# Document

# Of

# Mercedes-Benz Greener Manufacturing .
# Project 1

## DESCRIPTION - :

Reduce the time a Mercedes-Benz spends on the test bench.
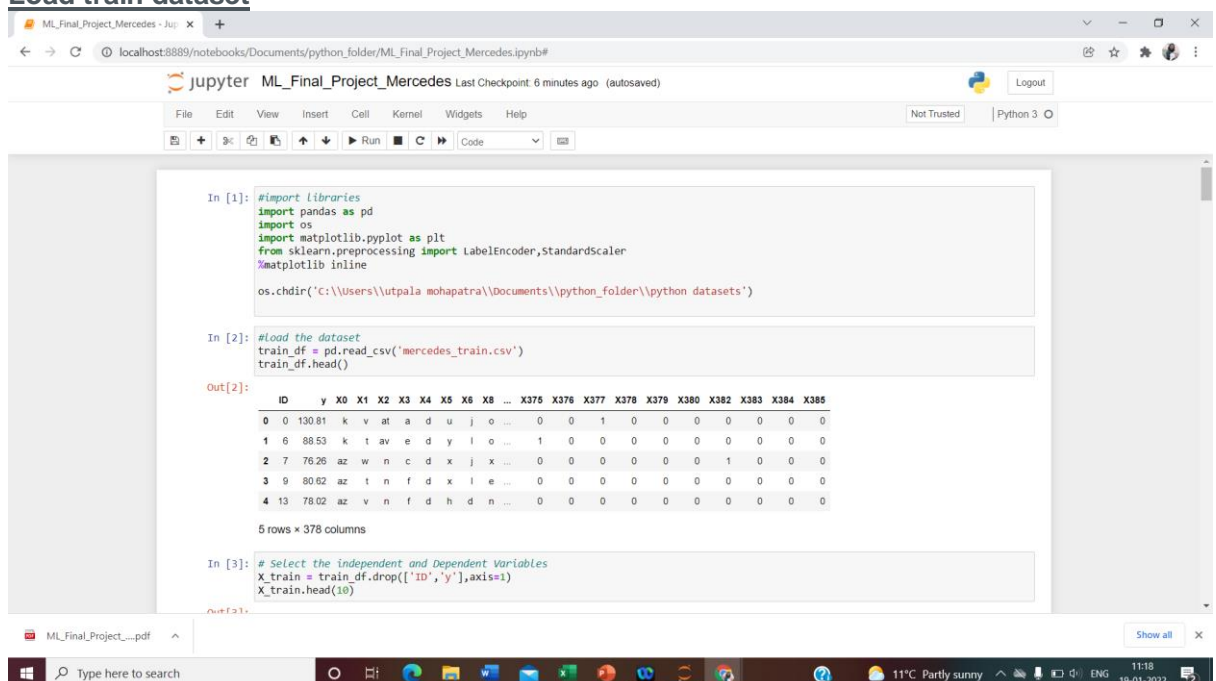
Problem Statement Scenario:
Since the first automobile, the Benz Patent Motor Car in 1886, Mercedes-Benz has stood for important automotive innovations. These include the passenger safety cell with a crumple zone, the airbag, and intelligent assistance systems. Mercedes-Benz applies for nearly 2000 patents per year, making the brand the European leader among premium carmakers. Mercedes-Benz is the leader in the premium car industry. With a huge selection of features and options, customers can choose the customized Mercedes-Benz of their dreams.

To ensure the safety and reliability of every unique car configuration before they hit the road, the company's engineers have developed a robust testing system. As one of the world's biggest manufacturers of premium cars, safety and efficiency are paramount on Mercedes-Benz's production lines. However, optimizing the speed of their testing system for many possible feature combinations is complex and time-consuming without a powerful algorithmic approach.

You are required to reduce the time that cars spend on the test bench. Others will work with a dataset representing different permutations of features in a Mercedes-Benz car to predict the time it takes to pass testing. Optimal algorithms will contribute to faster testing, resulting in lower carbon dioxide emissions without reducing Mercedes-Benz's standards.

## Tasks - :

1- **Load train dataset**

## 2- Select Dependent and independent variables



## 3- Check for Null and Unique Values

## 4- Drop the Zero variance Columns from Numeric Data



## 5- Apply LabelEncoder to the Categorical variable

## 6- Concatenate the Cleaned Numeric data and the Converted Categorical data



## 7- After Scaling Apply Dimensionality Reduction by using PCA

## 8- Apply Train Test Split



## 9- Load Test Dataset

## 10- Repeat all actions to get Clean Test Dataset



## 11- Check the RMSE and Predict the data by applying XGBoost to the Clean Test Dataset after scaling and PCA