

Used Car Price Prediction Using Cat Boost Gradient Machine and Light Gradient Boosting Machine

¹Sai Guptha Grandhi, ²Bangaru Suresh, ³Chittajallu Nagarjuna,

Junior Data Scientist, Social Tek AI & ML Solutions, Hyderabad

Junior Data Scientist, Social Tek AI & ML Solutions, Hyderabad

Junior Data Scientist, Social Tek AI & ML Solutions, Hyderabad

Submitted: 01-08-2021

Revised: 14-08-2021

Accepted: 17-08-2021

ABSTRACT: On going Recession in the world, the automobile is one of the industries that suffers a lot in the production of cars because the availability of technology, designing and manufacturing of the parts are not available in the market for best price. So, people in the India shows inclined towards used cars which have innumerable benefits. A used car prediction has been a high meticulous research area, as it requires notable efforts and knowledge of the field expert. Significant number of various attributes are examined for the reliable and accurate prediction. To build a model for predicting the worth of used cars in India. We applied seven machine learning techniques (Multiple linear Regression, Random Forest, Gradient Boost Machine, light GBM, Cat Boost regression, Ada Boost Regression and XGB). However, the mentioned algorithms were applied to work as an individual, ensemble and stacking. The data used for the prediction was collected from the web portal Car24.com using web scraper that was written in Python programming language with Selenium. Performance of different algorithms were then compared to find one that best suits the scrapped data set. The final prediction model was integrated into Python application. Furthermore, the model was evaluated using test data and the comparing of different evaluation metrics like coefficient of determination (r^2) and root mean square error (RME) the lowest MAE and R^2 obtained with Extreme Gradient Boosting train r^2 with 0.986 and test r^2 with 0.944. In terms of RMSE the train values are with 0.118 and test 0.120.

KEYWORDS: used car price prediction, regression, machine learning, linear regression, cat boost regression, random forest, GBDT, light GBM and GBM.

I. INTRODUCTION:

The field of Artificial Intelligence has allowed analysts to cover insights from historical data and past events to get the inference and to predict the future insights. One of the most famous case study used car price prediction. Over, rapid growing population in India, transportation plays a prominent role in which car is one kind of transportation, according to study conducted by Indian Blue Book, in FY 20 the used car market in India stands at 4.2 million units. It also claimed that used car sales increased by 5% YoY within the in the same fiscal, while the new car sales contracted by 17.8%, due to the pre-pandemic depression. The used car sales were actually 50% above the amount of latest car sales recorded at 2.8 million units.

Accurate car price prediction involves expert knowledge, because price usually depends on various attributes and factors. Typically, most significant ones are brand and model, age, Kilometers travelled, Insurance, Type of Transmission, fuel type used and ration of the car. The dataset collected using web scrapping with selenium from car24.com, the purpose is about price prediction of the car. We split the historical data set that we obtained from car24.com into three groups – train set, validation set and test set. The train data set are going to be input for the training model. The validation data set will be an to check model for the generations of predictions. This last data set will be used to check the correctness for the generations of predictions for the test data. There can be many flaws in the train set because of which training of the model on the basis of such data is not possible or can leads to an incorrect output.

This dataset has been studied and analyzed using various machine learning algorithms like Multiple linear regression. various languages and

tools are used to implement these algorithms including Python. The approach of the research paper is centered on Python for executing algorithms-, Random Forest, GBDT, light GBM, Non-Linear regression and XGB.

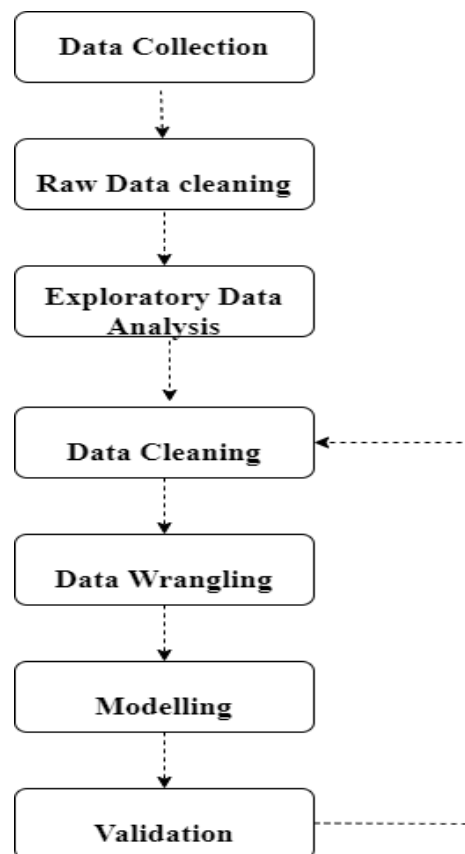
II. RELATED WORK:

Predicting price of a used cars has been studied extensively in various researches. Listian discussed, in her paper about regression model that was built using Support Vector Machines

(SVM) can predict the worth of a car that has been leased with better precision than multivariate regression or some simple multiple correlation. This is often on the grounds that Support Vector Machine (SVM) is best in handling datasets with more dimensions and it's less susceptible to overfitting and underfitting. The weakness of this research is that a change of straightforward regression with more advanced SVM regression wasn't shown in basic indicators like mean, variance or variance.

III. EXPERIMENTAL METHODS:

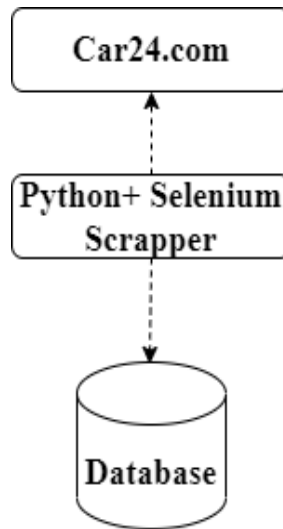
Approach for used car price prediction proposed in this paper is composed of several steps,



IV. DATA COLLECTION:

Data is collected from a local web portal for selling and buying cars car24.com, during post Pandemic, as a time interval itself has high impact on the price of the cars in India. The following attributes were captured for each car: brand, model, years of car

purchased, kilometers travelled, fuel, year of manufacturing, transmission type, city, since manual data collection is hectic and time-consuming task, especially when there are numerous records to process, a “web scraper” as a part of this research is created to get this job done automatically and reduce the time for data gathering.



Web scraping is the process of using automatic bots to extract content and data from a website and save data into required database.

Manual data extraction is time consuming and therefore web scrapers are used to do this job in a small amount of time.

RAW DATA CLEANING:

	Name	Price	Rating	city	Kilometers	Year_of_Purchase	Owner	Fuel_Type	Transmission	#TO	Insurance	Insurance_Type
0	2010 Mahindra LX	₹ 174,509	4.0 out of 5	3686	Kilometers34,854 km	Year of PurchaseMay 2010	OwnerFirst Owner	FuelPetrol	TransmissionMANUAL	RTOAP09	NaN	Insurance TypeInsurance Expect
1	2013 Mahindra Wagon R 1.0 VXX	₹ 333,989	4.2 out of 5	3686	Kilometers39,541 km	Year of PurchaseJuly 2013	OwnerSecond Owner	FuelPetrol	TransmissionMANUAL	RTOAP26	Insurance5/6/2021	Insurance TypeComp
2	2014 Mahindra Wagon R 1.0 VXX	₹ 353,199	4.3 out of 5	3686	Kilometers23,233 km	Year of PurchaseMarch 2014	OwnerSecond Owner	FuelPetrol	TransmissionMANUAL	RTOAP26	Insurance28/4/2021	Insurance TypeComp
3	2013 Hyundai Eon D LTE Plus	₹ 237,889	4.4 out of 5	3686	Kilometers27,740 km	Year of PurchaseNovember 2013	OwnerFirst Owner	FuelPetrol	TransmissionMANUAL	RTOAP13	Insurance1/1/2022	Insurance TypeComp
4	2017 Hyundai Eon ERA PLUS	₹ 300,889	4.4 out of 5	3686	Kilometers12,236 km	Year of PurchaseAugust 2017	OwnerSecond Owner	FuelPetrol	TransmissionMANUAL	RTOATS10	NaN	Insurance TypeExpired

After raw data has been collected and stored to local database, The primary Raw data contains 32158 samples with 12 features attributes like Name, Price, Rating, City, Kilometers, year of purchase, owner, fuel type, rating of the car, RTO, insurance and insurance type. Raw data cleaning step was applied using different function and different techniques like

split function, loops concepts and raw data converted into useful data which attributes contains Name of the model, Company of the car, year and month of purchased, kilometers travelled, owner type, transmission, fuel type, city, RTO insurance data and type. After the secondary staged data it looks like shown in a figure

	Price	Rating	city	Kilometers	Owner	Fuel_Type	Transmission	RTO	Insurance	Insurance_Type	year	month	Company	Model
0	174699	4	3686	34854	First	Petrol	MANUAL	RTOAP	NaN	Expired	2010	5	maruti	alto bi
1	333999	4.2	3686	39541	Second	Petrol	MANUAL	RTOAP	6/8/2021	Comp	2013	7	maruti	wagon r 1.0 vxi

V. EXPLORATORY DATA ANALYSIS:

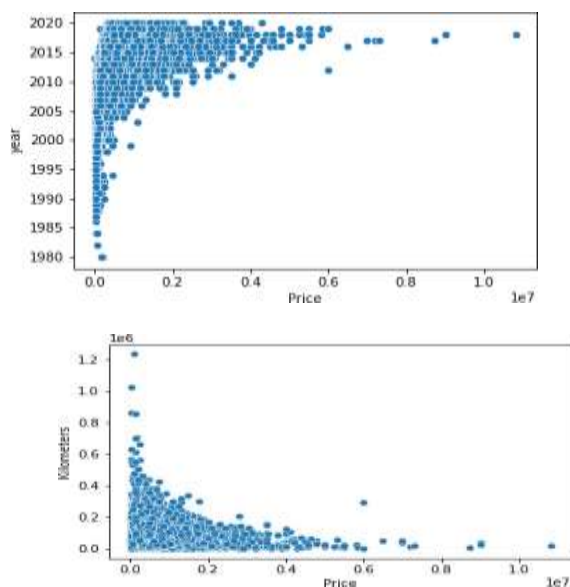
After raw data cleaning, the next foremost step is Exploratory Data Analysis (EDA) to get the data understanding of each column and how the column is useful for predicted the model. On a careful analysis, many of the features were sparse and they do not contain useful information for prediction.

Hence, it is decided to remove them from the dataset.

The attributes RTO and city were completely

removed. After the EDA data set contains 32158 samples with 12 features. Next, with the data set were performed different statistical operations like Mean, Median, Mode, Count, frequency, maximum, minimum, Quartile 1, Quartile 2, Quartile 3, Quartile 4, Skewness and Kurtosis and standard deviation.

Visualization: To check whether the distribution follows any linearity or not



The above graph shows that the main attributes are not following any linearity between them.

VI. DATA CLEANING:

After EDA our data is imperfectly manipulated now, in this data cleaning process we were mostly deals with Missing value and outliers treatment. On a careful data analysis, we find that some of the attributes are like Transmission, Insurance, kilometers and type of insurance have missing values. Transmission is a categorical variable so, we were used different techniques like Hot Deck Imputation, Simple imputer and Iterative Imputer.

Hot Deck Imputation: It is a method for handling missing values in which each missing data is replaced with an observed response from a "similar" unit. It involves replacing missing values of attributes with observed values from a respondent

(the donor) that is similar to the non-respondent with respect to characteristics observed by both cases.

Simple Imputer: Simple Imputer is a scikit-learn class. It is one of the most useful technics to handling the missing data in the predictive model dataset. It replaces the Nan values with a specified placeholder. It is implemented by the use of the Simple Imputer () method which takes the arguments like missing values, fill value and strategy.

Iterative Imputer: In this process where each feature is modelled as a function of the other features, where missing values are predicted using different regressor like Kenn, Extra tree, Bayes ridge and Decision tree. Each feature is imputed

sequentially, one by one, allowing prior imputed values to be used as part of a model in predicting subsequent features. This iterative process is repeated multiple times, allowing ever improved estimates of missing values to be calculated as missing values across all features are estimated. This is also known as fully conditional specification (FCS) or multivariate imputation by chained equations (MICE).

For categorical feature like Transmission, insurance and type of insurance we were used Simple imputer with mode and continuous variable like kilometers we were used simple imputer with mean, Hot deck imputation and iterative imputer compared the results with skewness. Overall, by the analysis Extra Regressor gives us a wonderful skewness with—

After solving the puzzle of Missing values, we here with biggest influencing factor is outliers. Clearing outliers is the heart for feature Engineering because this can influence the distribution and effect the model. There are different techniques to detect and handle outliers. To detect the outliers we were used Boxplot, Inter quartile range, Grubb's test, Isolation Forest and scatter plots. Treating outliers is the most important aspect in feature engineering we were used Fissurization, Arbitrary Outlier Capper and rectify using Inter quartile range.

Treatment of Outliers:

Fissurization: It is a method of averaging that initially replaces the smallest and largest values with the observations closest to them. It changes distribution to gaussian distribution.

Arbitrary Outlier Capper: Using Arbitrary capping methods we can cap the distribution with maximum and minimum values to get better output.

Setting Quartile range: Dealing missing values with quartile range gives a major impact of the outliers. Outliers will remove completely and replaced with maximum and minimum values. They were different kinds of ratio were given to remove the outliers they were 90-10, 0.75-0.25.

Dropping of Outliers: we check the percentage of outliers and drop the outliers to get maximum benefit.

We used above all methods and get the maximum output with setting quartile ranges. After the Data cleaning Process now, we don't have any outliers and missing values in the data now, we can go for

Detection Methods: Using Z- score or Extreme value analysis to detect the unanimous behavior in the dataset and this indicates how many standard deviations a data point is from the sample's mean, assuming a gaussian distribution. Using Z score I detected that in Kilometers travelled and rating attributes have an outlier.

Box-plot: Using Boxplot we were used to detect the skewed of the distribution and to check were the outliers present in the distribution.

Inter- quartile range: By calculating the inter quartile range of attributes in the data set I find that some of the data points are highly influencing for the dataset.

Grubb's test: Using Grubb's test and p value to find the alternative hypothesis and null hypothesis.

Isolation forest: The principle beyond the Isolation Forest algorithm outliers is few and far from the observations.

Using all the above methods, we were checked and come to conclusion that they are two outliers' attributes present in the dataset which is kilometers travelled and rating given.

data wrangling.

Data Wrangling: Basically, Data wrangling is divided into Encoding, Binning, train-test split, scaling and transformation.

Transformation: we were checked every attribute of the dataset and distribution of the transformation and we comes to know categorical variable are following Poison and Binomial distribution and continuous variable having the skewness greater than 1 and less than -1. So, we given log normal distribution to get skewness greater than 1 and less than -1.

Train-Test spilt: Using train test spilt we divide the data into 2 parts with 80% of the train and 20% of the test data. Our total data consists of 32158 records. In which 80% means 25726 records for training and 6430 records for testing.

Encoding: For both the train and test data we perform different types of methods to convert the data from categorical to numeric so, machine learning algorithms can easy to proceed. We used

One hot encoder, dummy variable trap, label encoder for the data. This all are one type of method and another method we just keep the categorical feature as same because it is useful for cat boosting algorithm.

One hot Encoding: This method is primary use for nominal data in which there is no comparing between the columns. It is a representation of categorical variables as binary vectors. This first requires that the categorical values be mapped to integer values. Then, each integer value is represented as a binary vector that is all zero values except the index of the integer, which is marked with a 1.

Dummy Variable: It is also known as indicator variable. This method feeds the data to machine to understand the behavior of categorical variable and gives the input in the values 1 or 2.

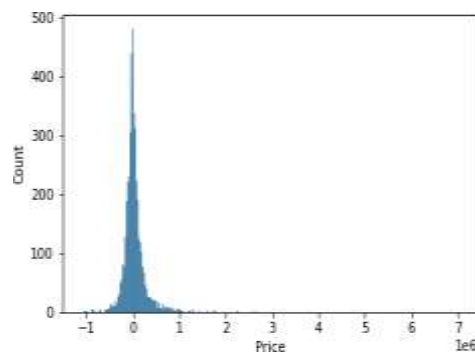
Label Encoder: It is useful for to encode the ordinal data to feed for machine. Using Label

encoder, I divide the attributes like company and model. The attributes like Transmission, Owner and fuel type convert using one hot encoding and dummy variable.

Feature Transformation: Attributes like Ratings and Kilometers are not following standard Normal distribution So; we can convert This attribute using log Normal transformation and Box cox transformation.

Modelling: We used different Algorithms like Multiple linear Regression, Random Forest, Gradient Boosting Machine (GBM) Regressor, light Gradient Boosting Regressor, Extreme Gradient Boosting Regressor and cat boost Regressor and Adaptive Boosting Regressor

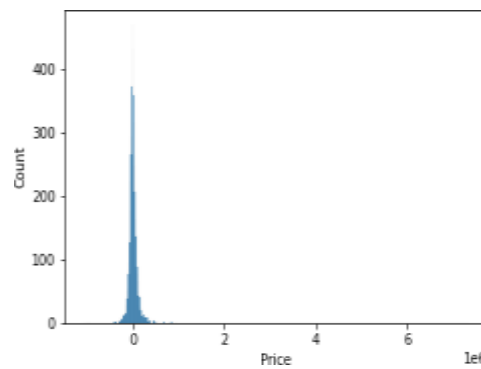
1. Multiple Linear Regression: Multiple Linear Regression was chosen as the first model due to its simplicity and comparatively small training time. We haven't used any regularization since results clearly showed low variance.



The above histogram shows the errors of the predicted variable. This clearly states that it is right skewed distribution. The more errors are closed to Zero and minimum of error in the range between -1 to 1.

Random Forest Regressor: Random Forest Regressor is an ensemble learning based on decision tree model. Specifically, this Random Forest is designed by Ho to reduce the overfitting

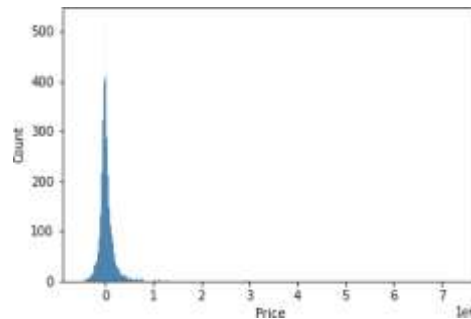
problem of Decision Tree. A random forest algorithm consists of many decision trees. The ‘forest’ generated by the random forest algorithm is trained through bagging or bootstrap aggregating. Bagging is an ensemble meta-algorithm that improves the accuracy of machine learning algorithms. This model was hence chosen to account for the large number of features in the dataset and compare a bagging technique with the following gradient boosting methods.



The above histogram exhibits the predictive variable errors. This clearly states that the distribution is Right Skewed. Most errors are closed to Zero with minimal error are less than 0.25.

Gradient Boosting Regressor Machine: Gradient Boosting is another type of decision tree-based method that is generally described as “a method of transforming weak learners into strong learners”. This means that like a typical boosting method, observations are assigned different weights and based on certain metrics, the weights of

difficult to predict observations are increased and then fed into another tree to be trained.; it allows for the optimization of arbitrary differentiable loss functions. In each stage a regression tree is fit on the negative gradient of the given loss function. This, allows for the optimization of arbitrary differentiable loss functions. In each stage a regression tree is fit on the negative gradient of the given loss function. This model was chosen to account for non-linear relationships between the features and predicted price.



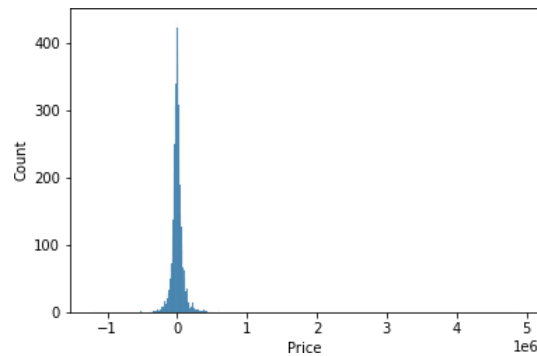
The above histogram exposes the predictive variable errors. Most errors are closed to Zero.

Extreme Gradient Boosting (XGBoost): XGBoost is quite similar to the gradient boosting

algorithm but features many additive features that significantly improve its performance such as built-in support for regularization, parallel processing as well as giving additional hyperparameters to tune

such as maxdepth, subsample and sampling method. A maximum depth of 16 was used and the

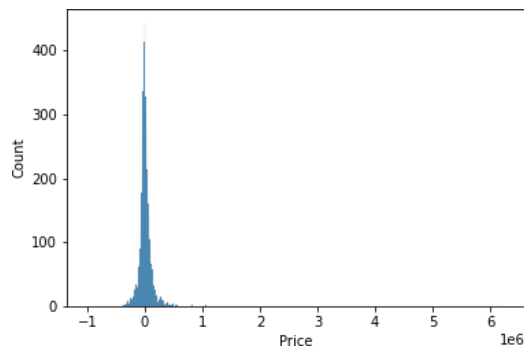
algorithm was run on all cores in parallel.



The above histogram shows predictive variable errors. This clearly states that Peakness (Kurtosis) is higher and it is right skewed. Most errors are closed to Zero. This is best Compared to

boosting based framework which is gaining popularity due it higher speed and accuracy compared to XGBoost or the original gradient boosting method. This Light GBM has a leaf-wise tree growth instead of a level-wise approach resulting in higher loss reduction.

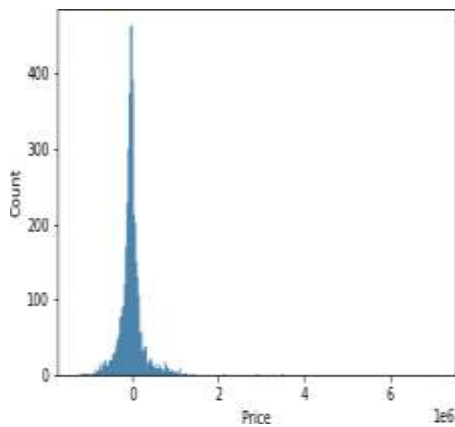
Light GBM : Light GBM is another gradient



The above histogram displays the predict variable error. It is clearly Visible that the range of errors is less than 0.5 to 0.5

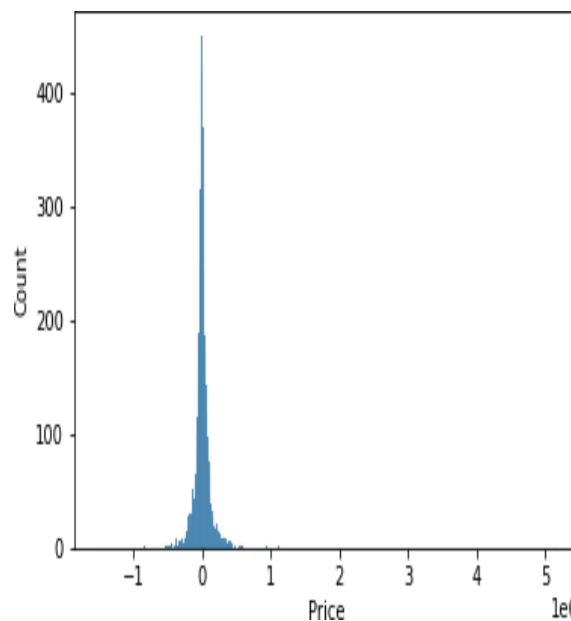
to reduce bias as well as the variance for supervised learning. It works on the principle where learners are grown sequentially. Except for the first, each subsequent learner is grown from previously grown learners. In simple words, weak learners are converted into strong ones.

Adaptive Boosting Machine: It is called Adaptive Boosting as the weights are re-assigned to each instance, with higher weights to incorrectly classified instances. Boosting is used



The above histogram demonstrates the predict variable error. It is clearly Visible that the range of errors are very minimal
 The overhead histogram describes the predict variable error. It is clearly Visible that the range of errors is minimal to zero.

Cat Boost Regressor: Cat Boost supports various types of data categories like numerical, categorical and text. The special feature with this categorical data will handle directly without using encoding techniques and having many features to handle. In Cat Boost “Boosting” is refers to the **gradientboosting machine learning**



The above histogram displays the predict variable and it error of the predicted variable. It is Visible that the peakness is achieved at 400 and the range of error is between -1 to 1. The special about this will deal with categorical variables as it is and not required any kind of change like encoding.

Coefficient of Determination (R^2): The coefficientdetermination is one of the evolution metrics for Regression that examines how differences in one variable can be explained by the difference in a second variable, when predicting the outcome of a given event. This can be denoted as R^2 .

Evaluation Metrics:

$$R^2 = 1 - \frac{SS_{res}}{SS_{tot}}$$

Where,
 SS_{res} is the sum of squares of the residual errors.
 SS_{tot} is the total sum of the errors.

Root Mean Squared Error: The most commonly used metric for regression tasks is RMSE (root-

Here, y_i denotes the true score for the i th data point, and \hat{y}_i Denotes the predicted value. One intuitive way to understand this formula is that it is the mean-square error), also known as RMSD (root-mean-square deviation). This is defined as the square root of the average squared distance between the actual score and the predicted score.

$$RMSE = \sqrt{\frac{\sum_i (y_i - \hat{y}_i)^2}{n}}$$

Euclidean distance between the vector of the true scores and the vector of the predicted scores, averaged by n , where n is the number of data points.

VII. RESULTS:

The results of our test are as follows

Name	Train(R2)	Test(R2)	Train (RMSE)	Test (RMSE)	Remarks
Linear Regression	0.719	0.727	0.449	0.444	Low R2
Random Forest	0.987	0.916	0.096	0.243	Over Fitted Model
Gradient Boosting	0.867	0.879	0.308	0.295	Moderate Model
Extreme Gradient Boosting	0.980	0.944	0.118	0.120	Good model
Light Gradient Boosting	0.923	0.925	0.234	0.232	Good Model
Cat Boost	0.948	0.922	0.192	0.192	Good Model
Ada Boost	0.660	0.655	0.502	0.506	Low R2

VIII. CONCLUSION:

The results shows that Compared to Linear Regression, most Boosting methods performed comparably well. We get high variance with the

Extreme Gradient Boosting Machine and very slight less with Light gradient boost, Gradient boosting and cat boost machine. The extreme gradient gives train r2 score with 0.980 and with test 0.944 which good and optimized model for predicting thesecond carprice.

FUTURE WORK: For better performance, we plan to judiciously design deep learning network structures with Artificial Neural Networks and train on clusters of data rather than the whole dataset and to use different procedures to select the foremost features for dataset

REFERENCES:

- [1]. Listiani, M. (2009). Support vector regression analysis for price prediction in a car leasing application (Doctoral dissertation, Master thesis, TU Hamburg-Harburg)
- [2]. https://www.temjournal.com/content/81/TEMJournalFebruary2019_113_118-JournalonCarPredictionusingClassification.
- [3]. sklearn.ensemble.RandomForestRegressor— Retrieved from: <https://scikit-learn.org/stable/modules/generated/sklearn.ensemble.RandomForestRegressor.html>
- [4]. Catboost regressor- https://catboost.ai/docs/concepts/python-reference_catboostregressor.html
- [5]. car24.com: [https://www.cars24.com-Data-retrieved fromExtraGradientBoosting-https://xgboost.readthedocs.io/en/latest/](https://www.cars24.com-Data-retrieved-fromExtraGradientBoosting-https://xgboost.readthedocs.io/en/latest/)
- [6]. Aprojectreportoncar24.com http://cs229.stanford.edu/proj2019aut/data/assignment_308832_raw/26612934
- [7]. Pudaruth, S. (2014). Predicting the price of used cars using machine learning techniques. Int. J. Inf. Comput. Technol, 4(7), 753-764.