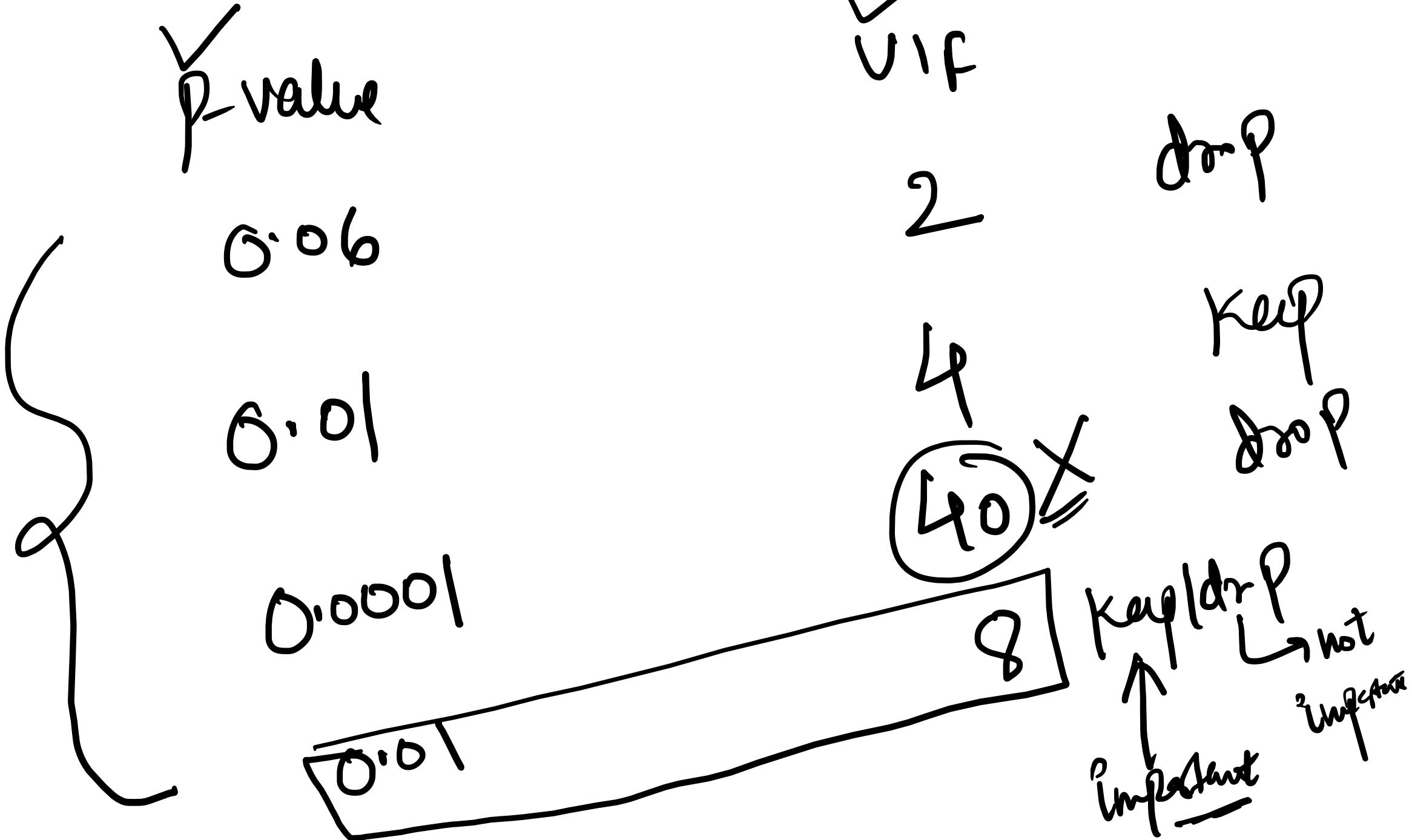
x_1
 x_3
 ~~0.05~~

$\cancel{x}^{0.5}$

$P\text{-value} > 0.05$ 3 exception



1
from business perspective you
can still keep it



Adj R^2

Residual Analysis

(i) Multi collinearity
(ii) Plot y_{pred} vs y_{actual} .
(iii) Plot the error term.

VIF

↳ Multicollinearity

↳ two independent features are related with each other

$$= \frac{1}{1 - R^2}$$

$R^2 = \text{vne}$

$$\text{Adj } R^2 = 1 - \frac{(1-R^2)(N-1)}{(N-P-1)}$$

Age
20
30
40
50
25

Transform

$$\frac{20-21}{10}$$
$$\frac{30-21}{10}$$
$$\frac{40-21}{10}$$

Test → Transform
21

Train

fit

$$\left\{ \begin{array}{l} \bar{x} = 21 \\ \sigma = 10 \end{array} \right.$$

fit

$$\frac{31-21}{10} = \frac{10}{10} = 1$$

Test
date

→ 1 data point

0 → Age
30

$$\bar{x} = 30$$

$$\frac{30 - 50}{\sigma} = \text{train data}$$

2 35
~~35 - 5~~
~~6~~
= 0

Final No of feature

6-10

Adjusted R² ~ 80%.

[78-82]

Data Preparation:

- Load the data and understand it using dictionary provided.
- Convert the columns to proper data types.
- Create dummies for categorical variables.

Model Building

- Divide the data to train and test.
- Perform scaling.
- Divide the data into X and y.
- Perform Linear Regression.
- Use mixed approach if you want.

How to go about selecting features for a good model?

- RFE
- Manual
- Mixed

Assignment Steps

Model Evaluation

- Check the various assumptions.
- Check the Adjusted R-Square for both test and train data.
- Report the final model.

Assignment-Subjective

Steps to answer subjective part

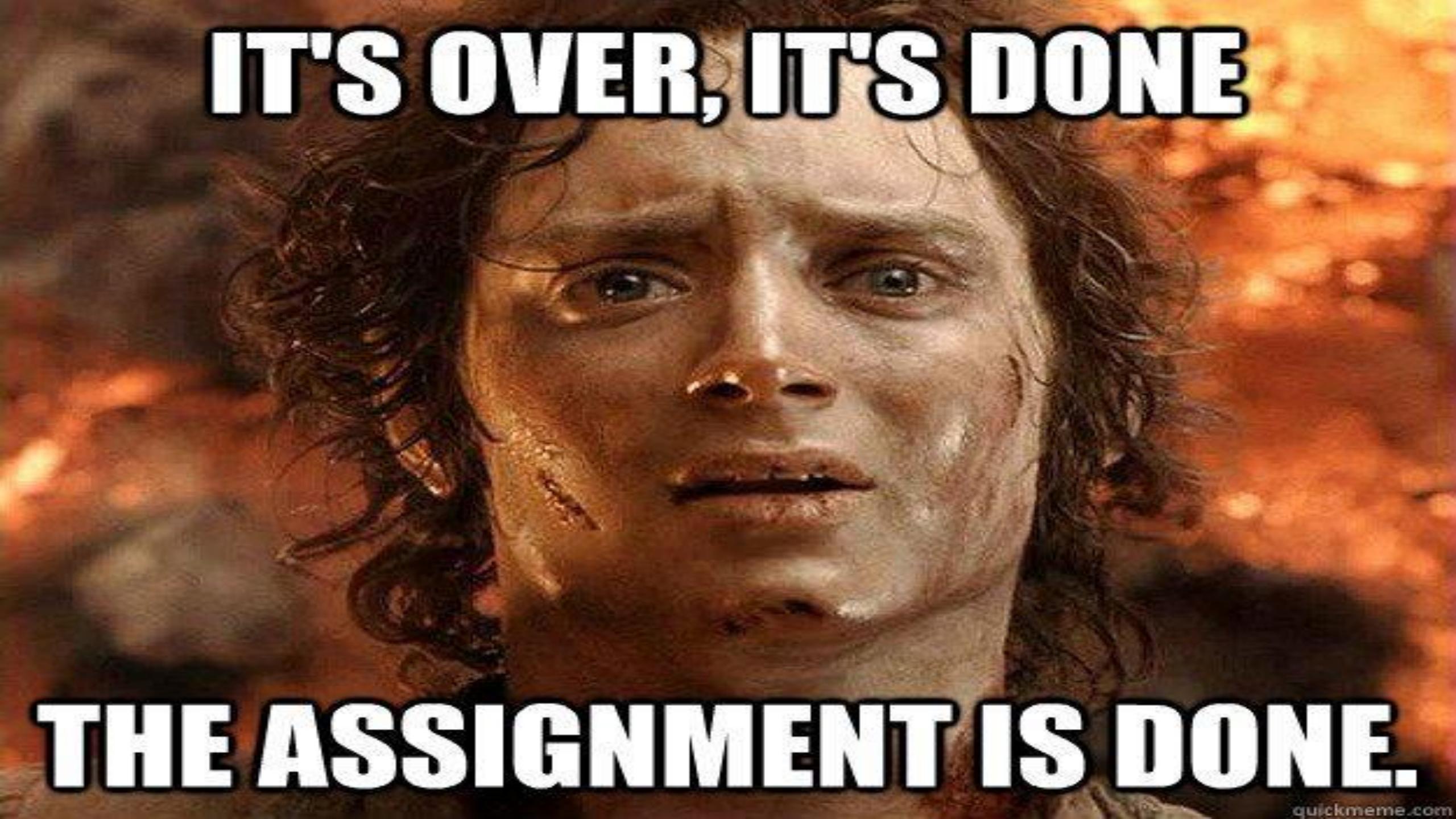
- Answer all the questions.
- You can write the answer using any software but submit the file in PDF format
- You can use images and plots to support your answer.
- Make sure the question is answered with sufficient number of word: No limit
- Please don't copy for any online available literature.

Assignment-Endnote

JN + PDF → 2i8
upl^{oad}

What to keep in mind

- Add comments after every cell of code. So that we can understand your approach and method.
- Describe the results.
- For subjective answers, use DOC and type on it, if you wish to add images you can. But convert it to PDF before submitting.
- Create only one Jupyter notebook.
- Submit one zip file with the code and the PDF.
- Use StackOverflow for dealing with syntax errors. Rather than being stuck at one place or waiting for someone to resolve your doubts, take action and use the resources available on the internet to save time.
- Post on the discussion forums for resolving any doubts you have
- Finally, write code manually instead of copy-pasting from the in-content notebooks provided. Builds a habit of writing code. It's okay to look and write, but don't just copy-paste under any circumstance. Because of just copy-pasting, a lot of our students have faced difficulties in the past when they had to write some code on their interview.



IT'S OVER, IT'S DONE

THE ASSIGNMENT IS DONE.