| In [1]: | <pre>import numpy as np import pandas as pd #using graphical representation visualization import matplotlib.pyplot as plt</pre> |
|----------------------------|--|
| Out[3]: | |
| In [4]: | <pre><class 'pandas.core.frame.dataframe'=""> RangeIndex: 892 entries, 0 to 891 Data columns (total 6 columns): # Column Non-Null Count Dtype</class></pre> |
| | 1 company 892 non-null object 2 year 892 non-null object 3 Price 892 non-null object 4 kms_driven 840 non-null object 5 fuel_type 837 non-null object dtypes: object(6) memory usage: 41.9+ KB |
| In [5]: | Data cleaning started df.head(3) name company year Price kms_driven fuel_type |
| | Hyundai Santro Xing XO eRLX Euro III Hyundai 2007 80,000 45,000 kms Petrol Mahindra Jeep CL550 MDI Mahindra 2006 4,25,000 40 kms Diesel Maruti Suzuki Alto 800 Vxi Maruti 2018 Ask For Price 22,000 kms Petrol |
| | <pre># isnumeric function to check numeric values df = df[df['year'].str.isnumeric()] # converting into integer (year) df['year'] = df['year'].astype(int)</pre> |
| | <pre># removing Ask for price df = df[df['Price']!='Ask For Price'] df['Price'].unique()</pre> |
| Out[9]: | array(150,000), 14,15,000, 13,25,000, 12,000, 14,15,000, 15,50,000, 14,15,000, 15,20,000, 14,15,000, 15,20,000, 14,15,000, 15,25,000, 14,15,000, 15,25,000, 14,15,000, 15,25,000, 14,15,000, 15,25,000, 14,15,000, 15,25,000, 14,15,000, 15,25,000, 14,15,000, 15,25,000, 14,15,000, 15,25,000, 14,15,000, 15,25,000, 14,15,000, 15,25,000, 14,15,000, 15,25,000, 14,15,000, 15,25,000, 14,15,000, 15,25,000, 14,15,000, 15,25,000, 14,15,000, 15,25,000 |
| In [10]: | # removing commas and add astype function df['Price'] = df["Price"].str.replace(',','').astype(int) kms_driven column cleaning |
| In [11]: Out[11]: | df.head(3) name company year Price kms_driven fuel_type Hyundai Santro Xing XO eRLX Euro III Hyundai 2007 8000 45,000 kms Petrol |
| In [12]: | 1 Mahindra Jeep CL550 MDI Mahindra 2006 425000 40 kms Diesel 3 Hyundai Grand i10 Magna 1.2 Kappa VTVT Hyundai 2014 325000 28,000 kms Petrol # removing kms in kms_driven column df['kms_driven'] = df['kms_driven'].str.split(' ').str.get(0) |
| In [13]: Out[13]: | df . head(5) name company year Price kms_driven fuel_type 0 Hyundai Santro Xing XO eRLX Euro III Hyundai 2007 80000 45,000 Petrol 1 Mahindra Jeep CL550 MDI Mahindra 2006 425000 40 Diesel |
| | 3 Hyundai Grand i10 Magna 1.2 Kappa VTVT Hyundai 2014 325000 Petrol 4 Ford EcoSport Titanium 1.5L TDCi Ford 2014 575000 36,000 Diesel 6 Ford Figo Ford 2012 175000 41,000 Diesel |
| | <pre># removing comma from kms_driven values df['kms_driven'] =df['kms_driven'].str.replace(',','') df.head(5) name company year Price kms_driven fuel_type</pre> |
| In [16]: | 0 Hyundai Santro Xing XO eRLX Euro III Hyundai 2007 80000 45000 Petrol 1 Mahindra Jeep CL550 MDI Mahindra 2006 425000 40 Diesel 3 Hyundai Grand i10 Magna 1.2 Kappa VTVT Hyundai 2014 325000 28000 Petrol 4 Ford EcoSport Titanium 1.5L TDCi Ford 2014 575000 36000 Diesel 6 Ford Figo Ford 2012 175000 41000 Diesel # let see fuel_type unique values df['fuel_type'].unique() |
| Out[16]: In [17]: | <pre>array(['Petrol', 'Diesel', nan, 'LPG'], dtype=object) # checking and removing null values. df = df[~df['fuel_type'].isnull()]</pre> |
| Out[18]: | |
| Out[19]: | df.shape |
| Out[20]: | |
| In [21]: Out[21]: | <pre># checking name df['name'] 0</pre> |
| | Hyundai Grand i10 Magna 1.2 Kappa VTVT Ford EcoSport Titanium 1.5L TDCi Ford Figo Maruti Suzuki Ritz VXI ABS Tata Indica V2 DLE BS III Toyota Corolla Altis Tata Zest XM Diesel |
| | Tata Zest XM Diesel Mahindra Quanto C8 Name: name, Length: 816, dtype: object # checking null values in company df['company'].isnull().sum() |
| Out[22]: In [23]: Out[23]: | <pre>#checking unique values df['company'].unique() array(['Hyundai', 'Mahindra', 'Ford', 'Maruti', 'Skoda', 'Audi', 'Toyota',</pre> |
| | 'Renault', 'Honda', 'Datsun', 'Mitsubishi', 'Tata', 'Volkswagen', 'Chevrolet', 'Mini', 'BMW', 'Nissan', 'Hindustan', 'Fiat', 'Force', 'Mercedes', 'Land', 'Jaguar', 'Jeep', 'Volvo'], dtype=object) #excepting only three words and removing other words df['name'] = df['name'].str.split().str.slice(start = 0, stop = 3).str.join(' ') |
| In [25]: Out[25]: | <pre># comming three words df['name'] 0</pre> |
| | Ford EcoSport Titanium Ford Figo 883 Maruti Suzuki Ritz 885 Tata Indica V2 886 Toyota Corolla Altis 888 Tata Zest XM |
| In [26]: Out[26]: | 889 Mahindra Quanto C8 Name: name, Length: 816, dtype: object df |
| | 1 Mahindra Jeep CL550 Mahindra 2006 425000 40 Diesel 3 Hyundai Grand i10 Hyundai 2014 325000 28000 Petrol 4 Ford EcoSport Titanium Ford 2014 575000 36000 Diesel 6 Ford Figo Ford 2012 175000 41000 Diesel |
| | W Fold Figo Fold 2012 173000 41000 Diesel 883 Maruti Suzuki Ritz Maruti 2011 270000 50000 Petrol 885 Tata Indica V2 Tata 2009 110000 30000 Diesel 886 Toyota Corolla Altis Toyota 2009 300000 132000 Petrol 888 Tata Zest XM Tata 2018 260000 27000 Diesel |
| | 888 Tata Zest XM Tata 2018 260000 27000 Diesel 889 Mahindra Quanto C8 Mahindra 2013 390000 40000 Diesel 816 rows × 6 columns # now change index |
| In [28]: Out[28]: | name company year Price kms_driven fuel_type |
| | 0 Hyundai Santro Xing Hyundai 2007 80000 45000 Petrol 1 Mahindra Jeep CL550 Mahindra 2006 425000 40 Diesel 2 Hyundai Grand i10 Hyundai 2014 325000 28000 Petrol 3 Ford EcoSport Titanium Ford 2014 575000 36000 Diesel 4 Ford Figo Ford 2012 175000 41000 Diesel |
| | III Maruti Suzuki Ritz Maruti 2011 270000 50000 Petrol 812 Tata Indica V2 Tata 2009 110000 30000 Diesel 813 Toyota Corolla Altis Toyota 2009 300000 132000 Petrol |
| | 814 Tata Zest XM Tata 2018 260000 27000 Diesel 815 Mahindra Quanto C8 Mahindra 2013 390000 40000 Diesel 816 rows × 6 columns #cleaned data csv file. |
| In [29]: In [30]: | <pre>#cleaned data csv file. df.to_csv('Cleaned_car_data.csv') df.info() <class 'pandas.core.frame.dataframe'=""> RangeIndex: 816 entries, 0 to 815</class></pre> |
| | Data columns (total 6 columns): # Column Non-Null Count Dtype 0 name 816 non-null object 1 company 816 non-null object 2 year 816 non-null int32 3 Price 816 non-null int32 |
| In [31]: | 4 kms_driven 816 non-null object 5 fuel_type 816 non-null object dtypes: int32(2), object(4) memory usage: 32.0+ KB # convert object into integer. |
| In [32]: | <pre>df['kms_driven'] = df['kms_driven'].astype(int) df.info() <class 'pandas.core.frame.dataframe'=""> RangeIndex: 816 entries, 0 to 815 Data columns (total 6 columns): # Column Non-Null Count Dtype</class></pre> |
| | # Column Non-Null Count Dtype 0 name 816 non-null object 1 company 816 non-null int32 2 year 816 non-null int32 3 Price 816 non-null int32 4 kms_driven 816 non-null int32 5 fuel_type 816 non-null object |
| In [33]: Out[33]: | dtypes: int32(3), object(3) memory usage: 28.8+ KB df.head(5) name company year Price kms_driven fuel_type |
| | Hyundai Santro Xing Hyundai 2007 Booo 4500 Petrol Mahindra Jeep CL550 Mahindra 2006 Hyundai Grand i10 Hyundai Grand i10 Hyundai Grand i10 Ford EcoSport Titanium Ford 5igo Ford 2012 Ford 2012 Ford 5igo Ford 2012 Ford 2012 Ford 5igo Ford 2012 Ford 2012 Ford 5igo Ford 2012 Ford 2013 Ford 5igo Ford 2012 Ford 5igo Ford 2012 Ford 5igo Ford 5igo Ford 2012 Ford 5igo Ford 2012 Ford 5igo Ford 2012 Ford 5igo Ford 2012 Ford 2012 Ford 5igo Ford 2012 Ford 2012 Ford 2012 Ford 5igo Ford 2012 Ford 2012 Ford 2012 Ford 5igo Ford 2012 Ford 2012 |
| In [34]: Out[34]: | df.describe() year Price kms_driven count 816.000000 8.160000e+02 816.000000 |
| | mean 2012.444853 4.117176e+05 46275.531863 std 4.002992 4.751844e+05 34297.428044 min 1995.00000 3.000000e+04 0.000000 25% 2010.00000 1.750000e+05 27000.00000 50% 2013.000000 2.999990e+05 41000.000000 |
| In [35]: | 50% 2013.000000 2.999990e+05 41000.000000 75% 2015.000000 4.912500e+05 56818.500000 max 2019.000000 8.500003e+06 400000.000000 # maximum value |
| In [36]: | df |
| Out[36]: | 0 Hyundai Santro Xing Hyundai 2007 80000 45000 Petrol 1 Mahindra Jeep CL550 Mahindra 2006 425000 40 Diesel 2 Hyundai Grand i10 Hyundai 2014 325000 28000 Petrol 3 Ford EcoSport Titanium Ford 2014 575000 36000 Diesel |
| | 4 Ford Figo Ford J 2012 175000 41000 Diesel 811 Maruti Suzuki Ritz Maruti 2011 270000 50000 Petrol 812 Tata Indica V2 Tata 2009 110000 30000 Diesel 813 Toyota Corolla Altis Toyota 2009 300000 132000 Petrol |
| | 814 Tata Zest XM Tata 2018 260000 27000 Diesel 815 Mahindra Quanto C8 Mahindra 2013 390000 40000 Diesel 816 rows × 6 columns |
| In [37]: In [38]: | <pre># Extracting training data X = df[['name', 'company', 'year', 'kms_driven', 'fuel_type']] y = df['Price'] from sklearn.model_selection import train_test_split X_train, X_test, y_train, y_test = train_test_split(X, y, test_size=0.2)</pre> |
| In [39]: In [40]: | <pre>X_train, X_test, y_train, y_test = train_test_split(X, y, test_size=0.2) from sklearn.linear_model import LinearRegression from sklearn.preprocessing import OneHotEncoder from sklearn.compose import make_column_transformer from sklearn.pipeline import make_pipeline</pre> |
| In [41]: Out[41]: | df name company year Price kms_driven fuel_type 0 Hyundai Santro Xing Hyundai 2007 80000 45000 Petrol |
| | 1 Mahindra Jeep CL550 Mahindra 2006 425000 40 Diesel 2 Hyundai Grand i10 Hyundai 2014 325000 28000 Petrol 3 Ford EcoSport Titanium Ford 2014 575000 36000 Diesel 4 Ford Figo Ford 2012 175000 41000 Diesel |
| | 811 Maruti Suzuki Ritz Maruti 2011 270000 50000 Petrol 812 Tata Indica V2 Tata 2009 110000 30000 Diesel 813 Toyota Corolla Altis Toyota 2009 300000 132000 Petrol 814 Tata Zest XM Tata 2018 260000 27000 Diesel |
| | 815 Mahindra Quanto C8 Mahindra 2013 390000 40000 Diesel 816 rows × 6 columns # use and create object(ohe) ohe = OneHotEncoder() |
| In [43]: Out[43]: | <pre>ohe = OneHotEncoder() #fit values in ohe ohe.fit(X[['name','company','fuel_type']]) OneHotEncoder()</pre> |
| In [44]: In [45]: | <pre># categorical column transform column_trans = make_column_transformer((OneHotEncoder(categories=ohe.categories_), ['name','company','fuel_type']), remainder='passthrough') column_trans ColumnTransformer(remainder='passthrough',</pre> |
| Out[45]: | ColumnTransformer(remainder='passthrough', |
| | volkswagen vento konekt , volvo sao summum], drype=object), |
| | <pre># using linaearregression algorithm model = LinearRegression() # construct a pipeline</pre> |
| In [48]: Out[48]: | <pre>pipe = make_pipeline(column_trans, model) # using simple fit pipe.fit(X_train, y_train) Pipeline(steps=[('columntransformer',</pre> |
| ~1. | ColumnTransformer(remainder='passthrough', transformers=[('onehotencoder', OneHotEncoder(categories=[array(['Audi A3 Cabriolet', 'Audi A4 1.8', 'Audi A4 2.0', 'Audi A6 2.0', 'Audi A8', 'Audi Q3 2.0', 'Audi Q5 2.0', 'Audi Q7', 'BMW 3 Series', 'BMW 5 Series', 'BMW 7 Series', 'BMW X1', 'BMW X1 sDrive20d', 'BMW X1 xDrive20d', 'Chevrolet Beat', 'Chevrolet Beat array(['Audi', 'BMW', 'Chevrolet', 'Datsun', 'Fiat', 'Force', 'Ford', 'Hindustan', 'Honda', 'Hyundai', 'Jaguar', 'Jeep', 'Land', |
| | |
| | ('linearregression', LinearRegression())]) y_pred = pipe.predict(X_test) y_pred array([5.28099773e+05, 2.32220189e+05, 2.34173713e+05, 2.32972000e+05, 3.63964050e+04, 5.84084761e+05, 3.07402753e+05, 3.64084210e+05, |
| | 1.26407959e+06, 2.98251178e+05, 2.25826584-06, 9.3795158e+05, 1.982515726e+05, 1.2116405e+05, 1.2161405e+05, 1.2161405e+05, 2.48730945e+05, 2.88730945e+05, 2.88730945e+05, 2.88730945e+05, 2.8813095e+05, 2.8813095e+05, 2.8813095e+05, 2.2813095e+05, 2.2813095e+05 |
| In [51]: Out[51]: | 4.32721238e+05, 9.37959158e+05, 4.46799675e+05, 1.11366865e+06, 3.45664218e+05, 1.06510242e+05, 2.0250361e+05, 2.84868155e+05, 4.27954810e+04, 2.32033117e+06, 3.86495225e+05, 4.5248094e+05, 3.51375202e+05, 2.08344998e+05, 3.59249745e+05, 4.5248094e+05, 3.51375202e+05, -7.29555601e+04, 4.05840940e+05, 3.16657452e+05, 2.53216752e+05, -7.29555601e+04, 4.05840940e+05, 3.16577452e+05, 3.05434870e+05, 4.31056332e+05, 3.93757797e+05, 2.32036795e+06, 4.08658000e+05, 2.40174795e+05, 4.19000969e+05, 5.39647217e+05, 1.93969024e+05, -1.35113406e+05, 4.19000969e+05, 2.7679824de+03, 2.45408369e+05, -1.97851686e+04, 1.46990530e+05, 6.1785344e+03, 2.45408369e+05, -1.97851686e+04, 1.46990530e+05, 6.1785344e+03, 2.45408369e+05, 3.53183121e+05, 8.78507624e+04]) r2_score(y_test,y_pred) 0.48253481800219367 Conclusion: So the model is ready |