



CAR PRICE PREDICTION

Team Presentation



Nguyen Hoang Vu
20190100



Dao Duc Manh
20194794

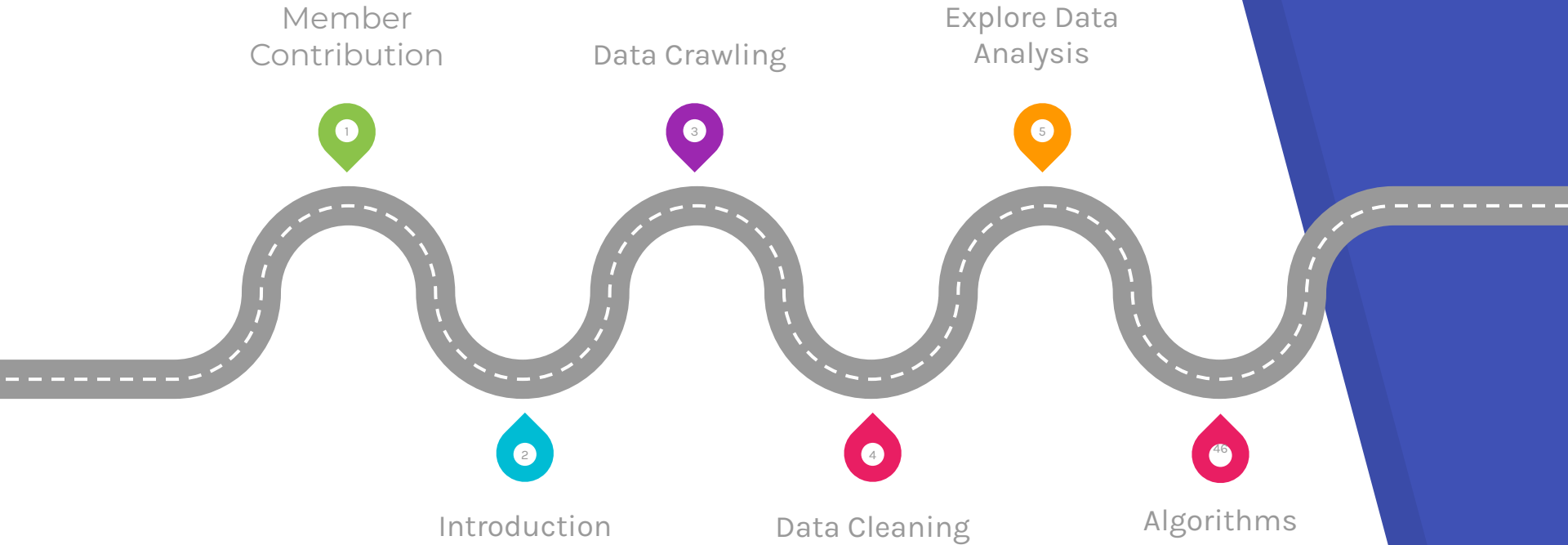


Dang Quang Minh
20194796



Vu Hoang Nam
20194809

Content



A thick yellow diagonal stripe runs from the top right corner towards the bottom left, separating the white background on the left from the solid yellow background on the right.

1.

MEMBER CONTRIBUTION



Member contribution

Member name	Task	Review by
Đào Đức Mạnh	<ul style="list-style-type: none"><input type="checkbox"/> Clean data<input type="checkbox"/> Explore data<input type="checkbox"/> Report, slide	All member
Đặng Quang Minh	<ul style="list-style-type: none"><input type="checkbox"/> Crawl data<input type="checkbox"/> Algorithms<input type="checkbox"/> Report, slide	All member
Nguyễn Hoàng Vũ	<ul style="list-style-type: none"><input type="checkbox"/> Crawl data<input type="checkbox"/> Algorithms<input type="checkbox"/> Report, slide	All member
Vũ Hoàng Nam	<ul style="list-style-type: none"><input type="checkbox"/> Clean data<input type="checkbox"/> Explore data<input type="checkbox"/> Report, slide	All member

2.

INTRODUCTION



INTRODUCTION

Approximately 40 million used vehicles are sold each year. Effective pricing strategies can help any company or individual to efficiently sell its products in a competitive market and make a profit, that is why we have chosen this topic for the data science project.

In this report, we will illustrate our process of doing the research, beginning from crawling car data to predicting car prices using some Machine Learning algorithm.

A thick yellow diagonal stripe runs from the top right towards the bottom left, separating the white background on the left from the solid yellow background on the right.

3.

**DATA
CRAWLING**



CRAWLING TOOLS



Extract data from sites with multiple levels of navigation



An open-source and collaborative framework for extracting the data



CRAWLING PROCESS

- ❖ Flow
 - Using Web Scraper, crawls all car links
 - 84, 000 links available
 - Using scrapy, extracts features from links
 - 37 features + more than 75, 000 data lines
- ❖ Difficulty
 - Number of lines of data is large
 - Takes nearly a day for crawling
 - Old links and old website

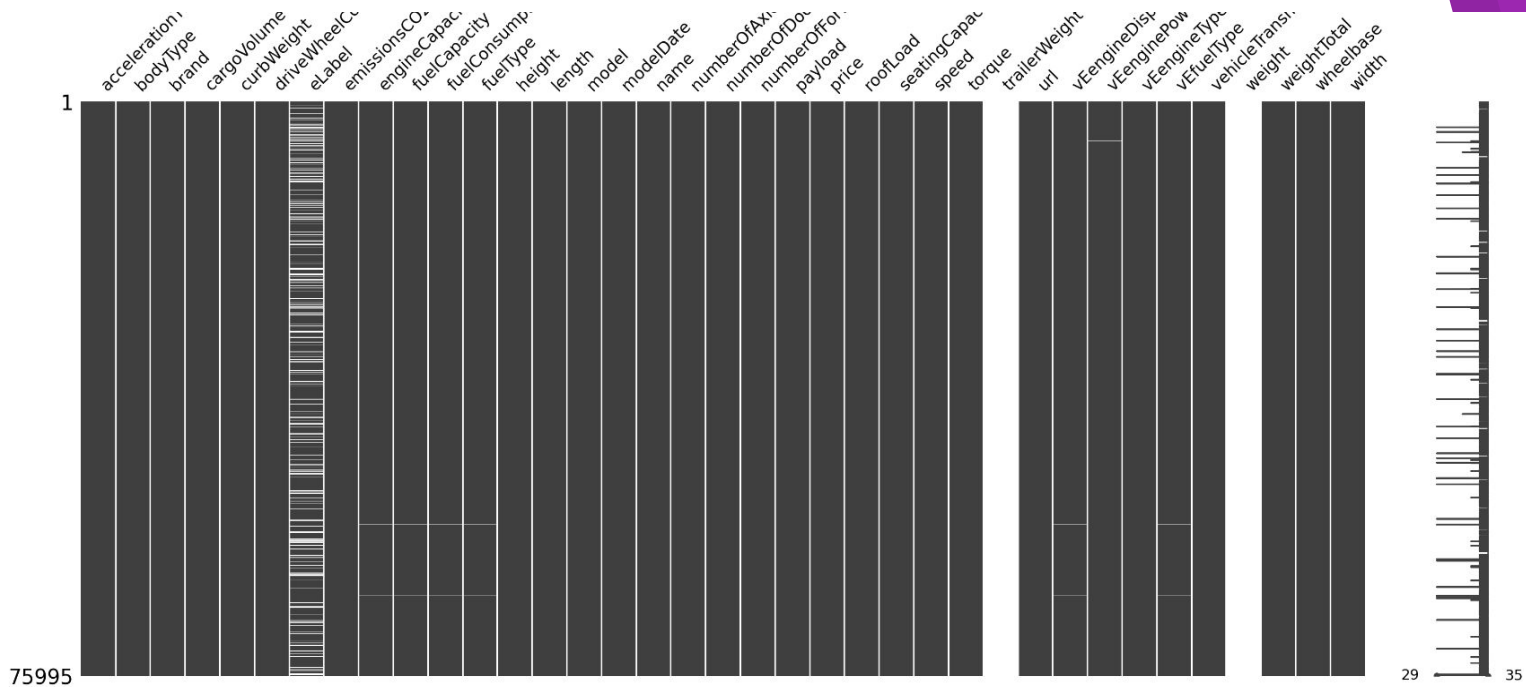


4.

Data Cleaning



Null value



METHOD HANDLING

Deleting rows

- ❖ Pros
 - Complete remove
 - Delete row/column not high value
- ❖ Cons
 - Loss info
 - High missing percentage

Mean/Median/Mode

- ❖ Pros
 - Prevent data loss
- ❖ Cons
 - Imputing variance and bias

Fill null With KNN

- ❖ Pros
 - Correlation is ignored
- ❖ Cons
 - Time-consuming
 - Choice distance not robust result



Cleaning Process

- ❖ Remove two null columns (trailerWeight, weight)
- ❖ Null data after cleaning takes $> 30\%$ of data
 - NO deleting rows
- ❖ Fill the NaN value by the most frequent value
- ❖ Use KNN in order to fill our missing data
 - 30 features and 75824 lines of data



5.

DATA

EXPLORING

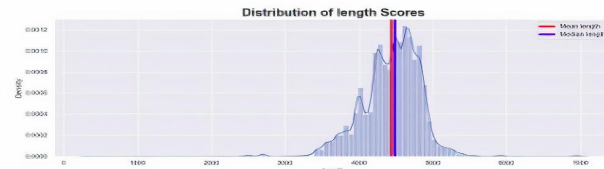
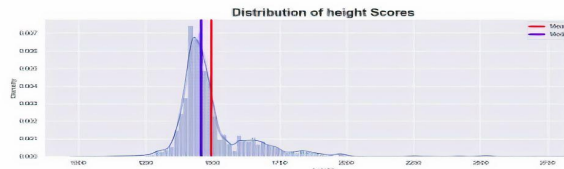
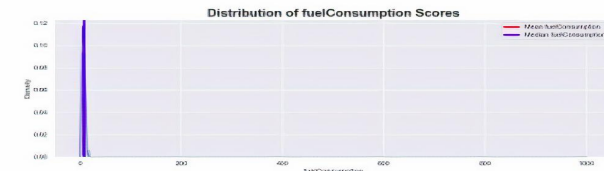
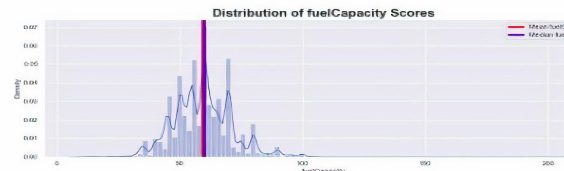
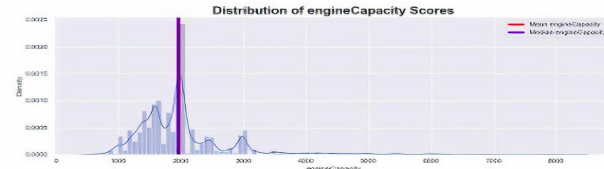
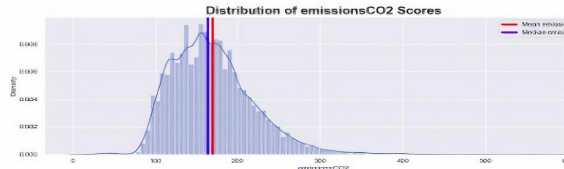
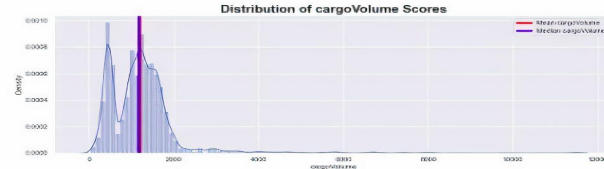
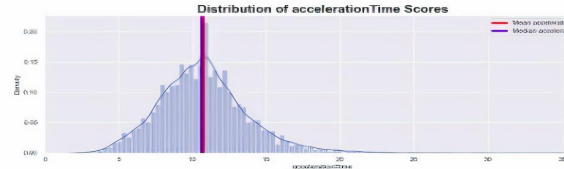


General Analysis

```
<class 'pandas.core.frame.DataFrame'>
RangeIndex: 75995 entries, 0 to 75994
Data columns (total 37 columns):
#   Column                                Non-Null Count  Dtype
---  -
0   accelerationTime                      73763 non-null  float64
1   bodyType                             75995 non-null  object
2   brand                                75995 non-null  object
3   cargoVolume                           74946 non-null  float64
4   curbWeight                            75769 non-null  float64
5   driveWheelConfiguration               75995 non-null  object
6   eLabel                                61841 non-null  object
7   emissionsCO2                          61705 non-null  float64
8   engineCapacity                       75818 non-null  float64
9   fuelCapacity                         75815 non-null  float64
10  fuelConsumption                      74862 non-null  float64
11  fuelType                              75858 non-null  object
12  height                               75925 non-null  float64
13  length                               75974 non-null  float64
14  model                                75995 non-null  object
15  modelDate                            75995 non-null  int64
16  name                                 75995 non-null  object
17  numberOfAxles                        75995 non-null  int64
18  numberOfDoors                        75995 non-null  int64
19  numberOfForwardGears                 74026 non-null  float64
20  payload                              74976 non-null  float64
21  price                                75824 non-null  float64
22  roofLoad                             64773 non-null  float64
23  seatingCapacity                      74026 non-null  float64
24  speed                                75215 non-null  float64
25  torque                               75888 non-null  float64
26  trailerWeight                        0 non-null     float64
27  url                                  75995 non-null  object
28  vEngineDisplacement                  75818 non-null  float64
29  vEnginePower                         75946 non-null  float64
30  vEngineType                          75995 non-null  object
31  vFuelType                            75858 non-null  object
32  vehicleTransmission                 75995 non-null  object
33  weight                               0 non-null     float64
34  weightTotal                          74983 non-null  float64
35  wheelbase                            75969 non-null  float64
36  width                               75970 non-null  float64
dtypes: float64(23), int64(3), object(11)
memory usage: 21.5+ MB
```

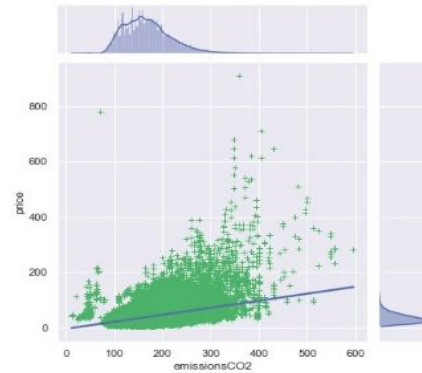
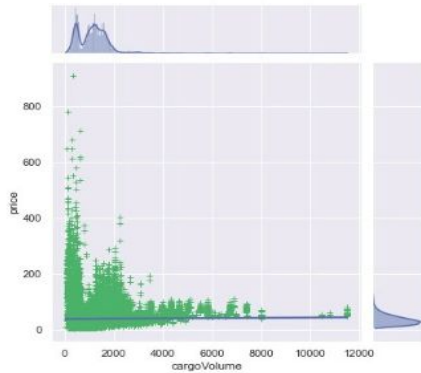
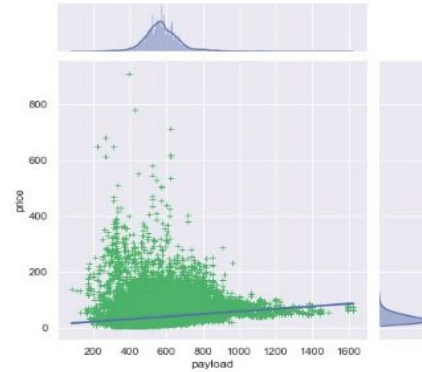
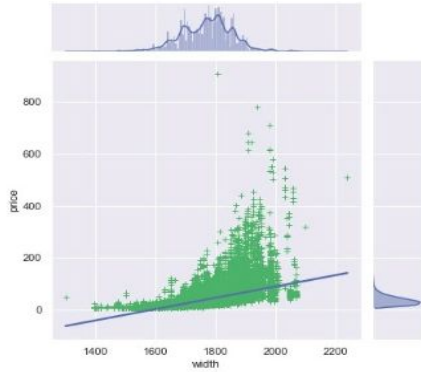

Numerical Data

❖ Explore Skewness



Numerical Data

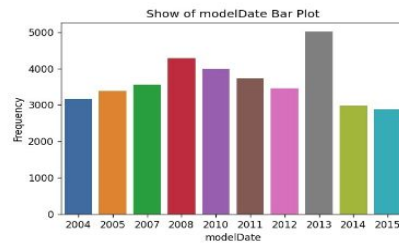
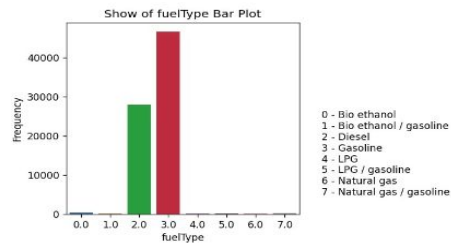
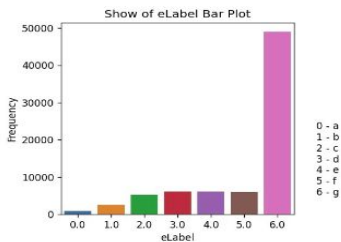
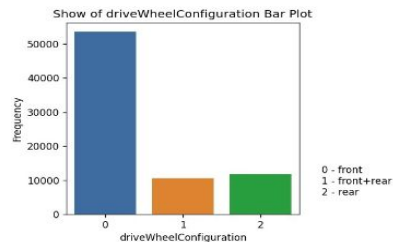
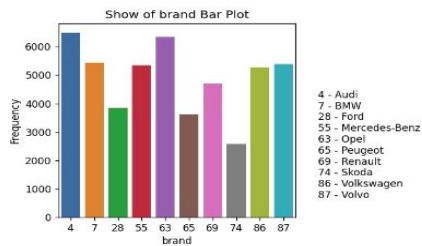
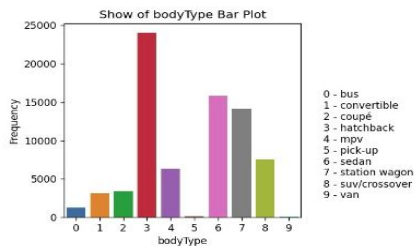
❖ Relationship with Target





Categorical Data

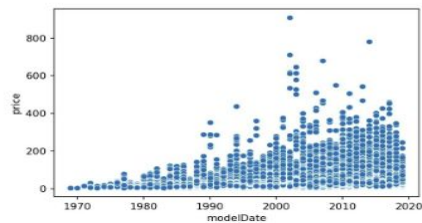
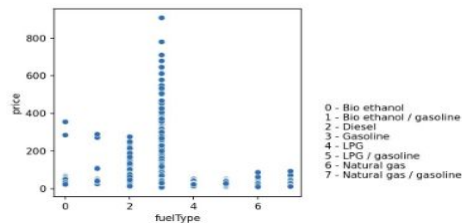
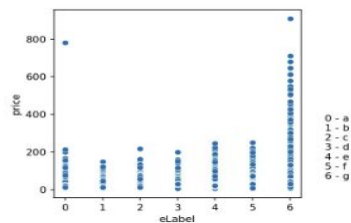
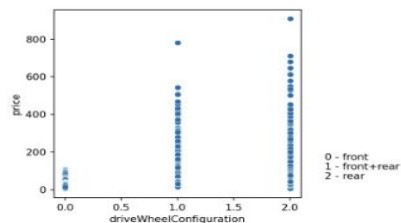
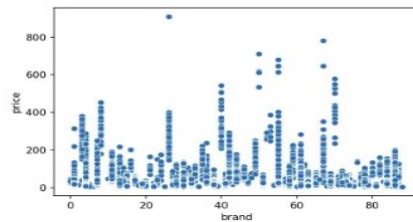
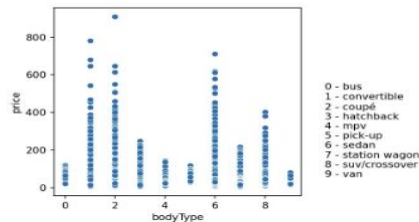
► Barplot





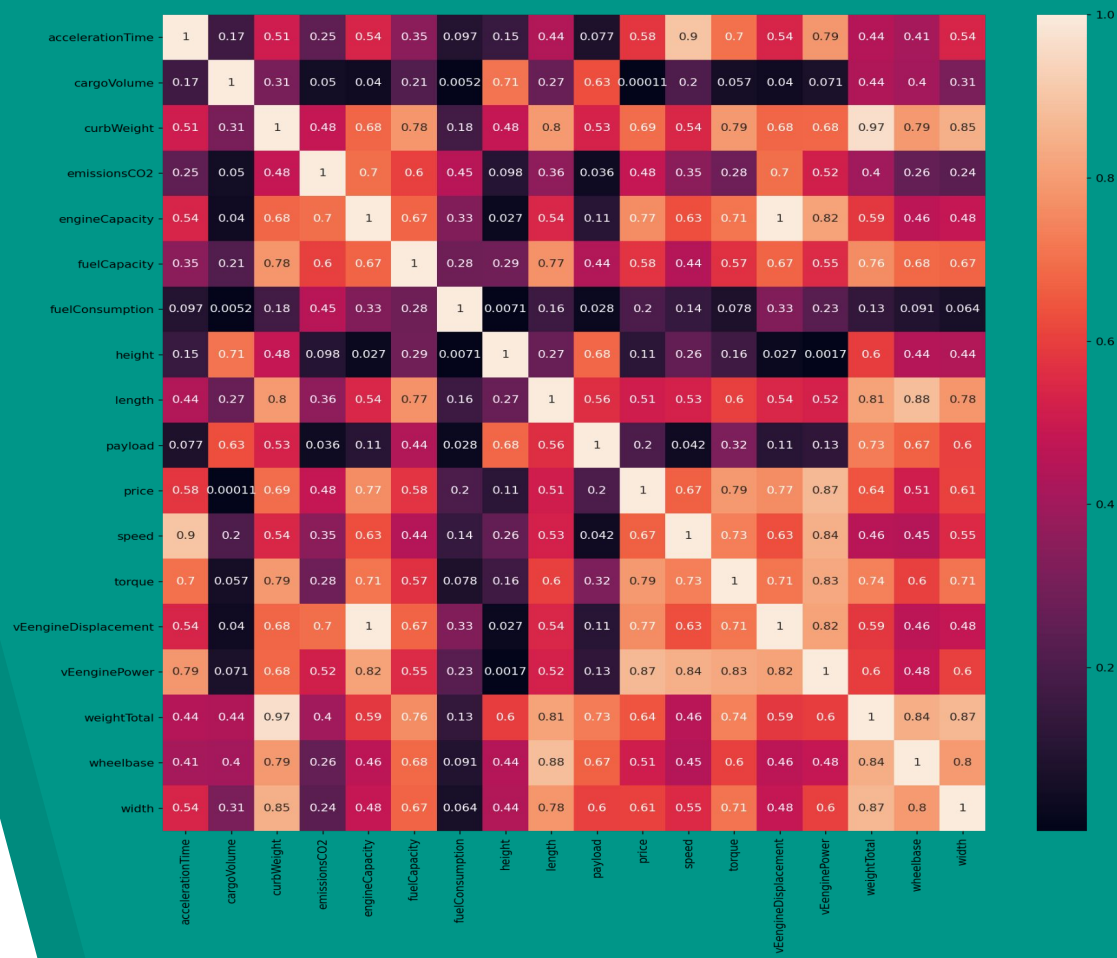
Categorical Data

➤ Scatter plot



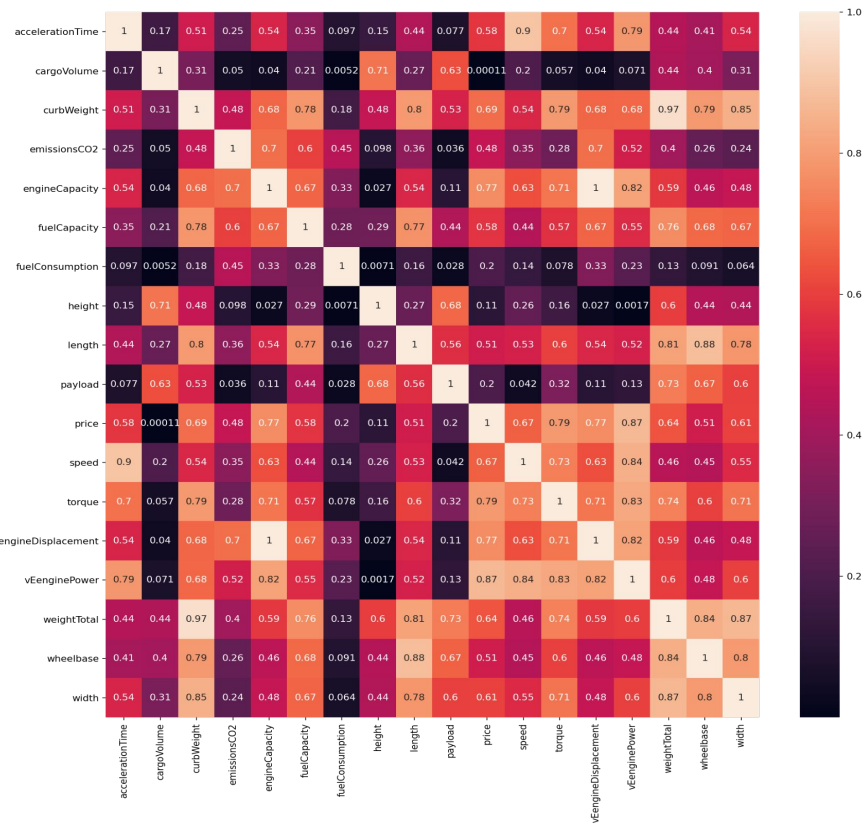


Heatmap





High correlation filter



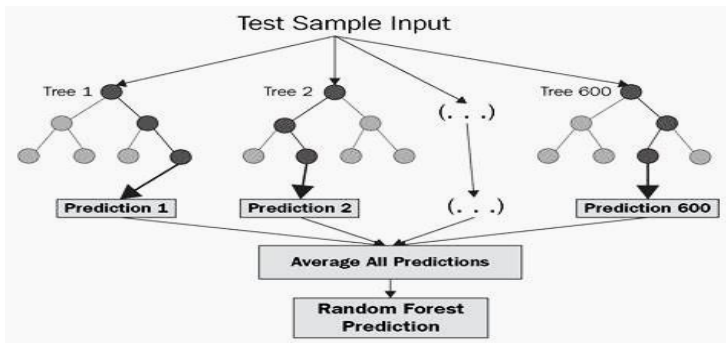
- ❑ If the correlation values between some features are big, we can consider them as the same in the dataset.
- ❑ Those similar features may make the dataset harder for the model to learn.

6.

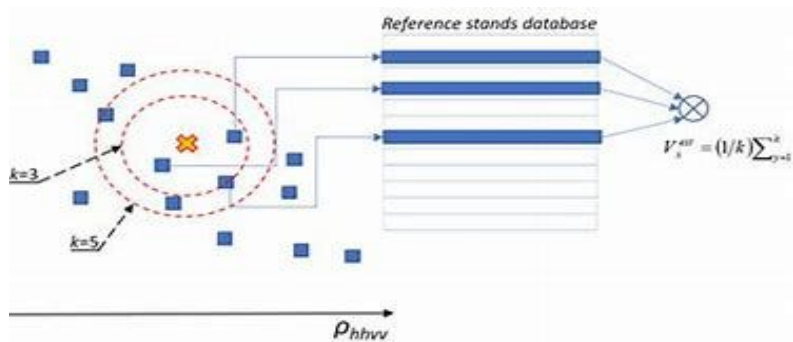
ALGORITHMS



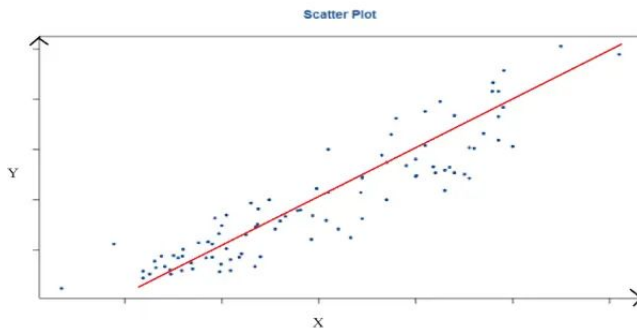
Algorithms



Random Forest



KNN



Linear Regression



Algorithms

- ❑ Besides Random Forest, KNN, Linear Regression, we also try some Boosting model:
 - ❑ Cat Boosting
 - ❑ XGB Boosting
 - ❑ LGBM Boosting

- ❑ To our knowledge, the basic idea of the Boosting algorithms is to use the combination of multiple weak models. Then, training weak learners sequentially, each will try to correct its previous model.



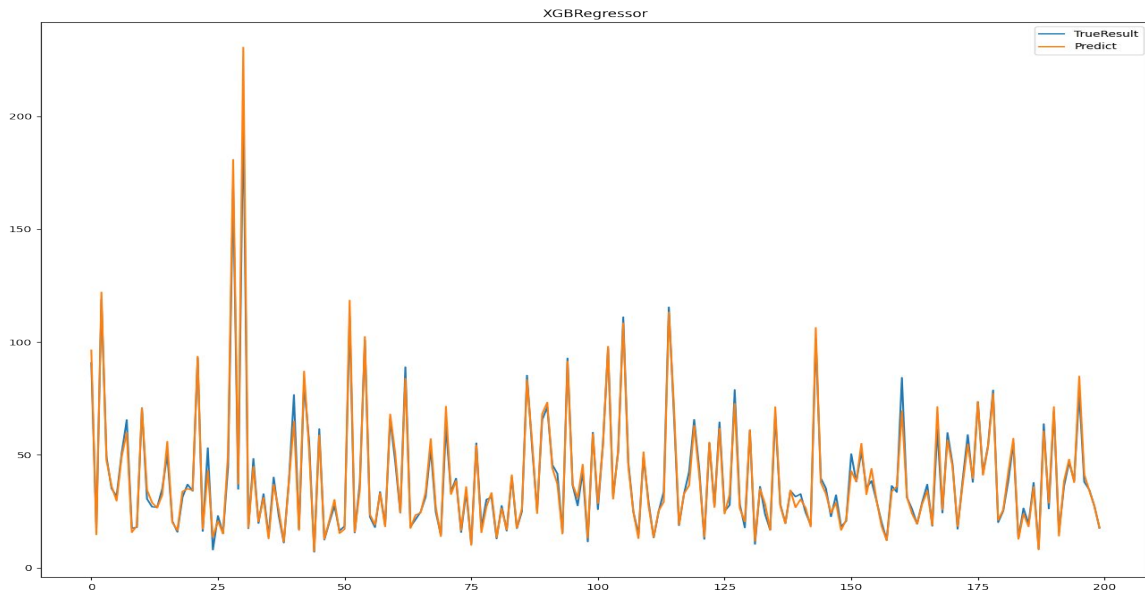
Results

Algorithm	RMSE loss
LinearRegression	11.31
KNN Regressor	7.90
CatBoostRegressor	7.19
RandomForestRegressor	6.19
XGBRegressor	6.05
LGBMRegressor	6.01

- ❑ We split the dataset into 2 parts:
 - ❑ Train set 70%
 - ❑ Test set 30%
- ❑ We use StandardScaler for normalizing data.

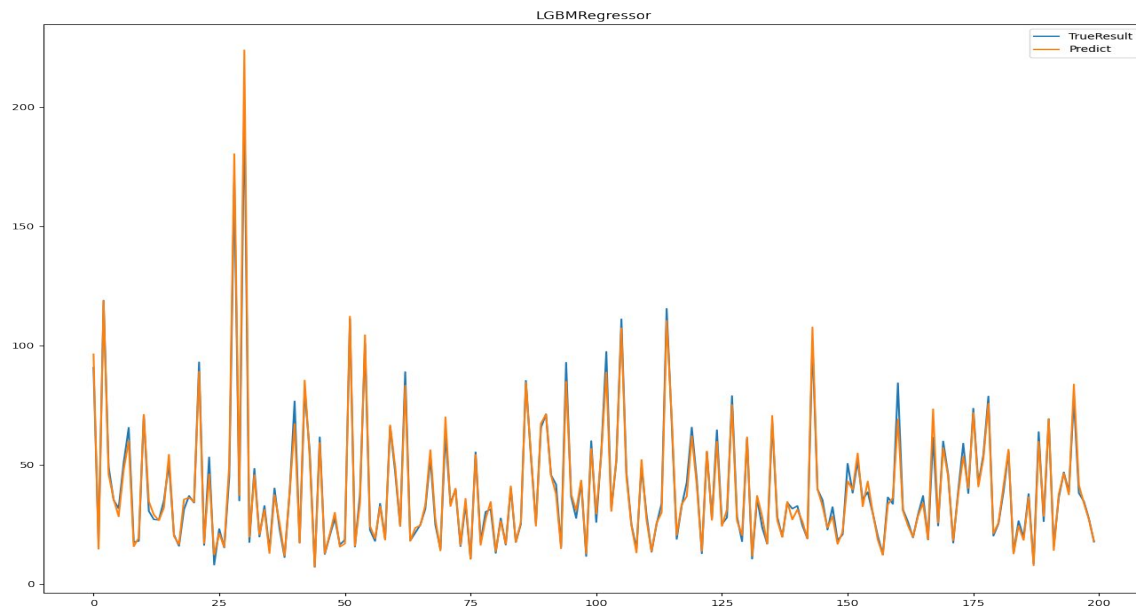


Two best models





Two best models



THANKS FOR LISTENING!