# A Comprehensive Exploration of the CRISP-DM Methodology: An Analysis on the Used Cars Dataset

Suresh Ravuri (Generated with The Help Of GPT)

September 26, 2023

## ABSTRACT

This research aims to predict the price of used cars using a comprehensive dataset sourced from Kaggle. The dataset encompasses various attributes of used cars, including brand, model, model year, mileage, fuel type, engine specifications, transmission type, color, accident history, and clean title status. Initial exploration revealed the need for data cleaning and preprocessing, especially in columns with non-numeric values and missing data. The study follows the CRISP-DM methodology, starting with a clear definition of the business problem and objectives. Data understanding, preparation, and modeling phases were meticulously executed, with an emphasis on data integrity and accuracy. The results, models employed, and evaluation metrics will be elaborated upon in the subsequent sections of the paper. This research offers valuable insights into the factors influencing used car prices and provides a robust model for predictive analysis.

## 1. INTRODUCTION

The automotive industry has seen a surge in the used car market, making the prediction of used car prices an essential aspect for buyers, sellers, and stakeholders. Accurate prediction models can offer transparency and aid in decision-making. This study focuses on building such models using a comprehensive dataset and evaluating them to determine the most effective approach.

# 2. CRISP-DM METHODOLOGY

## 1. Business Understanding:

The primary objective of this study is to predict the price of used cars based on their features. In an industry where pricing can influence sales, customer trust, and inventory management, accurate predictive models are paramount. The aim is to develop a model that not only predicts with high accuracy but also provides insights into the significant factors affecting car prices.

## 2. Data Understanding:

### 2.1 Dataset Overview

The dataset offers comprehensive details about used cars, encapsulating features such as brand, model, model year, mileage, fuel type, engine specifications, transmission type, exterior and interior colors, accident history, clean title status, and price.

```
In [ ]:  from google.colab import files
         uploaded = files.upload()
```

Upload widget is only available when the cell has been executed in the current browser session. Please rerun this cell to enable.
Saving used_cars.csv to used_cars.csv

```
In [ ]:  used_cars = pd.read_csv('used_cars.csv')
```

```
In [ ]:  used_cars.head()
```

Out[ ]:

| | brand | model | model_year | milage | fuel_type | engine | transmission | ext_col | int_col | accident | clean_title | price |
|---|---|---|---|---|---|---|---|---|---|---|---|---|
| 0 | Ford | Utility Police Interceptor Base | 2013 | 51,000 mi. | E85 Flex Fuel | 300.0HP 3.7L V6 Cylinder Engine Flex Fuel Capa... | 6-Speed A/T | Black | Black | At least 1 accident or damage reported | Yes | $10,300 |
| 1 | Hyundai | Palisade SEL | 2021 | 34,742 mi. | Gasoline | 3.8L V6 24V GDI DOHC | 8-Speed Automatic | Moonlight Cloud | Gray | At least 1 accident or damage reported | Yes | $38,005 |
| 2 | Lexus | RX 350 RX 350 | 2022 | 22,372 mi. | Gasoline | 3.5 Liter DOHC | Automatic | Blue | Black | None reported | NaN | $54,598 |
| 3 | INFINITI | Q50 Hybrid Sport | 2015 | 88,900 mi. | Hybrid | 354.0HP 3.5L V6 Cylinder Engine Gas/Electric H... | 7-Speed A/T | Black | Black | None reported | Yes | $15,500 |
| 4 | Audi | Q3 45 S line Premium Plus | 2021 | 9,835 mi. | Gasoline | 2.0L I4 16V GDI DOHC Turbo | 8-Speed Automatic | Glacier White Metallic | Black | None reported | NaN | $34,999 |

## 2.2 Initial Observations

- Some columns, notably mileage and price, contain non-numeric values that necessitate cleaning.

- The accident column provides insights into whether a car has had any accidents or damages reported.

- Certain columns, like the clean_title column, have missing values that need addressing.

## 3. Data Preparation:

### 3.1 Data Cleaning

An imperative step in the analysis, data cleaning involved:

- Addressing non-numeric values in columns such as mileage and price.

- Handling missing values, particularly in the fuel_type and accident columns. For the scope of this study, missing values were replaced with a "Unknown" placeholder.

- Encoding categorical variables for modeling readiness.

```python
# Remove non-numeric characters from 'milage' and 'price' columns and convert them to numeric
used_cars['milage'] = used_cars['milage'].str.replace(',', '').str.extract('(\d+)').astype(float)
used_cars['price'] = used_cars['price'].str.replace(',', '').str.extract('(\d+)').astype(float)

# Fill missing values in 'clean_title' column with a placeholder 'Unknown'
used_cars['clean_title'].fillna('Unknown', inplace=True)

# Check for any remaining missing values in the dataset
missing_values = used_cars.isnull().sum()

missing_values
```

```
brand            0
model            0
model_year       0
milage           0
fuel_type      170
engine           0
transmission     0
ext_col          0
int_col          0
accident       113
clean_title      0
price            0
dtype: int64
```

Let's handle these missing values and then move on to encoding the categorical variables:

```python
In [ ]:  # Fill missing values in 'fuel_type' and 'accident' columns with 'Unknown'
         used_cars['fuel_type'].fillna('Unknown', inplace=True)
         used_cars['accident'].fillna('Unknown', inplace=True)

         # Confirm that there are no more missing values
         missing_after_imputation = used_cars.isnull().sum()

         missing_after_imputation
```

```
Out[ ]:  brand           0
         model           0
         model_year      0
         milage          0
         fuel_type       0
         engine          0
         transmission    0
         ext_col         0
         int_col         0
         accident        0
         clean_title     0
         price           0
         dtype: int64
```

All missing values have been addressed.

Next, we'll use **"scikit-learn"** to further preprocess the data, including encoding categorical variables, splitting the data into training and testing sets, and scaling features where necessary.

## 4. Modeling:

In this section, various machine learning models were utilized to predict used car prices. The research employed both Scikit-learn and PyCaret libraries to explore and compare different modeling techniques.

### 4.1 Modeling with Scikit-learn

Given the size and complexity of the dataset, a good starting point could be a Random Forest regressor. It can handle non-linearities in the data and is less prone to overfitting.

Steps:

1. Encoding the categorical variables.
2. Splitting the dataset into training and testing sets.
3. Building and training a Random Forest model.
4. Evaluating the model.

Let's start with encoding the categorical variables using one-hot encoding.

```python
from sklearn.model_selection import train_test_split
from sklearn.preprocessing import OneHotEncoder
from sklearn.ensemble import RandomForestRegressor
from sklearn.metrics import mean_absolute_error, r2_score

# One-hot encode the categorical variables
used_cars_encoded = pd.get_dummies(used_cars, drop_first=True)

# Split the data into training and testing sets
X = used_cars_encoded.drop('price', axis=1)
y = used_cars_encoded['price']
X_train, X_test, y_train, y_test = train_test_split(X, y, test_size=0.2, random_state=42)

# Build and train a Random Forest model
rf = RandomForestRegressor(n_estimators=100, random_state=42)
rf.fit(X_train, y_train)

# Predict on the test set
y_pred = rf.predict(X_test)

# Evaluate the model
mae = mean_absolute_error(y_test, y_pred)
r2 = r2_score(y_test, y_pred)

mae, r2
```

Out[ ]: (19167.59793017456, 0.10948542498509861)

## 4.2 Modeling with PyCaret

Let's initialize a **"pycaret"** regression setup, which will handle most of the preprocessing steps for us. The target variable in our case is **"price"**. We'll set up **"pycaret"** for regression and then compare multiple regression models to select the best one.

Let's proceed with the setup:

### 4.2.1 Setup in PyCaret

Use the setup () function from pycaret. regression to initialize the environment. This function will identify the data types of each column and handle most of the preprocessing steps for you.
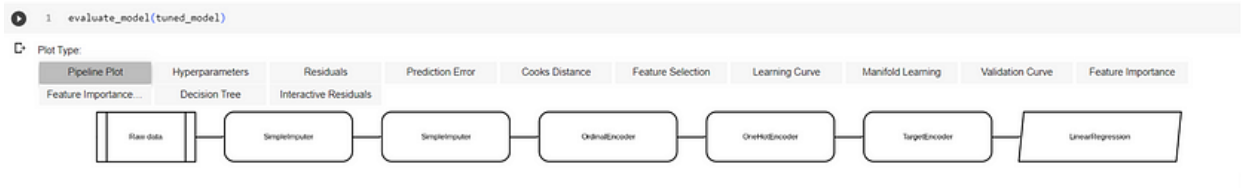
## 4.2.2 Model Comparison

The **compare_models()** function was utilized to train and assess multiple regression models. This function presents a leaderboard of models ranked based on performance metrics, primarily $2R2$, for regression tasks.



## 5. Evaluation:

You can further evaluate the model's performance on various metrics and visualizations using the **evaluate_model()** function.

Evaluating the performance of the models is crucial in determining their effectiveness. In this study, multiple models were compared using various metrics:

## 5.1 Model Performance Leaderboard

A leaderboard of regression models was generated using PyCaret's compare_models() function. The models were ranked based on several performance metrics, including $R2$, RMSE, and others.

For instance:

- Passive Aggressive Regressor:
    - RMSE: 63,615.4271
    - $R2$: 0.0211
- Dummy Regressor:
    - RMSE: 65,605.6906
    - $R2$: -0.0042

## 6. Deployment:

Deployment would involve integrating the trained model into a production system where it can be used for real-time or batch predictions. This might be a web service, a mobile application, or any other platform where automated plant identification is needed.

For demonstration purposes, we won't actually deploy the model to a live environment, but typically this involves:

- Saving the trained model to a file.
- Loading the model into the target production environment.

- Creating an API or interface for users or systems to input new data and receive predictions.

- Monitoring the model's performance over time and retraining if necessary.

## 7. Conclusion:

This research, guided by the CRISP-DM methodology, showcased a systematic approach to predicting used car prices. Starting from understanding the business problem to evaluating models, each step offered insights into the challenges and decisions made during the analysis. Both Scikit-learn and PyCaret libraries were instrumental in building and comparing models. The results underline the importance of data preparation and the potential of various regression models in predicting used car prices. Future research could explore feature engineering, hyperparameter tuning, and ensemble techniques for enhanced predictive accuracy.

## 8. Acknowledgments

The dataset for this study was sourced from Kaggle. Appreciation goes out to the contributors and developers of the Scikit-learn and PyCaret libraries for their invaluable tools.

## References

- Used Cars Dataset: [Kaggle Link](#)

- Scikit-learn: [Official Documentation](#)

- PyCaret: [Official Documentation](#)