

Branch: master [kaggle-competitions](#) / README.md[Find file](#) [Copy path](#) **albertovpd** Update README.md

7b4e204 15 minutes ago

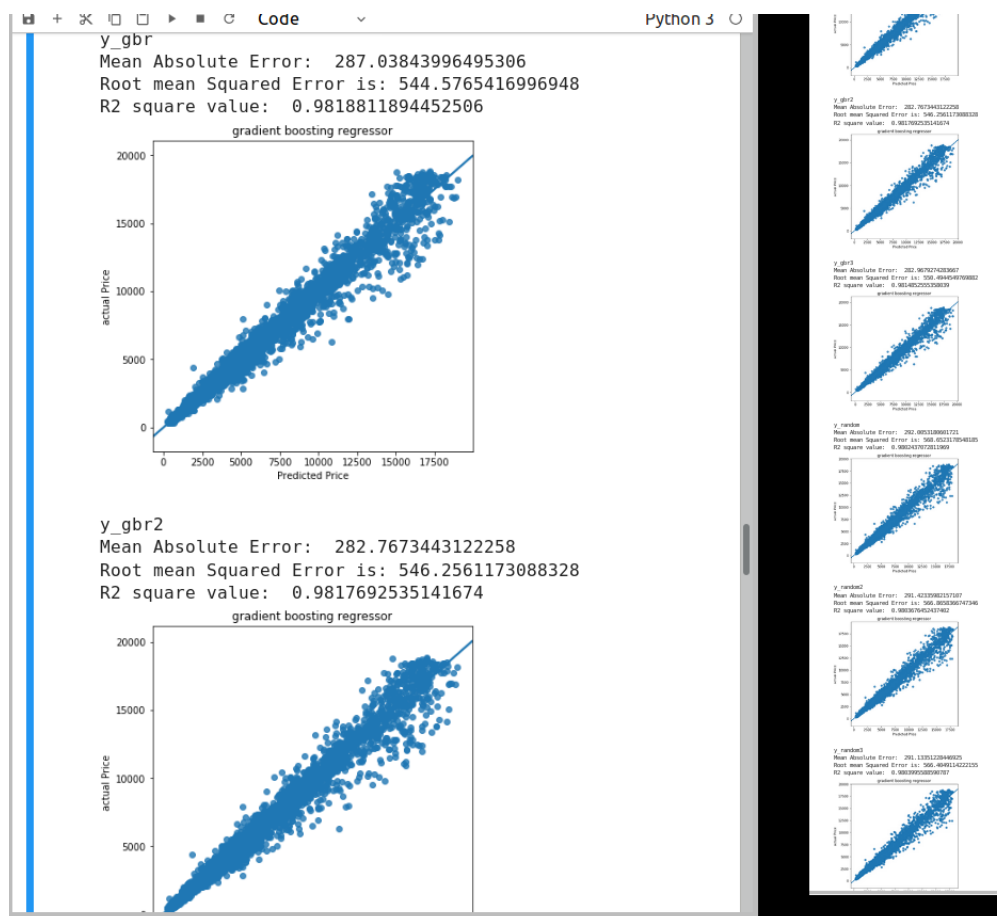
[1 contributor](#)

51 lines (34 sloc) 2.92 KB

[Raw](#) [Blame](#) [History](#)  

1st kaggle-competition

Or "How I spent 2 days learning a lot and getting frustrating low ranks"



The process

- I tried to apply some logic and don't try autoML or similar tools
- Cleaning job
 - Find documentation about the problem
 - Look in every column the type and elemnts

- Correlation matrix to find the proper features
- Turn categorical columns into numeric ones
 - At first I tried with numpy to give slight different weights to elements with numpy. In the end I just gave int types. -py to clean both datasets

```

#my_algorithms.py
import xgboost as xgb

2] > %a
from src.my_algorithms import results_algo

3] > %a
algo=dict(
    "y_gbr": GradientBoostingRegressor(n_estimators=1000).fit(X_train,
    y_train),
    "y_gbr2": GradientBoostingRegressor(n_estimators=2000).fit
    (X_train, y_train),
    "y_gbr3": GradientBoostingRegressor(n_estimators=3000).fit
    (X_train, y_train),
    "y_random": RandomForestRegressor(n_estimators=150).fit
    (X_train, y_train),
    "y_random2": RandomForestRegressor(n_estimators=300).fit
    (X_train, y_train),
    "y_random3": RandomForestRegressor(n_estimators=400).fit
    (X_train, y_train),
    "y_random4": RandomForestRegressor(n_estimators=550).fit
    (X_train, y_train),
    # "y_hist": HistGradientBoostingRegressor().fit(X_train,
    y_train),
    # "y_hist2": HistGradientBoostingRegressor(max_iter=200)
    .fit(X_train, y_train),
    # "y_hist3": HistGradientBoostingRegressor(max_iter=300)
    .fit(X_train, y_train),
    # "y_hist4": HistGradientBoostingRegressor(max_iter=400).fit
    (X_train, y_train)
    # "y_xgboost": xgb.XGBClassifier(random_state=1,
    learning_rate=0.01).fit(x_train, y_train)
)

] > %a

#my_results.py
13 from sklearn.metrics import r2_score
14
15 import numpy as np
16 import pandas as pd
17 import matplotlib.pyplot as plt
18 import seaborn as sns
19
20
21 def results_algo(dictio,X_test,y_test):
22     '''it receives a dictionary with name and
23     features of the model. cheesy gridsearch.
24     it needs X_test and y_test'''
25
26     for n,e in dictio.items():
27         score=e.score(X_test,y_test)
28         prediction = e.predict(X_test)
29
30         print(n)
31         print("Mean Absolute Error: ", mean_absolute_error(y_test, pre
32         diction))
33         print("Root mean Squared Error is: ", mean_squared_error(y_test
34         , prediction))
35         print("R2 square value: ", r2_score(y_test, prediction))
36         plt.figure(figsize= (6, 6))
37         plt.title("gradient boosting regressor")
38         sns.regplot(prediction, y_test)
39         plt.xlabel("Predicted Price")
40         plt.ylabel("actual Price")
41         plt.show()
42         print("")
43     print("Unlock plot in my_algorithms.py")
44
45
46
47

```

- <https://github.com/albertovpd/kaggle-competitions/blob/master/diamonds-datamad1019/output/pics/printing%20function.png>
- Remove outliers from the dataset to have a more accurate model

```

onds-datamad1019 > src > my_maths.py > remove_outlier_col

import numpy as np
def remove_outlier_col(df_in, col_name):
    q1 = df_in[col_name].quantile(0.25)
    q3 = df_in[col_name].quantile(0.75)
    iqr = q3-q1 #Interquartile range
    fence_low = q1-1.5*iqr
    fence_high = q3+1.5*iqr
    df_out = df_in.loc[(df_in[col_name] > fence_low) & (df_in[col_name] < fence_high)]
    return df_out

def remove_outliers(df):
    for e in list(df.columns):
        df=remove_outlier_col(df,"{}".format(e))
    return df

def standarizing(df):
    for col in df.columns:
        df[col] = (df[col] - np.mean(df[col])) / np.std(df[col])
    return df

```

- The price and weight of diamonds are exponentially related, so I give an exponential weight to the "carat" row, and removed columns strongly correlated with this one
- Standarize all my X columns with sklearn
- .py file to check 12 different configurations of algorithms, searching for the min Mean Absolute Error, Root mean Squared Error and R2 square value

alt

- .py file to clean automatically my test dataset exactly like the training dataset (yeah, I could have written a better code, deadline matters)

```
amonds-datamad1019 > src > cleaning_functions.py > ...
4 def cleaning_diamonds2(df):
5     # i already checked the nature of every row and its distribution
6
7
8
9     df.cut=df.cut.replace(to_replace={"Fair":1,"Good":2,"Very Good":3,"Ideal":4,"Premium":5})
10    color=list(df.color.unique())
11    color.sort(reverse=True)
12
13    letters=['D', 'E', 'F','G', 'H', 'I', 'J']
14    weight=list(range(1,8))
15    weight.sort(reverse=True)
16    color_dict=dict(zip(letters,weight))
17    for e in df.color:
18        for x,y in color_dict.items():
19            if e==x:
20                df.color.replace(e,y,inplace=True)
21
22    nomenclature=["I1","SI2","SI1","VS2","VS1","VVS2","VVS1","IF"]
23    weight=list(range(1,9))
24    weight.sort(reverse=True)
25    clarity_dict=dict(zip(nomenclature,weight))
26
27    for e in df.clarity:
28        for x,y in clarity_dict.items():
29            if e==x:
30                df.clarity.replace(e,y,inplace=True)
31
32    df['volume'] = df['x'] * df['y'] * df['z']
33    df.drop(["x","y","z"], axis=1,inplace=True)
34
35    # sorting rows randomly
36    #df.carat=(np.log(df.carat)).abs()
37    #df.carat=df.carat*math.exp(2)
38
39    #standarizing columns
40    #columns=['carat', 'cut', 'color', 'clarity', 'depth', 'table']
41    #for e in columns:
```

- Train the model with my best configuration (the first picture)

<https://raw.githubusercontent.com/albertovpd/kaggle-competitions/master/diamonds-datamad1019/output/pics/final%20results.png>

- Get the predicted "y"
- Submit results

Conclusions

I really wanted to try my skills as scientist in my first "fight" with ML. I miserably lost and learnt about it.

- 1. The most important task is to study the field and clean your dataset properly.
- 2. Do not give weight to columns without knowlegde. If you are going to do random stuff, you have way better tools for that purpose.
- 3. Do not use Scikit GridSearchCV randomly.
- 4. Avoid the bad use of H2O's autoML and tools like that. It Depends on the size and type of dataset. Nevertheless, if you don't have the proper time to study the field your working with, go back to 2 and do it. Also very useful when you think you can't optimize your result anymore.

Those are my principles, and if you don't like them... Well, I have others.