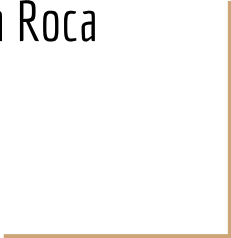




Group 1 Final Project

Kyle Schneider, Carolina Roca



Topic Selection, Source & Reasoning

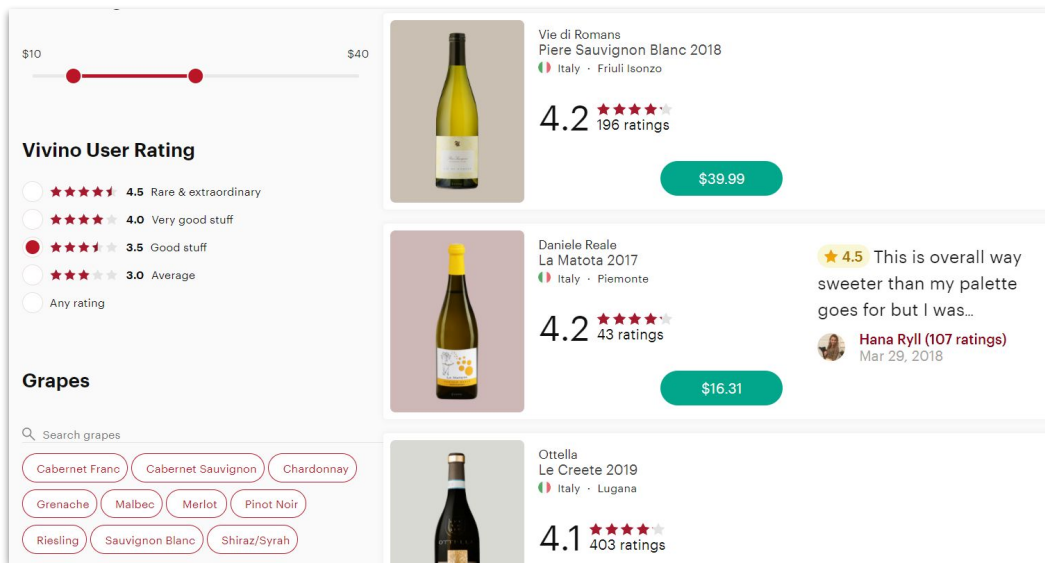
Wine Not?

The topic we chose to pursue was wine as we both enjoy it, and Kyle works in the CPG industry. The dataset we chose originated from Vivino and was hosted on Kaggle for extraction. It contains information about specific wines, where they come from, and how the public has rated them.

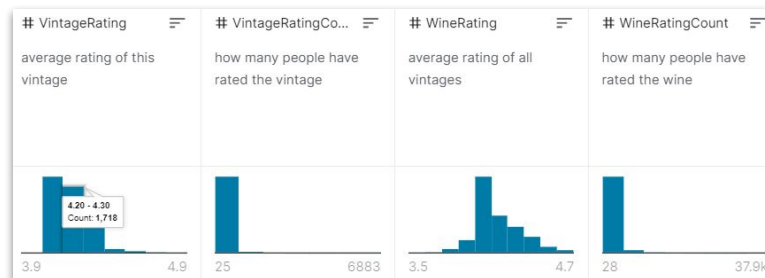
As wine consumers, we understand that vintage can greatly affect the quality, price, and taste of finished wine. As data scientists and machine learning practitioners, we want to understand if the features available in the dataset can lead to an accurate prediction of rating.

Our Data Source

vivino



kaggle



~5K rows of wine vintage ratings along with pricing, geographical, and other data

Question to answer...

Using the available data from Vivino wine ratings, can we accurately predict a given vintage's rating?

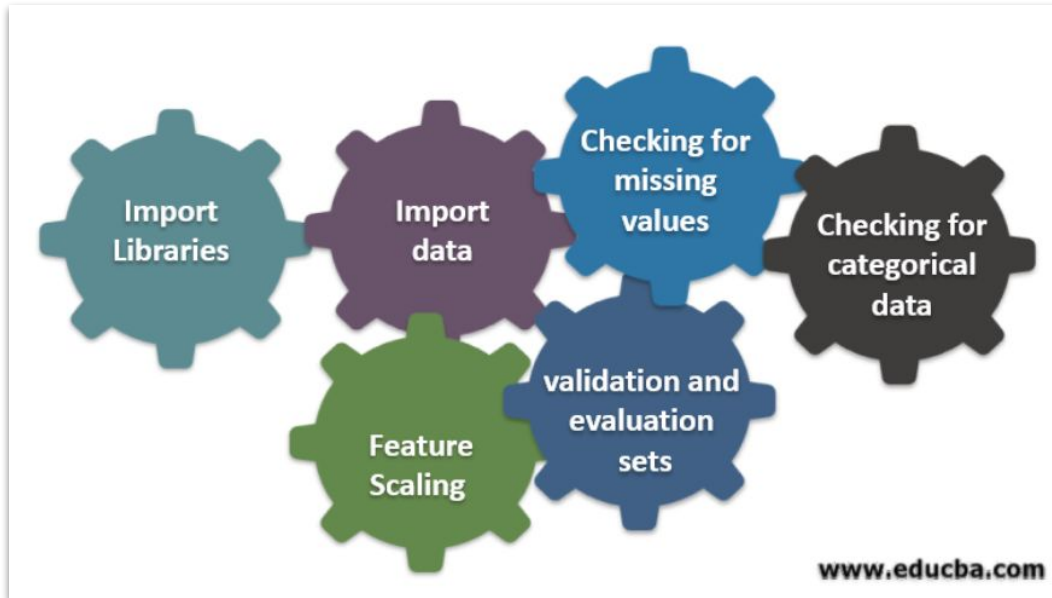


Understanding our data...

Tableau snips

Building the database...

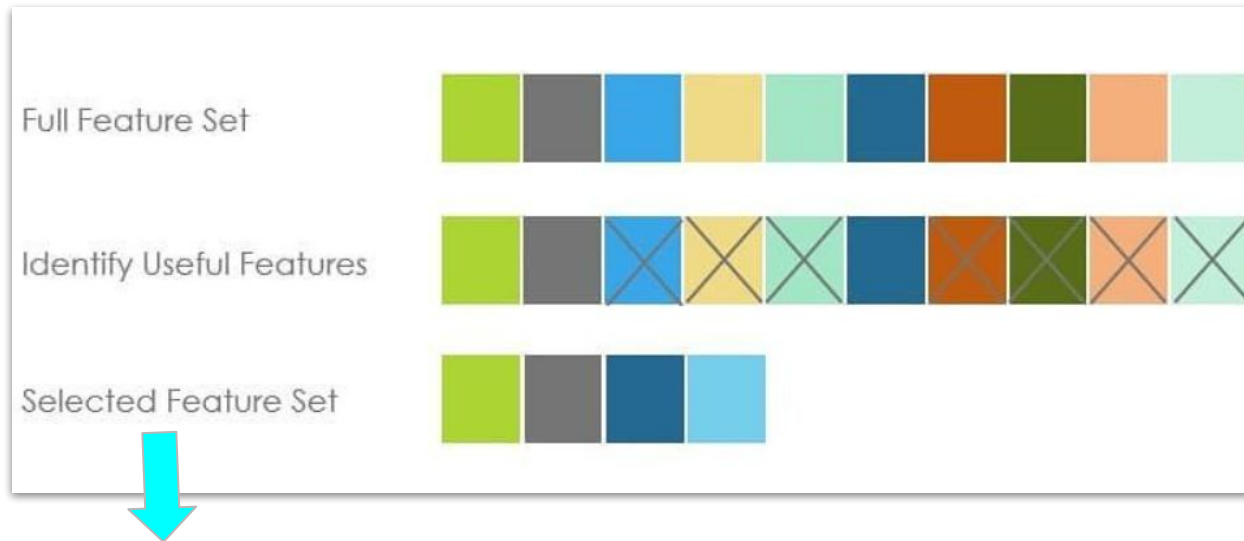
Machine Learning Model - Preprocessing



```
1 #Import initial dependencies
2
3 import pandas as pd
4 import pyodbc
5
6 #Establish connection to SQL database
7
8 conn = pyodbc.connect('Driver={ODBC Driver 17 for SQL Server};'
9                       'Server=tcp:group1-owner-nu.database.windows.net,1433;'
10                      'Database=final-project;'
11                      'Persist Security Info=False;'
12                      'Uid=GROUPDB1NU;'
13                      'Pwd=NU02282021!;'
14                      'MultipleActiveResultSets=False;'
15                      'Encrypt=Yes;'
16                      'TrustServerCertificate=No;'
17                      'Connection Timeout=30;')
18
19
20 cursor = conn.cursor()
21
22 #Read table from SQL
23
24 initial_df = pd.read_sql("SELECT * FROM dbo.final_table", conn)
25
26 initial_df.head()
```

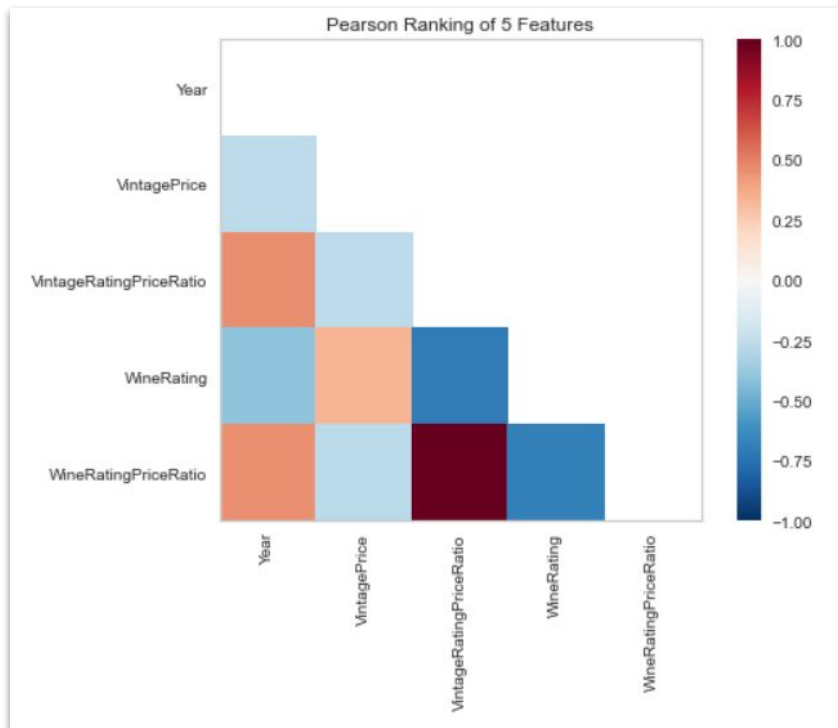
```
1 from sklearn.compose import ColumnTransformer
2 from sklearn.compose import make_column_transformer
3 from sklearn.compose import make_column_selector
4
5 ct = make_column_transformer(
6     (StandardScaler(), make_column_selector(dtype_include=np.float64)),
7     (OrdinalEncoder(), make_column_selector(dtype_include=object))
8 )
9 ct.fit_transform(X)
```

Machine Learning Model - Feature Selection



FullName	Winery	Year	Vintage Price	Vintage Rating Price Ratio	Wine Rating	Wine Rating Price Ratio
----------	--------	------	---------------	----------------------------	-------------	-------------------------

Description of Analysis Phase



```
# visualize pipeline
from sklearn import set_config
set_config(display="diagram")
model

Pipeline(steps=[('columntransformer',
                  ColumnTransformer(transformers=[('standardscaler',
                                                  StandardScaler(),
                                                  'WineRatingPriceRatio')],
                                      remainder='passthrough'),
                  ('standardscaler',
                   StandardScaler(),
                   'WineRating'),
                  ('ordinalencoder',
                   OrdinalEncoder(),
                   'Year'),
                  ('linearregression',
                   LinearRegression())])])
```

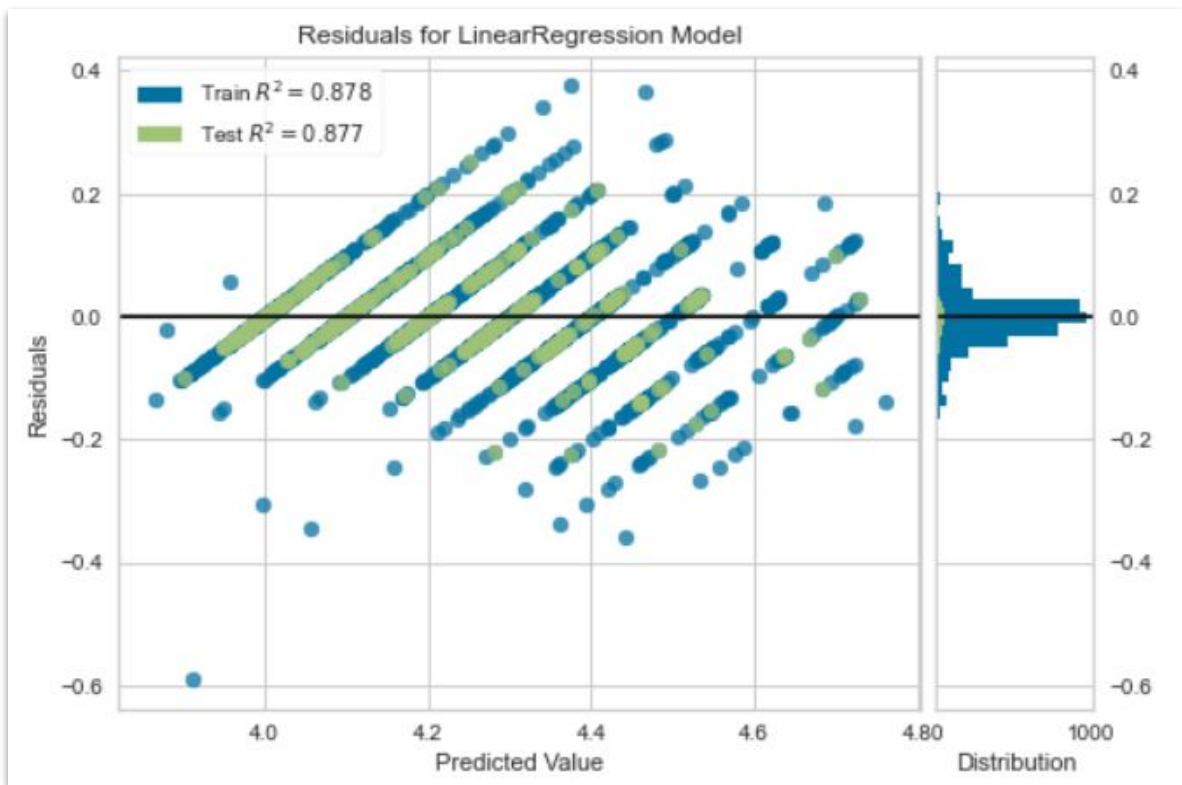
```
1 X_train, X_test, y_train, y_test = train_test_split(X, y, random_state=42)
2 model = make_pipeline(ct, LinearRegression())
3 model.fit(X, y)
4 model.score(X_test, y_test)
```

0.8640005437048607

```
1 X_encoded = ct.fit_transform(X)
2 y_encoded = y
3 lr = LinearRegression()
4 lr.fit(X_encoded, y_encoded)
5 lr.score(X_encoded, y_encoded)
```

0.877734987463609

Results

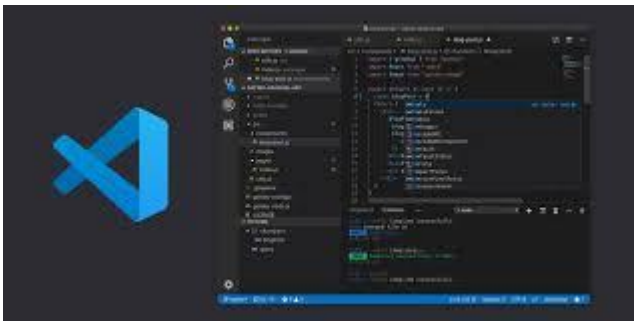
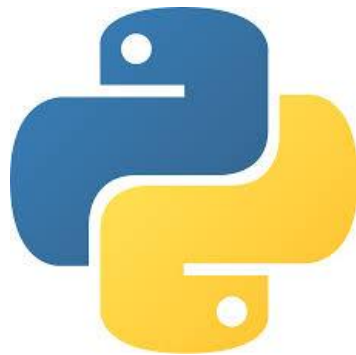


```
1 from sklearn.metrics import mean_squared_error
2
3 predictions = model.predict(X_test)
4 MSE = mean_squared_error(y_test, predictions)
5 r2 = model.score(X_test, y_test)
```

```
1 print(f"MSE: {MSE}, R2: {r2}")
```

MSE: 0.003395387554488735, R2: 0.877407772833871

Tools and Resources



Appendix