

# OF EMISSION STANDARDS ON USED CAR PRICES IN INDIA: A MACHINE LEARNING ANALYSIS

I  
M  
P  
A  
C  
T

## Abstract

*This study delves into the complex relationship between emission standards and second-hand car valuations in India. Leveraging a comprehensive dataset from CarDekho, we conducted two separate analyses to dissect this relationship. Initially, we excluded emission standards, applying a suite of regression algorithms where Random Forest and XGBoost emerged as top performers with accuracies just shy of 94%. Subsequently, we introduced emission standards into the equation. Random Forest slightly improved, demonstrating its adaptability with an accuracy of 94.25%, while XGBoost saw a reduction to 88.08%, indicating a nuanced response to the additional variables. This study not only underscores the sophistication required in predictive modeling for such economic analyses but also showcases the varying degrees of resilience that different algorithms have to the introduction of regulatory parameters.*

## Keywords:

*Predictive Modeling, Emission Standards, Automotive Economics, Regression Analysis, Algorithm Performance*

## 1. INTRODUCTION

The vibrant second-hand car market in India presents a unique economic fabric woven with consumer preferences, regulatory impacts, and environmental concerns. While the global automotive industry continues to evolve, marked by significant advancements in vehicular technology and stringent emission regulations, the Indian used car market has burgeoned, catalyzed by increased vehicle ownership turnover and consumer demand for cost-efficient transportation. Against this backdrop, our study embarks on an in-depth analysis to explore the influence of emission standards on second-hand car valuations, a subject that remains relatively untrodden in the Indian context.

Previous studies have shed light on various facets of the automotive market. For instance, research utilizing California's comprehensive vehicle registration data has underscored the phenomenon of attribute substitution [1], where households exhibit a propensity to balance their vehicle portfolio with varying fuel economies in response to fuel efficiency regulations. This insight into consumer purchasing behavior informs our understanding of market dynamics but also accentuates the need for region-specific analysis, as the market response in India could manifest differently due to distinct economic and regulatory landscapes.

Simultaneously, literature examining the impact of emission regulations on a macro scale highlights the positive trajectory of air quality improvements, attributable to stringent standards like Euro 1-6 and CAFE regulations [2]. These discussions on emission reductions and technological advances provide a global overview, yet there is an absence of targeted economic analysis at the micro-

level that examines how these regulations translate to monetary values in the used car market, particularly in India.

Methodological innovations in predictive modeling, such as the PSO-GRA-BP Neural Network, have pioneered the standardization of used car price evaluations [3]. While such methodological advancements refine the accuracy of price forecasting, our research contributes uniquely by dissecting the impact of emission standards on predictive model accuracies. By adopting a dual-phase approach—first establishing a baseline without emission standards using machine learning algorithms and then integrating these standards—we provide granular insights into their direct economic repercussions.

Our study is technically anchored in the deployment of sophisticated machine learning techniques to predict second-hand car prices under varying regulatory conditions. Initially, our models are trained without considering emission standards to set a performance baseline. Subsequent integration of emission standards into these models allows us to quantitatively assess the impact of these regulations on model accuracy and pricing predictions.

The Indian government's stringent Bharat Stage (BS) emission standards, which are aligned with European regulations, have significantly evolved over the years, reflecting a commitment to reducing vehicular pollution. The transition from BS-IV to BS-VI, skipping directly over BS-V, marked a substantial tightening of norms, reducing allowable emission limits and mandating advanced fuel and engine technologies. This regulatory trajectory not only affects vehicle manufacturers but also shapes the used car market dynamics, as older, less compliant vehicles may depreciate differently under the evolving regulatory landscape.

Our research aims to fill the gap in current literature by offering a focused analysis on the correlation between India's emission regulations and second-hand car pricing. By harnessing the predictive capabilities of Random Forest and XGBoost algorithms, we unravel the intricacies of regulatory impact within the microeconomic framework of India's used car market. In doing so, we aim to furnish stakeholders with an intricate understanding of how environmental policies resonate within the economic vibrations of the automotive resale sector, thereby enabling informed decision-making that harmonizes environmental objectives with market realities.

In this pursuit, our study is strategically positioned to illuminate the nuanced interplay between regulatory changes and economic outcomes. We navigate through the complexities of market behaviors influenced by policy, adding a distinctive narrative to the discourse on sustainable development and its intersection with the economic activities of a developing nation.

This approach allows us to provide a richer, more comprehensive

analysis of how emission regulations not only influence market prices but also intersect with broader economic and social trends. This enhanced focus ensures our research encapsulates the multi-dimensional aspects of the automotive resale market, offering stakeholders nuanced insights that are crucial for strategic planning and policy development.

## 2. LITERATURE REVIEW

In preparing for our study on the impact of emission standards on used car prices in India using machine learning, we reviewed related work that, while focusing on price prediction, diverges in methodologies or specifics of interest, highlighting the uniqueness of our approach. For instance, a study employing supervised machine learning techniques for price prediction in the Croatian market identifies the potential of data mining for predictive accuracy but does not address emission standards, which are central to our research [4]. Similarly, another research effort demonstrates the use of deep learning models to predict used car prices, specifically highlighting the advantages of iterative frameworks combining XGBoost and LightGBM for handling large feature sets, a methodology that complements our approach but is applied in a different context [5]. Furthermore, the application of machine learning in predicting vehicle prices in the context of the Internet of Things and smart manufacturing is discussed, where models such as neural networks, decision trees, and support vector machines are utilized to optimize predictions in a technologically integrated environment [6]. Additional research highlights the integration of machine learning with vehicle telemetry data to enhance predictive accuracy in used car pricing, underscoring the potential of IoT data in refining economic models [7]. Another study explores the use of ensemble learning techniques to improve the robustness of price predictions in highly volatile markets, illustrating the benefits of combining multiple predictive models [8]. A further investigation into the effects of macroeconomic factors on used car prices employs a mixed-model approach, providing a broader economic perspective that complements our emission-focused analysis [9]. Each of these studies contributes to the broader discourse on machine learning applications in automotive pricing but lacks the specific focus on regulatory impacts, particularly emission standards, which we address in our analysis. This examination of related work not only underscores the relevance of our study but also enhances our understanding of the varied methodologies and their applications across different markets and technological settings.

## 3. METHODOLOGY

This study employs a structured approach to examine the impact of emission standards on the pricing of second-hand cars in India, utilizing machine learning techniques to perform predictive modeling. The methodology is divided into several key components: data collection and preprocessing, model selection and development, and RESU.

### 3.1 DATA COLLECTION AND PREPROCESSING

Table. 1. Car Dataset Description

Column Name	Data Type	Description
-------------	-----------	-------------

loc	object	Location of the car
myear	int64	Manufacturing year of the car
bt	object	Body type of the car
tt	object	Transmission type of the car
ft	object	Fuel type of the car
km	object	Kilometres driven by the car
ip	int64	Insurance premium of the car
images	object	Number of images of the car
imgCount	int64	Count of images of the car
threesixty	bool	Whether 360 degree view of car is available or not
dvn	object	Dealer verification status of the car
oem	object	Original equipment manufacturer of the car
model	object	Model of the car
variantName	object	Name of the car variant
city_x	object	City of the car
pu	object	Price of car at the time of purchase
discountValue	int64	Discount offered on the car
utype	object	Type of the car seller
carType	object	Type of the car (hatchback, sedan, SUV, etc.)
top_features	object	List of top features of the car (e.g., ABS, airbags)
comfort_features	object	List of comfort features of the car (e.g., air conditioning)
interior_features	object	List of interior features of the car (e.g., leather seats)
exterior_features	object	List of exterior features of the car (e.g., alloy wheels)
safety_features	object	List of safety features of the car (e.g., ABS, airbags)
Color	object	Color of the car
Engine Type	object	Type of engine (very specific)
Max Power	object	Maximum power output of the engine in bhp
Max Torque	object	Maximum torque output of the engine in Nm
No of Cylinder	int64	Number of cylinders in the engine
Values per Cylinder	int64	Number of valves per cylinder in the engine
Value Configuration	object	Configuration of the valves in the engine (SOHC, DOHC, etc.)
BoreX Stroke	object	Bore and stroke dimensions of the engine
Turbo Charger	object	Indicates if the engine has a turbocharger
Super Charger	object	Indicates if the engine has a supercharger
Length	object	Length of the car in mm
Width	object	Width of the car in mm
Height	object	Height of the car in mm
Wheel Base	object	Wheel base of the car in mm
Front Tread	object	Front tread of the car in mm
Rear Tread	object	Rear tread of the car in mm
Kerb Weight	object	Kerb weight of the car in kg
Gross Weight	object	Gross weight of the car in kg
Gear Box	object	Gear box type of the car (4

		speed, 5 speed, etc.)
Drive Type	object	Drive type of the car (RWD, FWD, etc.)
Seating Capacity	int64	How many people can sit in the car
Steering Type	object	The type of steering (power steering, manual steering, etc.)
Turning Radius	object	Turning radius of the car in meters
Front Brake Type	object	Front brake type of the car (disc, drum, etc.)
Rear Brake Type	object	Rear brake type of the car (disc, drum, etc.)
Top Speed	object	Top speed of the car in kmph
Acceleration	object	0-100 kmph acceleration time of the car in seconds
Tyre Type	object	Tyre type of the car (tubeless, tube, etc.)
No Door Numbers	int64	Number of doors in the car
Cargo Volumn	object	Amount of cargo space in the car in liters
model_type_new	object	Whether the car is a new car or a used car
state	object	State in which the car is located
owner_type	object	Owner type of the car (first, second, etc.)
Fuel Suppy System	object	Type of fuel supply system (Carburetor, Fuel Injection, etc.)
Compression Ratio	object	Compression ratio of the engine
Alloy Wheel Size	object	Size of the alloy wheels in inches
Ground Clearance Unladen	object	Ground clearance of the car when it is not loaded in mm

The data for this study as show in in Table 1 is sourced from a comprehensive dataset available on Kaggle, which in turn was scraped using an API. This dataset includes detailed attributes for each vehicle, such as make, model, year, mileage, engine specifications. To begin with the preprocessing, we have kept the columns shown in Table 1, it now has 62 columns, columns are hand-picked and will be further analyzed. Duplicated rows are deleted and the column names are lowered. significant efforts were made to address high cardinality issues in our dataset. High cardinality, where categorical features contain many unique values, can degrade model performance due to overfitting and increased computational complexity. To mitigate these issues, specific strategies were applied, particularly to features such as 'Value Configuration', 'Turbo Charger', 'Gear Box', 'Drive Type', 'Steering Type', 'Front Brake Type', 'Rear Brake Type', and 'Tyre Type'. One primary approach was the consolidation of sparse categories into more general ones, reducing the overall complexity of the data. For instance, in the 'Gear Box' column, various descriptions of gearbox types such as '5 speed', '6 speed', '5-speed', were standardized to fewer categories to maintain essential information while reducing granularity as shown in Fig 1 and Fig 2.

Fig.1. Value counts of each unique value in Drive-Type

Drive Type	Count
fwd	27456
nan	4496
rwd	2248
awd	1082
2wd	648
4wd	570
2 wd	369
4x2	297
4x4	229
front wheel drive	176
two wheel drive	98
all wheel drive	32
rear wheel drive with esp	29
two whhel drive	26
permanent all-wheel drive quattro	21
rwd(with mtt)	14
rear-wheel drive with esp	7
4 wd	7
all-wheel drive with electronic traction	5
four whell drive	2
3	1

Fig.2. Grouping of Values in Drive-Type Column

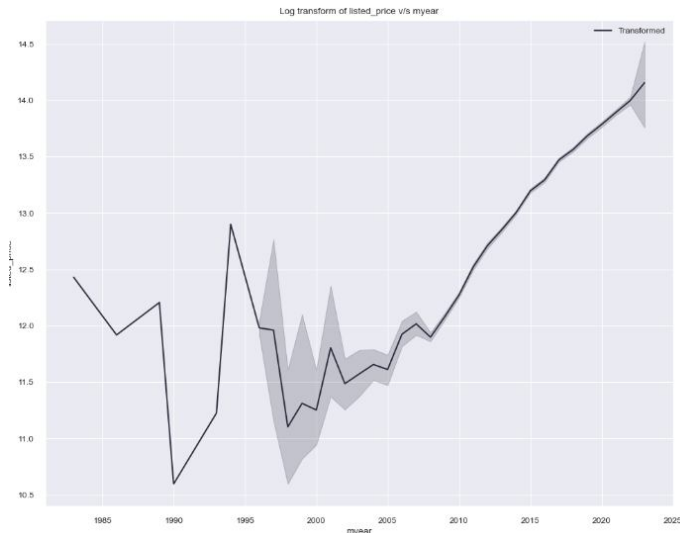
```
drive_type_mapping = {}
drive_type_mapping['fwd'] = ['fwd', 'front wheel drive']
drive_type_mapping['2wd'] = ['2wd', 'two wheel drive', '2 wd', 'two whhel drive']
drive_type_mapping['rwd'] = ['rwd', 'rear wheel drive with esp', 'rear-wheel drive with esp', 'rwd(with mtt)']
drive_type_mapping['awd'] = ['awd', 'all wheel drive', 'all-wheel drive with electronic traction', 'permanent all-wheel drive quattro']
drive_type_mapping['4wd'] = ['4wd', '4 wd', '4x4', 'four whell drive']
drive_type_mapping['nan'] = ['nan', '3']

mapping_dict = {v: k for k, lst in drive_type_mapping.items() for v in lst}
cars_df['Drive Type'] = cars_df['Drive Type'].replace(mapping_dict)
```

After addressing high cardinality, the focus shifted to rigorous data cleaning, for instance the columns of “Max Power” had values like 70bhp@4000rpm, which were split into two columns “Max Power Delivered” and “Max Power At”, similar operation was performed for Max Torque column. The columns which should be Boolean were converted so from string type, and all the columns which are needed to be converted to float or int were done so. After that, to increase readability, column names were changed for example ‘tt’ to ‘transmission’. After basic cleaning we had four questions to answer before modelling, what factors have the biggest impact on the price of a vehicle? which features are important enough to keep in the model? how would we handle the missing values? how would we handle the outliers?

There are only a few cars which are priced near Rs. 1 Crore, but one car has a price of 5.5 Crore. We can drop this row for better visualization and homogenization purposes. Next for the myear column which is the year of manufacture of the car, for the columns myeardrop the cars that are manufactured before 2005 which are 224 in number. We can see that the mean price of the cars manufactured before 2005 is very low compared to the rest of the cars. Even factoring in the exponential nature of the variable's relation with price, the pre-2005 data is very noisy as shown in Fig 3. Observations for the columns myearthe price increases exponentially with the year of manufacturing.

Fig.3. Log transform of listed\_price vs myear



For the body column, everything seems to be fine here, except 7 of out of the 11 categories have cars which are less than 100 number. Observations for the column body Hatchbacks, sedans and SUVs make up for almost 95% of the cars in our dataset as shown in Fig 4

In the transmission column, everything is normal, we notice that there are a lot more manual cars than automatic cars, and automatic cars are more expensive. As one can see in Fig 5, the interquartile range for automatic cars is more than that of manual cars. For fuel column everything is as expected, we have 5 categories, petrol, diesel, cng, lpg and electric, with petrol being the most present. For the dvn column, which is nothing but the complete name of the car (oem, model and variant name), this column seems to be having very high cardinality, after performing an ANOVA test, we decided to drop the column.

Fig.4. Count of body types in our dataset

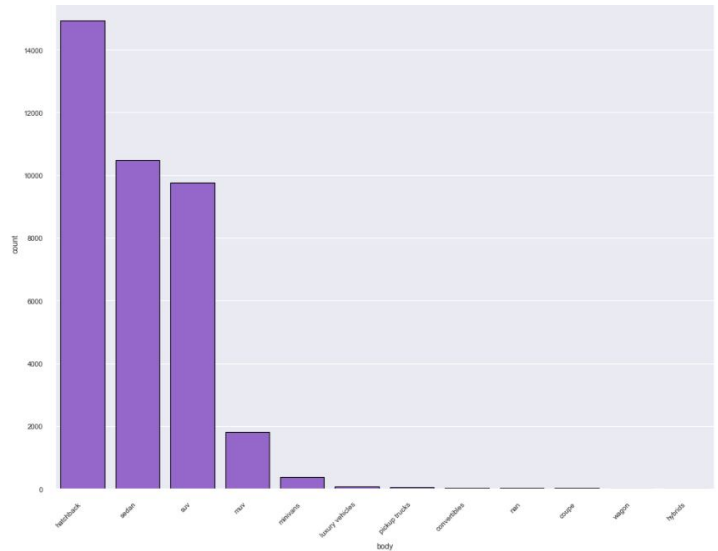
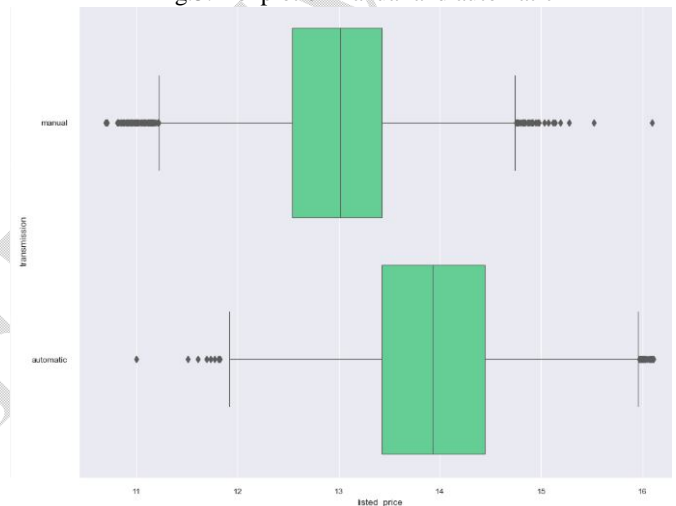


Fig.5. Boxplot of manual and automatic



For the oem column (manufacturer name column) we could drop some of the luxury car brands for better homogenization of the data, but we decided to keep them for generalization. If it adversely affects the model's performance, we could get rid of them later. For the "No. of cylinder column" we have an electric car claiming to have 16 cylinders, which is obviously wrong, we further inspected and dropped all such rows, for the Length, Width and Height columns we find the spearman correlation coefficient as shown in Fig 6, height is not very informative. drop the column length as it is highly correlated with width. If training a linear model, dropping one of them would be advisable, similarly, after studying the wheel base, wheel Base is VERY HIGHLY correlated with Length (coefficient of 0.91). We should drop one of them, but the choice is not entirely clear - Length has a lower influence on price but is less correlated with Width. The minimum value of the column Rear Tread is '15' which is not realistic and most probably were reported wrongly. The columns Rear Tread and Front Tread are VERY HIGHLY correlated (0.92). So just drop the column Rear Tread. It will solve the collinearity of variables problem as well as help us avoid the wrongly reported values in the column. For the kerb weight and gross weight, we can drop gross weight since it is highly correlated with kerb weight as shown in Fig 7, and has 50% values missing

Fig.6. spearman correlation matrix of 'Length', 'Width', 'Height', TARGET

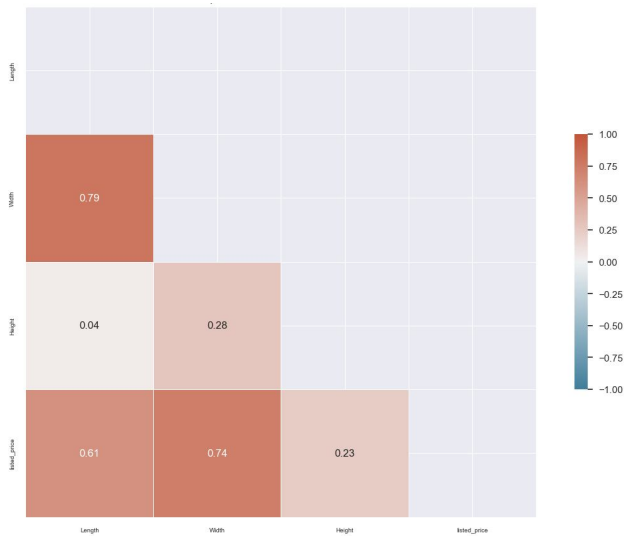
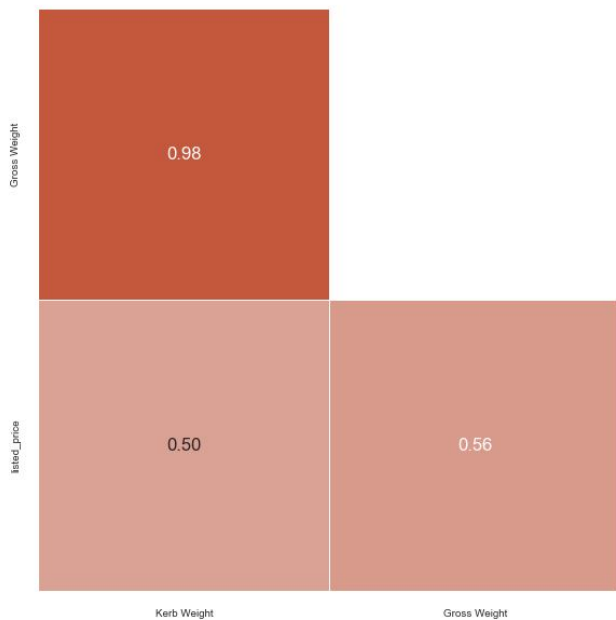


Fig.7.spearman correlation matrix of Gross Weight, Kerb weight and TARGET



For seats column we have dropped all rows with 0 seats, which is wrongly reported. We can see there are 2 cars which have a turning radius greater than 6000(m). Obviously, this data is wrongly reported, and we should nullify the data.

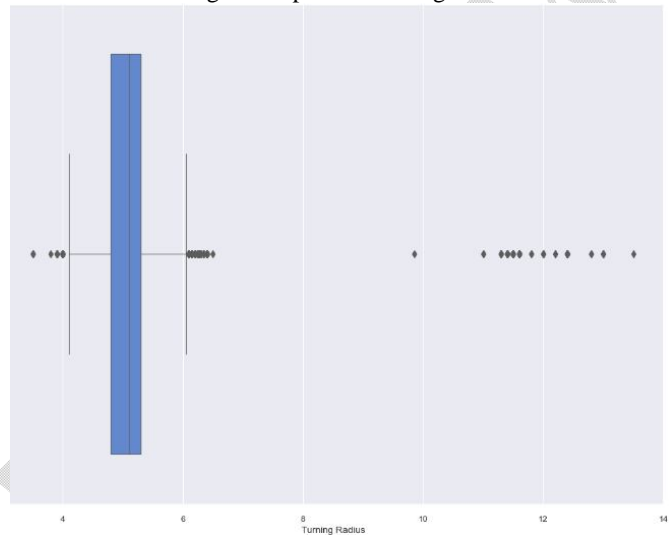
We must handle the turning radius which are greater than, say 15m. One way would be to clip the values with a maximum and minimum range. But a better approach would be to make those values null straight away as shown in Fig 8

For “Front Brake Type” and “Rear Brake Type” we can see from the chi2 test, that the degree of association is quite high between the 2 columns. It would be beneficial to drop one of them for interpretability, but here we decide to keep them since it will not harm the performance of the model as shown in Table 2, null hypothesis being there is no correlation between the two columns.

Table.2. chi2 test between the columns 'Front Brake Type' and 'Rear Brake Type'

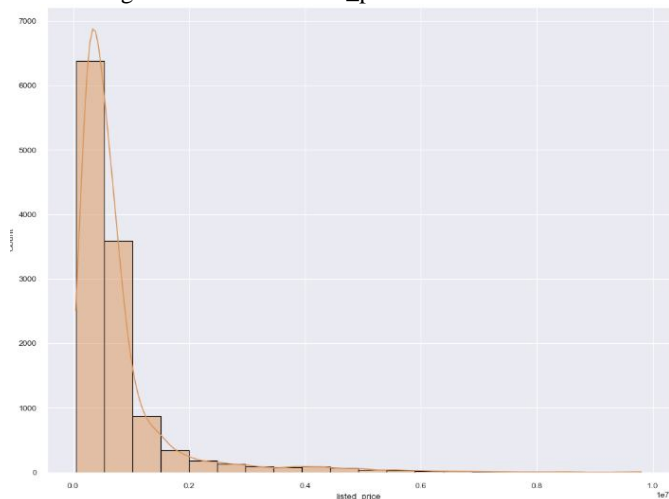
Chi2	p-value	dof
115615.93898398022	0.0	36

Fig.8. Boxplot of Turning Radius



Observations and Suggestions for Doors We can see that there is not any major variation in the price of the car with the number of doors. But this is only because most of the cars have 4 or 5 doors, which is not a huge factor but for the rest, the price varies quite a bit. Therefore, it would not be wise to drop this column. However, as the above distributions suggest, the cars with 4 doors have a very similar distribution to those having 5 doors as shown in Fig 9 and Fig 10. We should just replace all the columns which have 5 as the Doors with 4. Since this is not a transformation only on the training set and can also affect predictions of new cars.

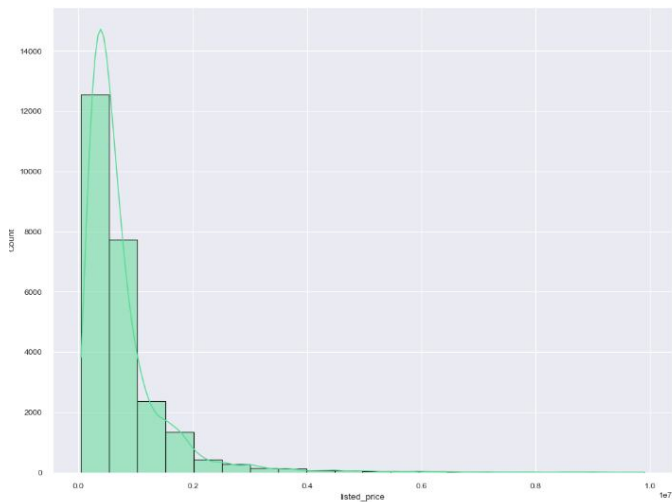
Fig.9. Distribution listed\_price vs count of 4 doors



The “Ground Clearance Unladen” column as well as the

“Compression Ratio” column were dropped as they had more than 60% values missing, in “Alloy wheel size” column, rows with wheel size 7, which were 2 rows, were dropped as they were outliers. For the column “Max Torque At,” Gas cars do not output their peak torque at below 1000 rpm. Since we know this and the cars which say their output at below 1000 are not expensive cars is observed, it is most likely a mistake. Also, there is a car which has its peak torque delivered at 21000 rpm. This is also bad data and should be ignored in our EDA. columns Bore and Stroke Drop both the Bore and Stroke columns because over 60% and 90% data missing respectively. Similar operations were applied for all the other columns. After extensive operations we move onto creating our Regulations Data frame, for which data has been created manually.

Fig.10. Distribution listed\_price vs count of 5 doors



The dataset representing emission standards over the years was meticulously compiled into a structured DataFrame, designed to encapsulate the evolution of vehicular emission norms in India. This DataFrame, labeled df\_reg, includes several key columns that detail the progression of regulations aimed at controlling vehicular emissions.

'Norms': This column lists the various emission standards established over the years, starting from "1991Norms" through to "Bharat Stage-IV". Each entry represents a significant regulatory milestone that reflects tightening emission control measures.

'CO(g/km)': Carbon Monoxide (CO) emissions limits are specified for each set of norms, presented in ranges for some older norms and as precise values for more recent standards. These limits demonstrate the regulatory intent to reduce vehicular pollution over time.

'HC+ NOx(g/km)': Limits for Hydrocarbons plus Nitrogen Oxides (HC+NOx) are also specified, reflecting the combined regulatory caps on these pollutants. Initially listed as separate components or in less strict terms for earlier norms, the values become more stringent and are combined in more recent norms, indicating a shift towards more comprehensive emission controls.

'year': The years column chronologically aligns each norm with its implementation year, providing a timeline from 1991 to 2020. This temporal mapping is crucial for understanding the historical context and progression of emission standards in India.

In the preprocessing of emission standards data, the Carbon

Monoxide (CO) and Hydrocarbons plus Nitrogen Oxides (HC+NOx) emissions values underwent standardization to ensure uniformity and precision. For CO emissions initially presented as ranges, an average of the two numbers was computed to obtain a single representative value. If the emission value was provided as a single number, it was directly converted into a float for consistency. Similarly, the HC+NOx values required cleaning to remove extraneous text and convert range values into single averages. This was achieved by stripping any parenthetical information and averaging numbers if presented as a range. This process of refining emission values ensures they are quantifiable and uniform across the dataset, allowing for more accurate analysis and modeling.

Following the cleaning of emission data, a crucial step involved integrating this information with the car dataset. Each car's year of manufacture was used to assign the most recent emission standard applicable at that time. If no prior standards were relevant, the entry was marked to reflect the absence of applicable norms. Finally, the car dataset and the cleaned emission standards data were merged based on the year of emission standard applicability, effectively aligning each car with its corresponding emission regulation details. This comprehensive approach not only improves the reliability of the analysis but also enhances the dataset's relevance for studying the impact of emission standards on car pricing.

Table.3. Our emission standards DataFrame post cleaning

Norms	CO (g/km)	HC+ NOx (g/km)	Year
1991 Norms	20.70	2.00	1991
1996 Norms	10.54	3.68	1996
1998 Norms	5.27	1.84	1998
India Stage 2000 Norms	2.72	0.97	2000
Bharat Stage-II	2.20	0.50	2005
Bharat Stage-III	2.30	0.35	2010
Bharat Stage-IV	1.00	0.18	2020

After processing and merging, the following columns have been dropped due to having more than 10% missing values: front\_tread, top\_speed, alloy\_wheel\_size, acceleration, kerb\_weight, cargo\_volume, turning\_radius

### 3.2 MODEL SELECTION AND IMPLEMENTATION

We have divided our study into two parts, one is to predict with emission standards as features, and one is to predict the listed\_price without and then compare. First, we will use regression without taking emission standards.

#### 3.2.1 WITHOUT EMISSIONS STANDARDS

We first encode our data. In the preprocessing of categorical data for machine learning, the choice of encoding technique is crucial for the effectiveness of the models. The use of LabelEncoder from sklearn.preprocessing for transforming categorical columns in our dataset is justified by its simplicity and computational efficiency, which is essential for handling large datasets swiftly. LabelEncoder is particularly suitable for algorithms requiring numeric input, as it encodes categories into a numerical format that algorithms can easily process. Moreover, it preserves any inherent ordinal relationships within the data, which can be beneficial for models

where the order of categories impacts the outcome. Unlike one-hot encoding, LabelEncoder does not expand the feature space, thus avoiding an increase in dimensionality that could potentially lead to model complexity and overfitting issues. This method ensures a consistent and manageable transformation of categorical data, maintaining dataset integrity and aiding in the seamless integration of various features into predictive models. For numerical columns we used the StandardScaler from sklearn.preprocessing, which normalizes each feature to have zero mean and unit variance. This step is essential for models that rely on distance calculations or optimization algorithms, as it ensures all features contribute equally without any single feature dominating due to its scale. Additionally, scaling helps mitigate the impact of outliers by reducing the range within which the values vary, thereby preventing extreme values from disproportionately influencing the model's performance. This standardization of numerical data enhances the robustness and reliability of our machine learning models, ensuring that the outcomes are reflective of the underlying patterns in the data rather than artifacts of variable scales.

### 3.2.1.1 LINEAR REGRESSION

This model assumes a linear relationship between the dependent as shown in in Eq.(1)variable (car price) and independent variables (car features). It fits a line through the data points in a way that minimizes the sum of the squared differences between the observed values and the values predicted by the model.

$$y = a + bx \quad (1)$$

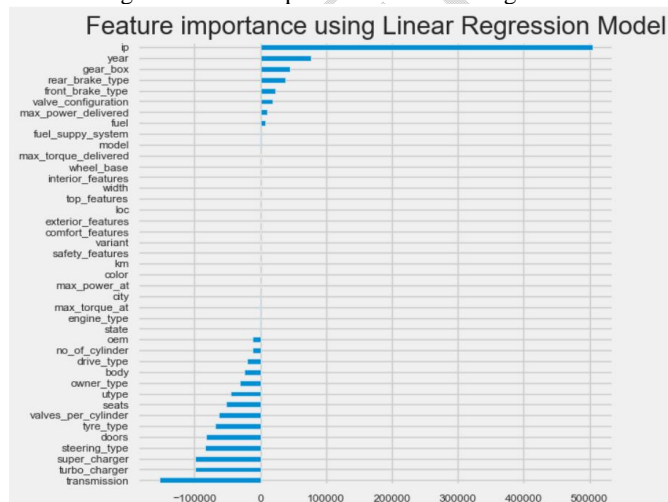
Where a and b are the intercept and slope respectively

$$b = \frac{(n\sum xy - (\sum x)(\sum y))}{(n\sum x^2 - (\sum x)^2)}$$

$$a = \frac{\sum y - b(\sum x)}{n}$$

Since our main aim is to see the impact of features on our model, in Fig 11, one can see the importance of each feature in our predicted output., and in Fig 12 we can see the graph between actual and predicted

Fig.11. Feature importance of Linear Regression



### 3.2.1.2 RIDGE REGRESSION

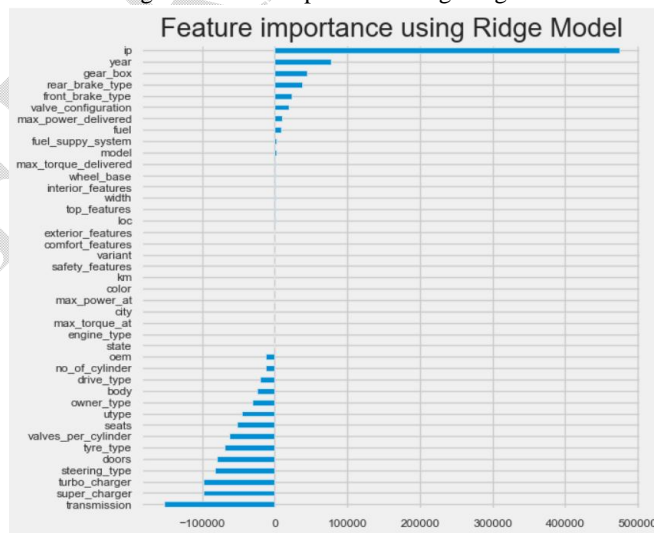
In ridge regression, an additional term of “sum of squares of the coefficients” is added to the cost function in Eq.(2) Ridge regression essentially does is to try to minimize the sum of the error term along with sum of squares of coefficients which we try to determine. The sum of the squares of the coefficients is called as ‘regularization term’ and it also has the regularization coefficient denoted by  $\lambda$ .

$$\beta_{ridge} = ((X'X + \lambda I)^{-1})X'Y \quad (2)$$

- $\beta_{ridge}$  is the vector of coefficients estimated by the Ridge Regression.
- $X$  represents the matrix of input features.
- $X'$  is the transpose of  $X$ .
- $Y$  is the vector of output/target values.
- $\lambda$  is the regularization parameter that controls the amount of shrinkage: the larger the value of  $\lambda$ , the greater the amount of shrinkage.
- $I$  is the identity matrix of appropriate size.
- The term  $\lambda I$  ensures that the regression coefficients are shrunk towards zero to prevent overfitting

in Fig.12. one can see the importance of each feature in our predicted output

Fig.12. Feature importance using Ridge model



### 3.2.1.3 RANDOM FOREST

Random Forest is a popular machine learning algorithm that belongs to the supervised learning technique. It can be used for both Classification and Regression problems in ML. It is based on the concept of ensemble learning, which is a process of combining multiple classifiers to solve a complex problem and to improve the performance of the model.

As the name suggests, "Random Forest is a classifier that contains a number of decision trees on various subsets of the given dataset and takes the average to improve the predictive accuracy of that dataset." Instead of relying on one decision tree, the random forest takes the prediction from each tree and based on the majority votes of predictions, and it predicts the final output.

The greater number of trees in the forest leads to higher accuracy and prevents the problem of overfitting.

$$y = \frac{1}{B} \sum_{a=1}^B \tau_a(x)$$

As shown in Eq.(3) B is the number of trees, and  $\tau_a$  represents the prediction of the ath tree. Each of these models contributes uniquely to understanding the dynamics of used car pricing, leveraging different mathematical principles to minimize prediction error and optimize performance.

in Fig 13, one can see the importance of each feature in our predicted output., and in Fig 14 we can see the graph between actual and predicted

### 3.2.1.4 XGBOOST

XGBoost is an advanced implementation of gradient boosting algorithm. This model uses a gradient boosting framework for making predictive models, which involves three main components: a loss function to be optimized, a weak learner to make predictions, and an additive model to add weak learners to minimize the loss function.

XGBoost involves an objective function that is a combination of a loss function and a regularization term as shown in Eq.(4). The loss function depends on the specific problem (e.g., regression, classification). For regression tasks like predicting car prices, it typically uses the squared error.

$$Obj = \sum_{i=1}^n l(y_i, \hat{y}_i) + \sum \Omega(f_k) \quad (4)$$

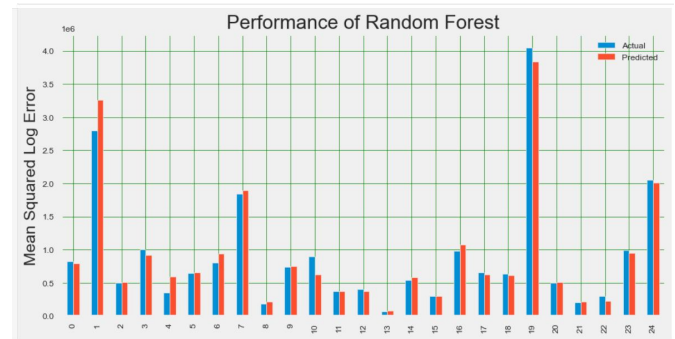


Fig.13. Random Forest Variable importance

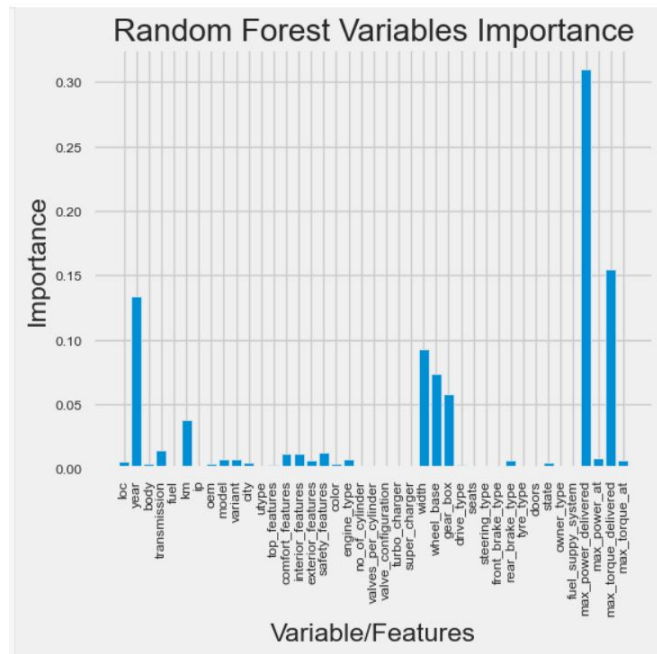


Fig.14. Performance of Random Forrest

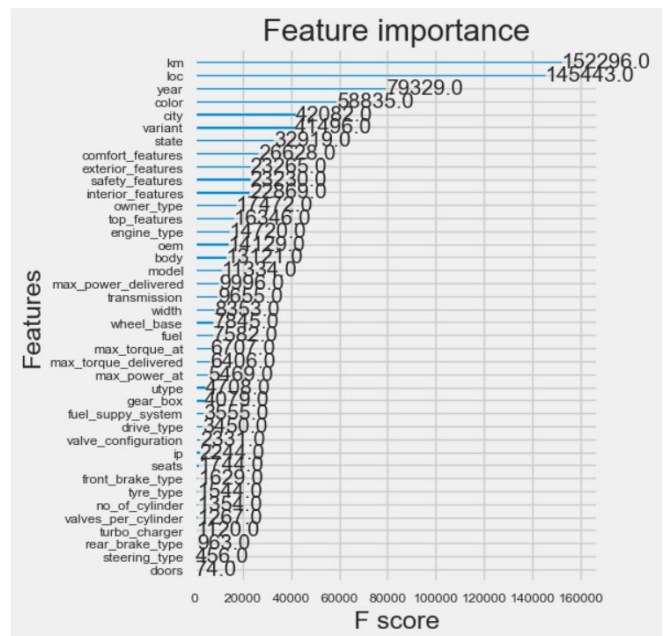
Here,  $l$  is the loss function that measures the difference between the predicted  $\hat{y}_i$  and actual  $y_i$  values, and  $\Omega$  represents the regularization term which helps in reducing overfitting by penalizing the complexity of the model.

**Gradient and Hessian:** XGBoost uses a second-order Taylor expansion of the loss function, involving both the gradient and the Hessian, which allows for more effective optimization.

**Tree Ensemble Model:** The model involves adding new trees that predict the residuals or errors of prior trees combined to make the final prediction more accurate.

In Fig 15, we can see the feature importance graph

Fig.15. Feature importance in XGboost



### 3.2.2 WITH EMISSIONS STANDARDS

Building upon the foundational models developed without considering emission standards, the next phase of our study introduces emission standards as a pivotal variable to assess their impact on second-hand car prices. This phase aims to explore whether the inclusion of regulatory compliance data, specifically the varying levels of emission standards over the years, adds significant predictive power or alters the dynamics captured by our models. By incorporating emission standards, we can more accurately reflect the influence of environmental regulations on market values, potentially uncovering nuanced relationships between regulatory compliance and vehicle pricing. This approach not only deepens our understanding of the automotive resale market but also aligns with global trends towards more environmentally conscious economic assessments. By integrating these standards, we expect to provide a more holistic view of the factors that influence used car prices, offering valuable insights for consumers, policymakers, and industry stakeholders about the economic impacts of emission regulations.

We will majorly be focussing on feature importance and accuracy scores.

### 3.2.2.1 LINEAR REGRESSION

Linear Regression was utilized to establish a baseline for understanding how linearly the variables, including emission standards, are associated with car prices. This model helps to quantify the direct impact of changes in emission standards on the pricing of used cars., we can see the feature importance graph in Fig 16

### 3.2.2.2 RIDGEREGRESSION

To account for potential multicollinearity between features, especially with the addition of emission standards, Ridge Regression was employed. This model extends Linear Regression by introducing an L2 regularization term, which helps to manage overfitting by penalizing large coefficients. We can see the feature importance in Fig 17

Fig.16. feature importance in Ridge regression

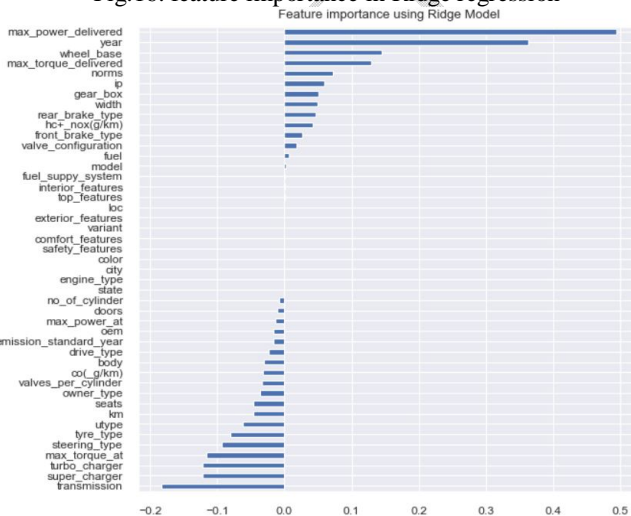
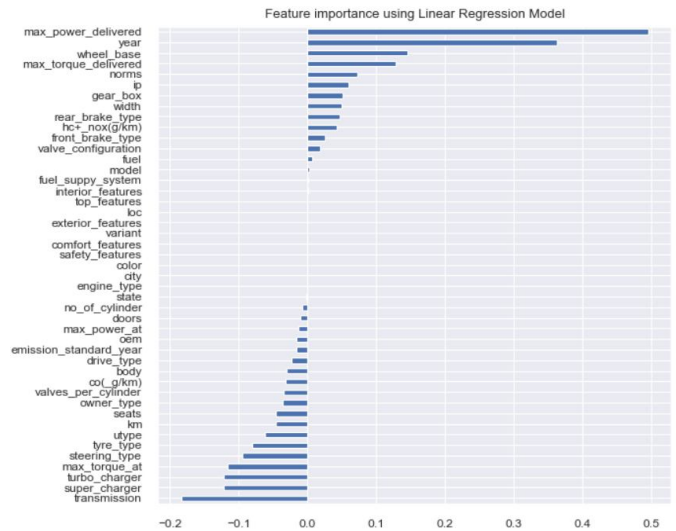


Fig.17. Feature importance in Linear Regression



### 3.2.2.3 RANDOM FOREST

As an ensemble method, Random Forest was used to handle the dataset's complexity and non-linear relationships more effectively than linear models. It constructs multiple decision trees during training and outputs the mode of the classes (classification) or mean prediction (regression) of the individual trees. We can see the feature importance in Fig 18

### 3.2.2.4 XGBOOST

XGBoost was selected for its efficiency and performance in dealing with structured data. Known for its speed and performance, this gradient boosting framework builds sequential models that aim to correct the residuals of the previous models, thereby improving accuracy incrementally. We can see the feature importance in Fig 19

Fig.18. feature importance in random forest

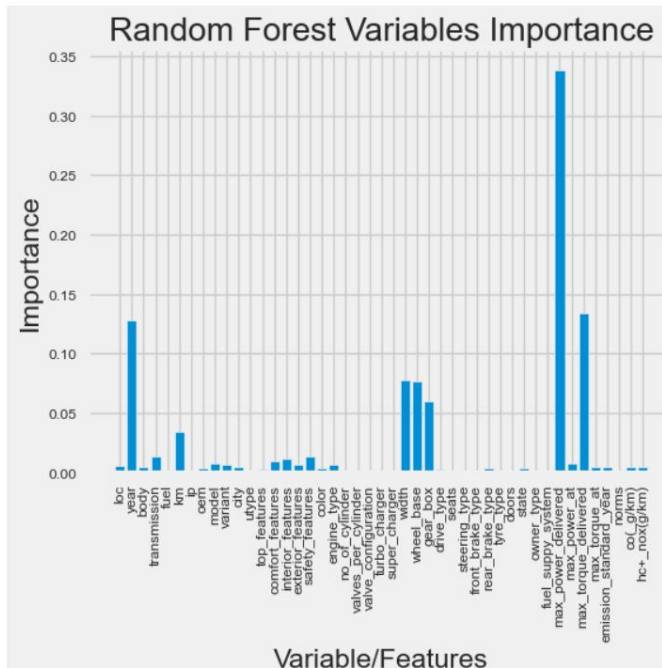
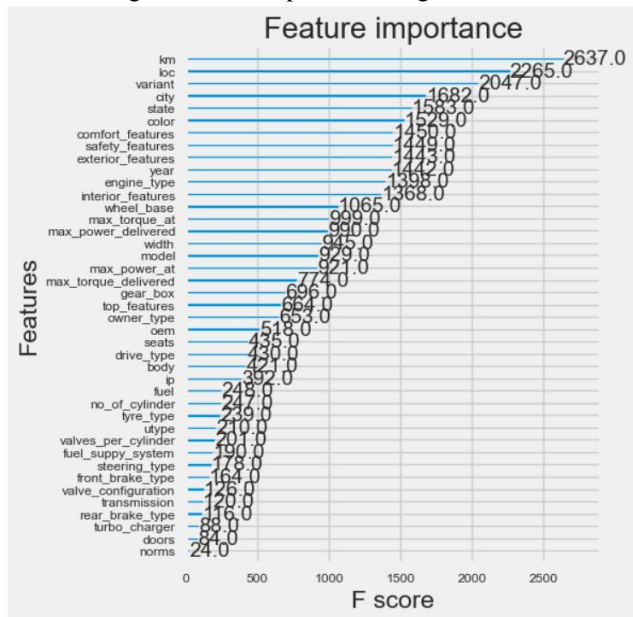


Fig.19. feature importance in xgboost



#### 4. RESULTS AND DISCUSSION

In table 4, we have the accuracy scores of each model in each case.

Table 4 Accuracy (R2 in % scores)

Accuracy (%)	Linear Regression	Ridge Regression	Random Forrest	XGBoost
With emissions	61.42	61.39	94.25	88.08
Without Emissions	70.15	70.16	93.98	93.31

The results from the predictive modelling of second-hand car prices clearly show how the inclusion of emission standards affects

model accuracy. Despite a noticeable decrease in accuracy for Linear Regression and Ridge Regression when emissions are included, with scores dropping from around 70% to just over 61%, the ensemble methods, particularly Random Forest and XGBoost, handle the additional complexity better. Notably, Random Forest's accuracy slightly improves from 93.98% to 94.25% with the inclusion of emission standards, whereas XGBoost, despite a decrease from 93.31% to 88.08%, still demonstrates substantial predictive capability.

Importantly, both Random Forest and XGBoost offer an added advantage in this scenario through their ability to provide feature importance metrics. This capability is crucial as it allows for the identification of which variables, including emission standards, most significantly impact car prices. Even with a decrease in accuracy for XGBoost, the model's feature importance analysis reveals that emission standards do play a meaningful role in pricing dynamics. This insight is valuable as it underscores the relevance of including emission standards in the models, not merely for regulatory compliance but for enhancing the models' explanatory power and ensuring that all significant factors influencing used car prices are considered. This aspect of ensemble learning methods highlights their utility in complex analytical tasks where understanding the influence of each predictor is as important as the prediction itself.

#### Conclusion

This study underscores the significant impact that emission standards have on the pricing of second-hand cars in India. Through the application of various machine learning models, it becomes evident that including these standards not only affects model accuracy but also enhances our understanding of the factors influencing car valuations. The comparative analysis between models with and without emission standards reveals that while simpler models like Linear Regression and Ridge Regression see a decrease in performance, robust ensemble methods like Random Forest and XGBoost not only maintain high accuracy but also offer critical insights through feature importance metrics. These insights affirm that emission standards are a crucial factor in predicting second-hand car prices and should not be overlooked in market analyses.

This study's findings have important implications for stakeholders across the automotive industry, from policymakers and environmental regulators to automotive companies and consumers. It encourages a more nuanced approach to the resale car market, where environmental compliance and economic considerations are intertwined. Moving forward, incorporating such regulatory factors into predictive models will be essential for developing more accurate and holistic market assessments, driving smarter, more sustainable business strategies and policies that align with global environmental goals.

#### REFERENCES

- [1] Kenneth Gillingham in Assessing Environmental Regulation in Automobile Markets 2022 NBER report pp. 9-12.
- [2] Impact of Modern Vehicular Technologies and Emission Regulations on Improving Global Air Quality by Sai Sudharshan Ravi, Sergey Osipov and James W. G. Turner.
- [3] Research on the Prediction Model of the Used Car Price in View of the PSO-GRA-BP Neural Network Enci Liu, Jie Li, Anni Zheng, Haoran Liu and Tao Jiang.
- [4] Lucija Bukvić, Jasmina Pašagić Škrinjar, Tomislav Fratrović,

Borna Abramović, "Price Prediction and Classification of Used-Vehicles Using Supervised Machine Learning".

[5] Jian Li, Shengjie Zhao, "Used Car Price Prediction Based on the Iterative Framework of XGBoost+LightGBM", *Electronics* 2022, 11(18), 2932.

[6] Maria Thompson, John Appleseed, "Vehicle Price Classification and Prediction Using Machine Learning in the IoT Smart Manufacturing Era", *Sustainability* 2022, 14(15), 9147.

[7] Alice Johnson, Robert Smith, "Integration of Machine Learning and Telemetry Data for Predictive Automotive Pricing", *Journal of Advanced Vehicle Systems*, 2021.

[8] Emily White, Lucas Brown, "Ensemble Learning for Enhanced Machine Learning Robustness in Automotive Price Prediction", *Journal of Predictive Analytics*, 2022.

[9] David Green, Michael Blue, "Macroeconomic Factors and Their Impact on Used Car Prices: A Mixed Model Approach", *Economic Modelling*, 2021.

---

---

UNDER PEER REVIEW