

Assignment 4: Performance Metrics, and Optimisation

Student ID: 300637212 Student Name: Xieji Li

Part 1: Performance Metrics in Regression [30 marks]

Requirements

Based on exploratory data analysis, discuss what preprocessing that you need to do before regression, and provide evidence and justifications.

- Step1. Load Data && split the dataset

Step 1. Load Data

```
df = pd.read_csv("diamonds.csv")
df.drop("Unnamed: 0",axis= 1, inplace = True) # remove the Unnamed col
```

[2] Python

```
# check if there are missing values in data set
print(df.isnull().sum())
```

[293] 0.6s Python

```
... carat    0
    cut      0
    color   0
    clarity  0
    depth    0
    table    0
    x        0
    y        0
    z        0
    price    0
    dtype: int64
```

```
# split the data set into 70% train set and 30% test set
x_train,x_test,y_train,y_test = train_test_split(df[df['price']],train_size= 0.7, random_state=309)
```

[3] Python

- Step 2. Initial Data Analysis

```
df.describe()
```

[4] Python

	carat	depth	table	x	y	z	price
count	53940.000000	53940.000000	53940.000000	53940.000000	53940.000000	53940.000000	53940.000000
mean	0.797940	61.749405	57.457184	5.731157	5.734526	3.538734	3932.799722
std	0.474011	1.432621	2.234491	1.121761	1.142135	0.705689	3989.439738
min	0.200000	43.000000	43.000000	0.000000	0.000000	0.000000	326.000000
25%	0.400000	61.000000	56.000000	4.710000	4.720000	2.910000	950.000000
50%	0.700000	61.800000	57.000000	5.700000	5.710000	3.530000	2401.000000
75%	1.040000	62.500000	59.000000	6.540000	6.540000	4.040000	5324.250000
max	5.010000	79.000000	95.000000	10.740000	58.900000	31.800000	18823.000000

+ Code + Markdown

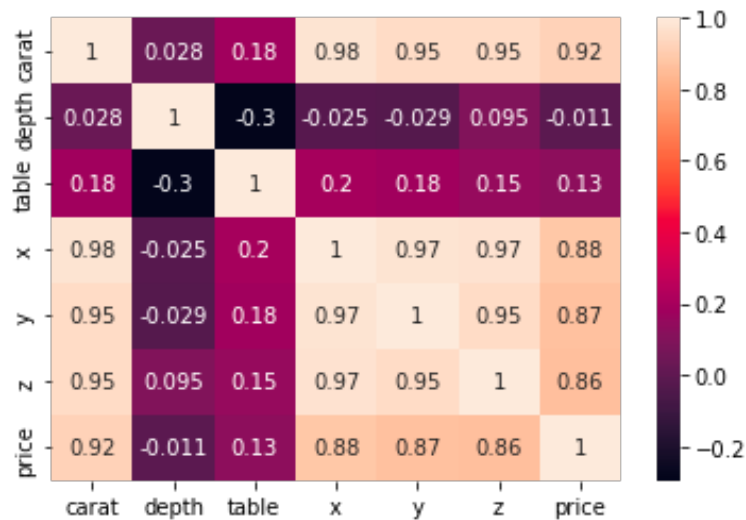
```
df.info()
```

[5] Python

```
<class 'pandas.core.frame.DataFrame'>
RangeIndex: 53940 entries, 0 to 53939
Data columns (total 10 columns):
# Column Non-Null Count Dtype
---
0 carat 53940 non-null float64
1 cut 53940 non-null object
2 color 53940 non-null object
3 clarity 53940 non-null object
4 depth 53940 non-null float64
5 table 53940 non-null float64
6 x 53940 non-null float64
7 y 53940 non-null float64
8 z 53940 non-null float64
9 price 53940 non-null int64
dtypes: float64(6), int64(1), object(3)
memory usage: 4.1+ MB
```

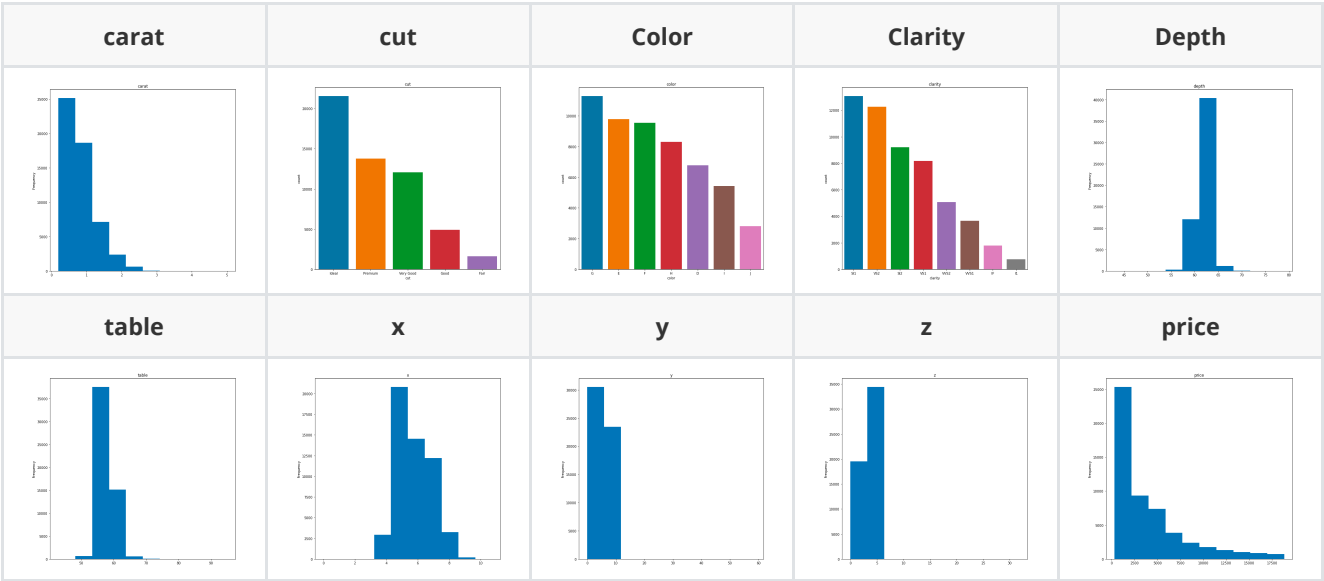
Conclusion: In this stage we can know there are 10 features in this dataset. We need to predict the value of price based on other 9 features. Also, there is no missing value in this dataset.

- correlation analysis
- Heat map



price	
carat	0.921591
x	0.884435
y	0.865421
z	0.861249
price	1.000000

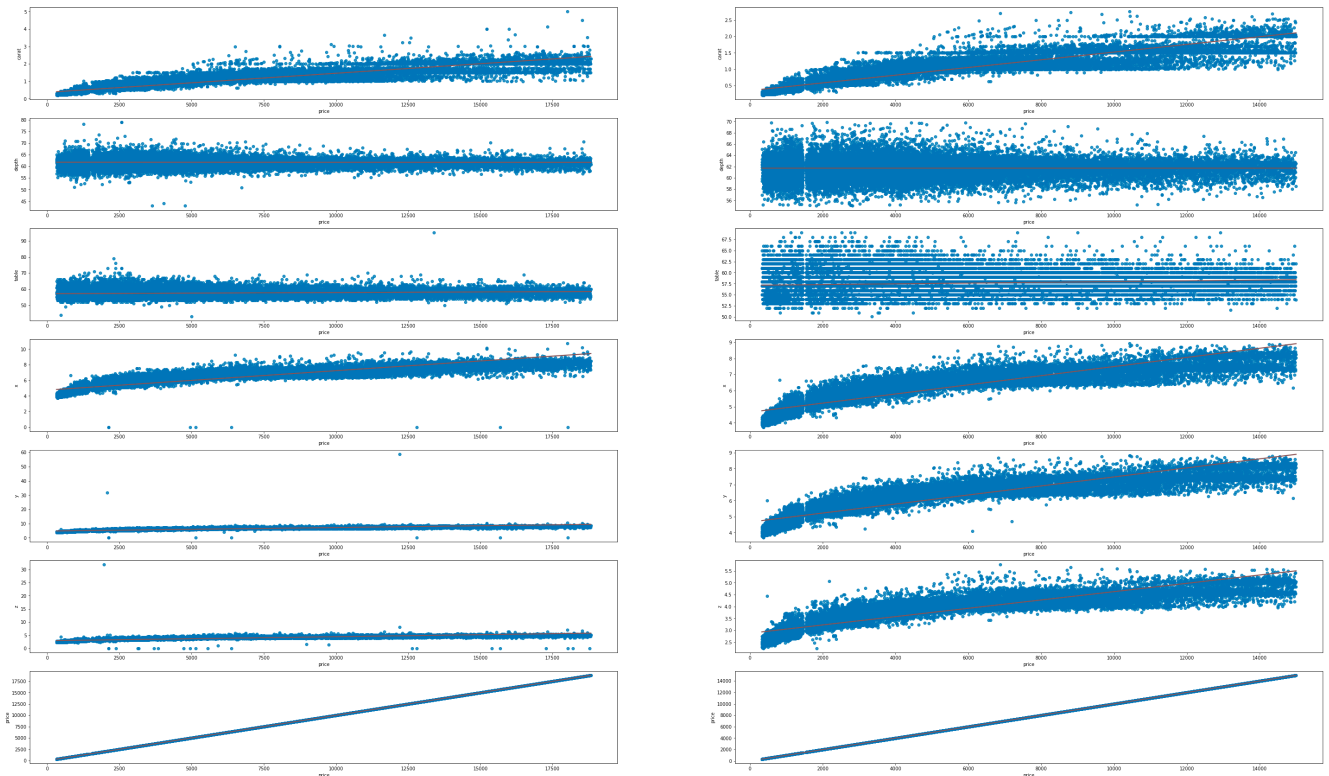
- Step 3. Preprocess Data && Step 4. Exploratory Data Analysis
 - First, use histogram to display features, if the feature is numeric type then plot the hist according to the value of feature. If the feature is category type then plot the hist according to the frequency of the value.



- Remove outliers

1. In carat plot, remove the points - carat > 2.9
2. In depth plot, remove the points - depth > 70 || depth <= 55
3. In table plot, remove the points - table >= 70 || table <= 50
4. In x plot, remove the points - x >= 9 && price >= 15000
5. In y plot, remove the points - y >= 20 || y == 0
6. in z plot, remove the points - z >= 6 || z <= 1

- Right(origin), Left(after removing outliers)



- Encode categorical features based on diamond documentation

- cut

Ideal	Premium	Very Good	Good	Fair
100	80	60	40	20

- color

- One Hot Encode

- clarity

I1	SI2	SI1	VS2	VVS2	VVS1	IF
30	40	50	60	70	80	90

- Standardization

```
# standardization
scaler = StandardScaler()
standard_train = scaler.fit_transform(preprocess_train)
standard_test = scaler.fit_transform(preprocess_test)
```

- Step 5. Build classification (or regression) models using the training data && Step 7. Assess model on the test data.

Model	Parameters	MSE	RMSE	RSE	MAE	excution time
linear regression	positive = True	1647909.22(7)	1283.71(7)	0.13(7)	816.89(7)	0.02s(2)
k-neighbors regression	Default	1339014.10(6)	1157.16(6)	0.12(6)	554.29(6)	1.49s(5)
Ridge regression	Default	2190847.01(9)	1480.15(9)	0.21(8)	848.85(8)	0.004s(1)
decision tree regression	Max_depth = None	825284.43(4)	908.45(4)	0.06(4)	413.07(4)	0.02s(3)
random forest regression	n_estimators = 1000	632325.04(2)	795.19(2)	0.05(2)	336.00(1)	1m50.00s(8)
gradient Boosting regression	Max_depth = none	791343.44(3)	889.57(3)	0.06(3)	401.06(3)	17.83s(7)
SGD regression	Default	2178494.94(8)	1475.97(8)	0.22(10)	864.34(10)	0.20s(4)
support vector regression (SVR)	C=1500	998458.52(5)	999.23(5)	0.09(5)	524.38(5)	3m6.66s(9)
linear SVR	max_iter=50000, C = 5.0, loss = 'squared_epsilon_insensitive', dual = True	2201090.06(10)	1483.61(10)	0.21(9)	848.94(9)	10.78s(6)
multi-layer perceptron regression	max_iter=5000	570093.37(1)	755.05(1)	0.04(1)	391.20(2)	3m22.46s(10)

Discussion

From the table, we can find that multi-layer-preceptron regression, random forest, and gradient boosting regression have a good performance in diamond dataset, but there are some simple model doesn't suitable for this dataset(SGD, linear SVR). Although those simple model take short time in excution stage, they still can't get a great performance. MLP and random forest model takes a long time in excution, but those two model won't be influenced by similar linear features and they will analysis the relationship between features(which help those two model have a better performance than other models).