



A REPORT ON

DATA SCIENCE Mini Project

**FARE AMOUNT PREDICTION**

SUBMITTED TO THE SAVITRIBAI PHULE PUNE UNIVERSITY, PUNE IN  
THE PARTIAL FULFILLMENT FOR THE AWARD OF

THE DEGREE OF BACHELOR OF ENGINEERING IN  
COMPUTER ENGINEERING

BY

- MISS. SANYOGITA DATTATRAY LONDHE. ROLL NO -358
- MISS.PRATIKSHA VIJAY RAUT. ROLL NO-374
- MISS.PAYAL PRAMOD BHOSALE ROLL NO-309
- MISS.POOJA KAILAS MEHER.ROLL NO-361

UNDER THE GUIDANCE OF PROF. S.N.NIMBALKAR

DEPARTMENT Of

COMPUTER ENGINEERING

# **Contents:**

- 1.Introduction
- 2.Problem Statement
- 3.Software Requirement Specification
- 4.Graphical User Interface
- 5.Conclusion

## INTRODUCTION:-

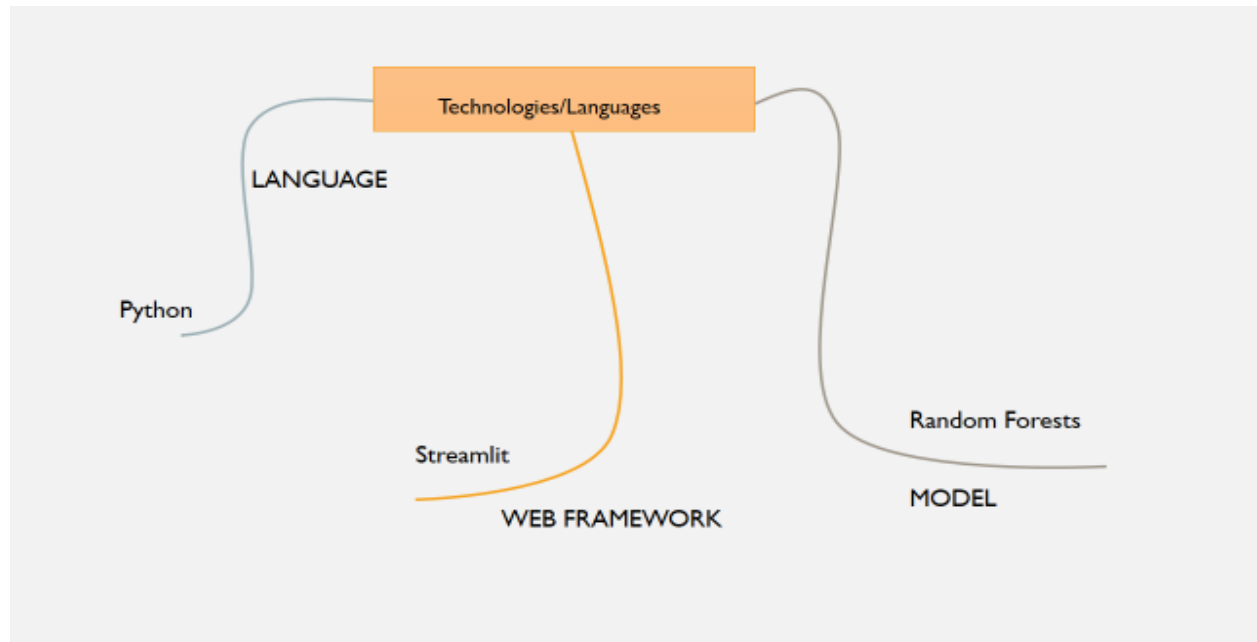
The taxi fare problem is one of several real-world problems that are used as case studies in the series of courses. Random Forest is far more flexible than a Linear Regression model. This means lower bias, and it can fit the data better. Complex models can often memorize the underlying data and hence will not generalize well. Parameter tuning is used to avoid this problem.



## **PROBLEM STATEMENT: -**

We take cab rides on a regular basis (sometimes even daily!), and yet when we're hitting that 'Book now' button, we rely on manual on-the-fly calculations rather than hardcore ML ones. And that's what I aim to demonstrate here.

# SOFTWARE REQUIREMENT SPECIFICATION:



**TAXI DATA ANALYSIS**

In this project we are predicting the fare amount

## Taxi dataset

dataset we are taken to do analysis (head of dataset).....

	VendorID	tpep_pickup_datetime	tpep_dropoff_datetime	passenger_count	trip_distance
0	1	01-01-2021 00:30	01-01-2021 00:36	1	0.5
1	1	01-01-2021 00:51	01-01-2021 00:52	1	0.1
2	1	01-01-2021 00:43	01-01-2021 01:11	1	2.6
3	1	01-01-2021 00:15	01-01-2021 00:31	0	1.6
4	2	01-01-2021 00:31	01-01-2021 00:48	1	1.7

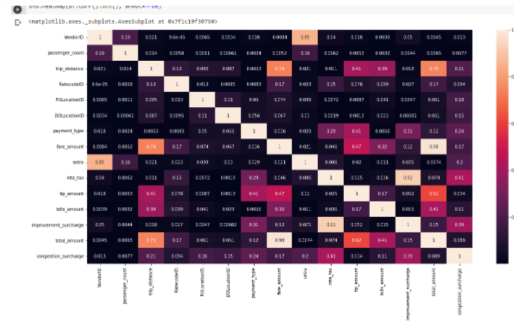
**Dealing with duplicate data**

Click is you want to check duplicate data

Click if you want to delete the duplicants rows

Windows taskbar showing search bar, taskbar icons (Edge, File Explorer, Store, Teams, Mail, WhatsApp, VS Code, Task View, Chrome, Firefox), system tray (32°C Clear, 21:35, 15-04-2022, ENG, 4 notifications).

## The features I created



Correlational matrix to check relations between variables

- **first feature:** From above correlational matrix we clearly know that our target variable is fare\_amount
- **second feature:** fare\_amount is depend on DOLocationID, payment\_type, extra, mta\_tax, tip\_amount, tolls\_amount, improvement\_surcharge, total\_amount, congestion\_surcharge

## Time to train model

Here you get to choose the hyperparameter of model and see how the performance change

## Time to train model

Here you get to choose the hyperparameter of model and see how the performance change

what should be the max\_depth of the model?



how many trees should there be?

100

here is a list of features in my data

0	VendorID
1	tpep_pickup_datetime
2	tpep_dropoff_datetime
3	passenger_count
4	trip_distance
5	RatecodeID
6	store_and_fwd_flag
7	PULocationID
8	DOLocationID
9	payment_type

which feature should be used as input feature?

PULocationID

mean absolute error of model is

27.36624393992797

mean squared error of model is

1622.1387772967344

r2 score of model is

-5.136961691562553





## **CONCLUSION:**

Random Forests significantly improved the predictive ability of our Machine Learning model. Another way to improve the model's accuracy is to increase the amount of training data, and/or building ensemble models. If we have a lot of dimensions (features) in the data, dimensionality reduction techniques can also help improve the model's accuracy.