

Used Cars Price Prediction Model Training

Roxy Zhang

3/25/2022

Introduction

Since the market of used cars is growing, there is an increasing need of the information of used cars prices. Therefore, I want to build a model predicting the price using attributes of the car as predictors, and thus provide information for the potential buyers for decision making.

Our data is from Kaggle. This dataset contains 22 (1998-2019) years of used cars information. There are 5872 rows and 11 predictors, among which 5 are numeric and 6 are categorical.

Data Cleaning

For better illustration, I transformed the original unit (100000 Indian Rupee) of response variable **price** to USD, using the exchange rate of 100000 Indian Rupee equals to 1309.75 USD on March 21, 2022. For numeric variables with unit, the units are deleted for model-building. I also delete the car model from the original **name** variable, keeping only the brand name, since the model can be greatly explained by the predictors I am using. The type **Fourth** in **owner_type** is actually Fourth and above. For analysis completeness, I dropped NA values. I also filtered out a row containing an extremely large value of **kilometers_driven**.

Below are the variables used:

- name**: The brand of the car.

- location**: The location in which the car is being sold or is available for purchase.

- year**: The year or edition of the model.

- kilometers_driven**: The total kilometres driven in the car by the previous owner(s) in KM.

- fuel_type**: The type of fuel used by the car.

- transmission**: The type of transmission used by the car.

- owner_type**: Whether the ownership is Firsthand, Second hand or other.

- mileage**: The standard mileage offered by the car company in kmpl or km/kg.

- engine**: The displacement volume of the engine in cc.

- power**: The maximum power of the engine in bhp.

- seats**: The number of seats in the car.

- price**: The price of the used car in US Dollar.

Exploratory Data Analysis

First, I plot the price against year. For better scaling, the price is log-transformed. There is an increasing trend of price from 1998-2014, while the price goes down rapidly since 2015.

Then I look at the price distribution for 5 categorical predictors respectively. There are some interesting findings:

1. There are several attributes associated with higher price: fuel type diesel, automatic transmission, first-hand owner, and cars with 2 seats.

2. The price range for each brand varies a lot. For example, Lamborghini has a condensed price range

slightly above 150k (which is not surprising), while the majority stays below 25k.

3. There is a relatively large number of outliers observed, indicating the big variance of our response variable. This may be because of the volatility of used cars market, which in turn suggests the significance of our project.

For better visualization, categorical variables with over ten terms (`name`, `location`) are filtered out in the correlation plot. From the upper-right corner of the plot we can see a strong positive correlation between `engine` and `power`. Besides, `milage` is negatively associated with both `engine` and `power`.

Data Partitioning

The data is split into training data and testing data by the proportion of 0.8:0.2.

Models

Model Building

4 models are used to fit the training data: Lasso, Elastic Net, Partial least squares (PLS), and Multivariate Adaptive Regression Spline (MARS). Cross validation is used to select the best parameter or parameter combination for each model.

Lasso

The tuning parameter λ controls the L1 regularization, as λ increases, the number of predictors in the model decreases. Setting the candidate values of λ to be from 0.1353353 to 54.59815 with 100 steps, the best-tune λ is 2.069. All five numeric predictors are included in the final model.

Elastic Net

Elastic net is a regularized regression method that linearly combines the L1 and L2 penalties of the Lasso and Ridge methods. In addition to setting and choosing a lambda value, elastic net also allows us to tune the alpha parameter where $\alpha = 0$ corresponds to Ridge and $\alpha = 1$ to Lasso. The optimum alpha chosen is 0.25, and the optimum λ is 5.136.

Partial least squares (PLS)

Instead of finding hyperplanes of maximum variance between the response and independent variables, PLS finds a linear regression model by projecting the predicted variables and the observable variables to a new space. The tuning parameter - number of components in the final model is 18.

Multivariate Adaptive Regression Spline (MARS)

MARS is a non-parametric regression technique and can be seen as an extension of linear models, which can be used to model nonlinear relationships and interactions between a set of predictor variables and a response variable. The best tune metrics are `nprune` = 12, `degree` = 2. There are 12 terms in the final model. The plots for tuning parameters are consistent with the best tuning parameters reported above. From the variance importance plots above, `power` and `year` exist in all four plots, indicating that they are important predictors.

Model selection

By plotting the root mean square error (RMSE) of training data, MARS model can be selected since it has the smallest RMSE.

Test model performance

For the 3 linear models, the test error is calculated using the test data. The test error of Lasso, Elastic net, PLS, MARS is 38896926.01, 38710263.67, 38602924.70, respectively. The PLS model has the least test error, so it is more preferred among the linear models.

Model limitations

The penalty function of **Lasso** has several limitations. For example (not the case of this data), in the “large p, small n” case (high-dimensional data with few examples), the Lasso selects at most n variables before it saturates. Also if there is a group of highly correlated variables, then the Lasso tends to select one variable from a group and ignore the others. As is mentioned above, there are some association among **power**, **engine** and **mileage**, thus the lasso model’s performance is not so well.

To overcome this limitation, the **Elastic Net** adds a quadratic part to the penalty, which when used alone is ridge regression. However, one disadvantage of Elastic Net is the computational cost. We need to cross-validate the relative weight of L1 vs. L2 penalty, and the number of values in the α grid increases computational cost. Besides, with greater flexibility comes increased probability of overfitting.

For **PLS**, the major limitations are a higher risk of overlooking ‘real’ correlations and sensitivity to the relative scaling of the descriptor variables.

A major disadvantage of **MARS** is that if the underlying function involves a pure high-order interaction effect with absolutely no corresponding main or lower-order interaction effects, then at best MARS will have to enter several more or less randomly selected low-order basis functions before it can capture this effect, or at worst it will fail to detect it.

Conclusions

Here are some conclusions drawn from data visualization and model training:

1. Some interesting facts are found from the data, such as having 2 seats and fueled by diesel are associated with higher price, and the various price distribution of different brands.
2. The MARS model performed best with the lowest training RMSE.
3. The maximum power of the engine, and the year of the car serve as important predictors in all models, which is as expected since these two attributes are the main consideration of buyers intuitively.
4. If the model interpretability is a major concern, the PLS model can be considered since it has an acceptable RMSE and better interpretation than the MARS model.