

# A Data mining approach to predict human wine taste preferences

A [dataset](#) with white and red wino verde samples is considered

The **Feature vector**:

- fixed acidity
- volatile acidity
- citric acid
- residual sugar
- chlorides
- free sulfur dioxide
- total sulfur dioxide
- density
- pH
- sulphates
- alcohol

Final Prediction:

- Quality

The original [kaggle](#) uses MR(multiple regression), SVN (Support Vector Machines) and NN (Neural Networks) for data mining.

Categorization of the problem:

- By Input: We have labelled data, which means we will do supervised learning
- By Output: We want to predict a class, which means we will do classification

In this project on the other hand we will utilize Random Forests for classification.

We achieve an accuracy of 96.22% without any model improvement.

## Game Plan

The ML Pipeline:

- Question and the required data
- Acquire the data
- Data Analysis
- Preparation -> data for the ML model
- Train the Model
- Test the Model
- Evaluate the Model
- Interpret the Model and report results visually and numerically
- Adjusting the Model

```
In [1]: import pandas as pd
import numpy as np
import random
import collections
import matplotlib.pyplot as plt
import seaborn as sns
import time
from pprint import pprint

from sklearn.model_selection import train_test_split, RandomizedSearchCV, GridSearchCV
from sklearn.ensemble import RandomForestClassifier
from sklearn.metrics import confusion_matrix, accuracy_score

# Can be used to export a decision tree in DOT format
from sklearn.tree import export_graphviz

import pydot
# Visualizing decision trees
from IPython.display import Image

sns.set()
```

## Data Analysis: Cleaning & Preparing

Data analysis usually occupies 80% of the time of any Data Scientist and Machine Learning Engineer. Without this crucial phase all the state-of-the-art models offered by sklearn would be completely in vain.

- Tidy Data
- Missing Data
- Outliers
- Class Imbalance Problem
- Feature Correlation

```
In [2]: features = pd.read_csv("data/winequality-white.csv")
features.head()
```

```
Out [2]:
```

	fixed acidity	volatile acidity	citric acid	residual sugar	chlorides	free sulfur dioxide	total sulfur dioxide	density	pH	sulphates	alcohol	quality		
0								7.027	0.3620	7.0045	45.170	1.0013	30.45	8.86

## Tidy Data

As can be clearly noticed this dataset is messy: **Multiple variables are stored in one column.**

Therefore the first step would be to [tidy](#) the data.

```
In [3]: # Get the columns
features_list = features.columns[0].split(" ")
features_list = [feature.lstrip(" ").rstrip(" ") for feature in features_list]

# Get the values
values = np.array(features[features.columns])
new_values = []
for row in values:
    new_values.append([list(map(float, row[0].split(" "))))
new_values = np.array(new_values)

# Form the new dataframe
features = pd.DataFrame(data=new_values, columns=features_list)
print("We have {} records and {} features".format(features.shape))

features.sample(3)
```

We have 4898 records and 12 features

```
Out [3]:
```

	fixed acidity	volatile acidity	citric acid	residual sugar	chlorides	free sulfur dioxide	total sulfur dioxide	density	pH	sulphates	alcohol	quality
1071	8.3	0.22	0.38	14.8	0.054	32.0	126.0	1.0002	3.22	0.50	9.7	5.0
1617	6.2	0.20	0.49	1.6	0.065	17.0	143.0	0.9937	3.22	0.52	9.2	6.0
649	7.1	0.26	0.34	14.4	0.067	35.0	189.0	0.9986	3.07	0.53	9.1	7.0

After these little tweaks we can tell that:

- Each feature forms a column
- Each observation/sample forms a row

## Identifying anomalies/missing data

- Missing Data:** Imputing or dropping values. Generally speaking, if an observation should have been made but it was not made, then you do not drop the sample. If it could not be made then you drop it.
- Outliers:** Can be dropped

One way to spot different anomalies as well as missing data is to compute a summary statistics.

The following will give us a descriptive statistics

```
In [4]: features.describe()
```

```
Out [4]:
```

	fixed acidity	volatile acidity	citric acid	residual sugar	chlorides	free sulfur dioxide	total sulfur dioxide	density	pH	sulphates	alcohol	quality
count	4898.000000	4898.000000	4898.000000	4898.000000	4898.000000	4898.000000	4898.000000	4898.000000	4898.000000	4898.000000	4898.000000	
mean	8.654788	0.278241	0.334192	6.391415	0.045772	35.308085	138.360657	0.994027	3.188267	0.4896		
std	0.843868	0.100795	0.121020	5.072058	0.021848	17.007137	42.498065	0.002091	0.151001	0.1141		
min	3.800000	0.060000	0.000000	0.600000	0.009000	2.000000	9.000000	0.987110	2.720000	0.2200		
25%	4.300000	0.210000	0.270000	1.700000	0.036000	23.000000	108.000000	0.991723	3.090000	0.4100		
50%	6.800000	0.260000	0.320000	5.200000	0.043000	34.000000	134.000000	0.993740	3.180000	0.4700		
75%	7.300000	0.320000	0.390000	9.900000	0.050000	46.000000	167.000000	0.996100	3.280000	0.5500		
max	14.200000	1.100000	1.660000	65.800000	0.346000	289.000000	440.000000	1.038960	3.820000	1.0800		

The dataset contains no missing values. However, it looks like we might have some outliers.

Let's visualize the data so that we can better tell this.

```
In [5]: # Set the style
plt.style.use('fivethirtyeight')

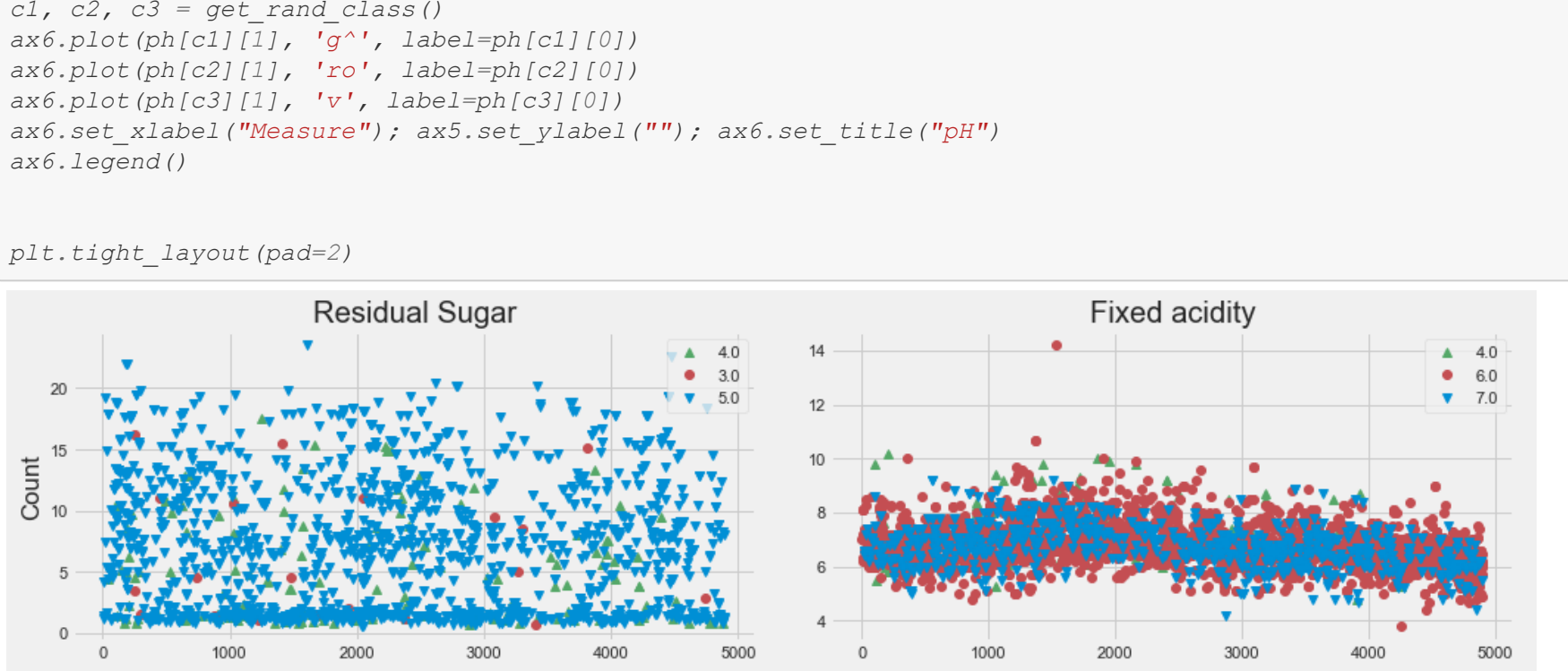
# Plotting layout
fig, (ax1, ax2, ax3) = plt.subplots(nrows=1, ncols=3, figsize=(15,5))

# residual sugar
ax1.plot(features['residual sugar'], 'ro')
ax1.set_xlabel("Measure"); ax1.set_ylabel("Count")
ax1.set_title("Residual Sugar")

# free sulfur dioxide
ax2.plot(features['free sulfur dioxide'], 'ro')
ax2.set_xlabel(""); ax2.set_ylabel("Count")
ax2.set_title("Free Sulfur Dioxide")

# total sulfur dioxide
ax3.plot(features['total sulfur dioxide'], 'ro')
ax3.set_xlabel(""); ax3.set_ylabel("Count")
ax3.set_title("Total Sulfur Dioxide")

plt.show()
```



We can certainly see outliers for these attributes. One way to deal with outliers is to remove and if we don't have many we can do so.

The outlier in the first figure belongs to quality 6.0 for which we do have plenty of samples. However, the two other outliers belong to 3.0 for which we do not have too many samples. Dropping them means dropping 10% of the data for that kind of quality.

Let's set a threshold of 125 for outliers in Free Sulfur Dioxide and 300 for outliers in Total Sulfur Dioxide in order to further investigate this.

```
In [6]: features[features['free sulfur dioxide']>300]
```

```
Out [6]:
```

	fixed acidity	volatile acidity	citric acid	residual sugar	chlorides	free sulfur dioxide	total sulfur dioxide	density	pH	sulphates	alcohol	quality
325	7.5	0.27	0.31	5.80	0.057	131.0	313.0	0.99460	3.18	0.59	10.5	5.0
1417	8.6	0.55	0.35	15.55	0.057	35.5	366.5	1.00010	3.04	0.63	11.0	3.0
1931	7.1	0.49	0.22	2.00	0.047	146.5	307.5	0.99240	3.24	0.37	11.0	3.0
2127	9.1	0.33	0.38	1.70	0.062	50.5	344.0	0.99580	3.10	0.70	9.5	5.0
2654	6.9	0.40	0.22	5.95	0.081	76.0	303.0	0.99705	3.40	0.57	9.4	5.0
4745	6.1	0.26	0.25	2.90	0.047	289.0	440.0	0.99314	3.44	0.64	10.5	3.0

```
In [7]: features[features['total sulfur dioxide']>125]
```

```
Out [7]:
```

	fixed acidity	volatile acidity	citric acid	residual sugar	chlorides	free sulfur dioxide	total sulfur dioxide	density	pH	sulphates	alcohol	quality
325	7.5	0.27	0.31	5.8	0.057	131.0	313.0	0.99460	3.18	0.59	10.5	5.0
1931	7.1	0.49	0.22	2.0	0.047	146.5	307.5	0.99240	3.24	0.37	11.0	3.0
2334	7.5	0.230	0.35	17.8	0.058	126.0	212.0	1.00241	3.44	0.43	8.9	5.0
3050	6.2	0.255	0.24	1.7	0.039	138.5	272.0	0.99452	3.53	0.53	9.6	4.0
4745	6.1	0.260	0.25	2.9	0.047	289.0	440.0	0.99314	3.44	0.64	10.5	3.0

We can see that it seems that quality=3 seems to have more outliers in both of these attributes and since we don't have much data of it (20 samples) in total it does not make sense to drop the outliers for these two features.

So we will stick with dropping a single outlier for the feature Residual Sugar.

```
In [8]: features = features.drop(index=features[features['residual sugar']>40].index, axis=0)
```

## Class-Imbalance Problem

Class Imbalance problem is a really nasty issue which is often difficult to tackle. Therefore a thumb up rule here is: data, data and some more data. When we have made sure that the phenomenon comes as a result of nature itself (from the way things are) and not from low amount of data we can proceed with investigating it.

```
In [9]: print("Classes are indeed imbalanced")
print(round(100*(features['quality'].value_counts()/features.shape[0]), 2))

Classes are indeed imbalanced
6.0 45.0
5.0 30.0
7.0 16.0
8.0 4.0
4.0 3.0
3.0 0.0
9.0 0.0
Name: quality, dtype: float64 2
```

We can see that a whopping 75% of the data belong to only two classes!

However, if our classes are separated nicely then the class-imbalance is not really a problem. Although there are tools and techniques to check this here we will just visualize for a couple of attributes and classes at a time.

```
In [10]: def get_rand_class():
    """Get 3 random classes from a total of 6. """
    rand = []
    while True:
        if len(rand) == 3:
            return rand
        r = random.randint(0,6)
        if r not in rand:
            rand.append(r)

In [11]: fig, (ax1, ax2), (ax3, ax4), (ax5, ax6) = plt.subplots(nrows=3, ncols=2, figsize=(15,12))

# Feature: residual sugar
rs = list(features['residual sugar'].groupby(features['quality']))

# Randomly pick 3 classes to display
c1, c2, c3 = get_rand_class()
ax1.plot(rs[c1][1], 'g', label=rs[c1][0])
ax1.plot(rs[c2][1], 'ro', label=rs[c2][0])
ax1.plot(rs[c3][1], 'v', label=rs[c3][0])
ax1.set_xlabel(""); ax1.set_ylabel("Count"); ax1.set_title("Residual Sugar")
ax1.legend()

# Feature: Fixed acidity
fa = list(features['fixed acidity'].groupby(features['quality']))

# Randomly pick 3 classes to display
c1, c2, c3 = get_rand_class()
ax2.plot(fa[c1][1], 'g', label=fa[c1][0])
ax2.plot(fa[c2][1], 'ro', label=fa[c2][0])
ax2.plot(fa[c3][1], 'v', label=fa[c3][0])
ax2.set_xlabel(""); ax2.set_ylabel("Count"); ax2.set_title("Fixed acidity")
ax2.legend()

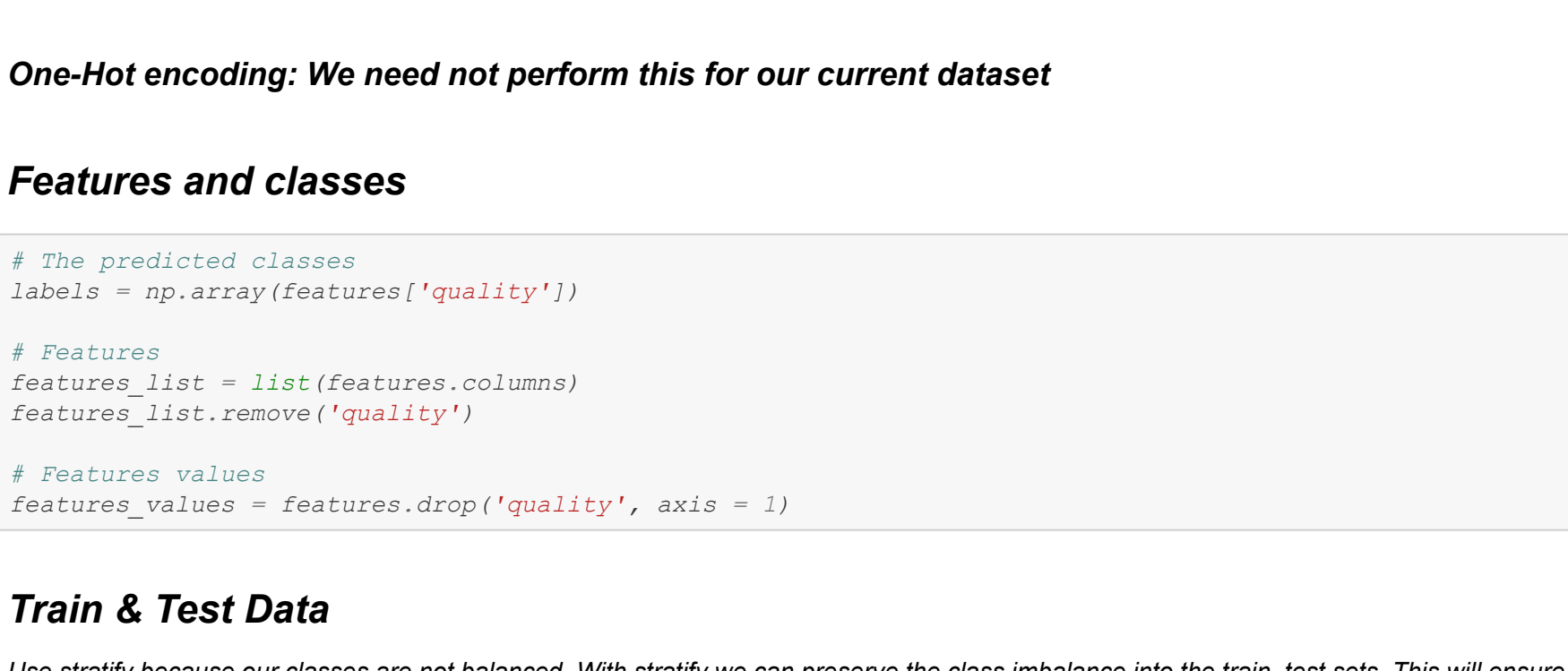
# Feature: citric acid
ca = list(features['citric acid'].groupby(features['quality']))
c1, c2, c3 = get_rand_class()
ax3.plot(ca[c1][1], 'g', label=ca[c1][0])
ax3.plot(ca[c2][1], 'ro', label=ca[c2][0])
ax3.plot(ca[c3][1], 'v', label=ca[c3][0])
ax3.set_xlabel(""); ax3.set_ylabel("Count"); ax3.set_title("Citric Acid")
ax3.legend()

# Feature: total sulfur dioxide
tsd = list(features['total sulfur dioxide'].groupby(features['quality']))
c1, c2, c3 = get_rand_class()
ax4.plot(tsd[c1][1], 'g', label=tsd[c1][0])
ax4.plot(tsd[c2][1], 'ro', label=tsd[c2][0])
ax4.plot(tsd[c3][1], 'v', label=tsd[c3][0])
ax4.set_xlabel(""); ax4.set_ylabel("Count"); ax4.set_title("Total Sulfur Dioxide")
ax4.legend()

# Feature: residual sugar
rs2 = list(features['residual sugar'].groupby(features['quality']))
c1, c2, c3 = get_rand_class()
ax5.plot(rs2[c1][1], 'g', label=rs2[c1][0])
ax5.plot(rs2[c2][1], 'ro', label=rs2[c2][0])
ax5.plot(rs2[c3][1], 'v', label=rs2[c3][0])
ax5.set_xlabel("Measure"); ax5.set_ylabel("Count"); ax5.set_title("Residual Sugar")
ax5.legend()

# Feature: pH
ph = list(features['pH'].groupby(features['quality']))
c1, c2, c3 = get_rand_class()
ax6.plot(ph[c1][1], 'g', label=ph[c1][0])
ax6.plot(ph[c2][1], 'ro', label=ph[c2][0])
ax6.plot(ph[c3][1], 'v', label=ph[c3][0])
ax6.set_xlabel("Measure"); ax6.set_ylabel("Count"); ax6.set_title("pH")
ax6.legend()

plt.tight_layout(pad=2)
```



We can tell from this plots that the features are not separated at all (some exceptions are there but they do not change the rule). These means that the class-imbalance problem for this dataset is real!

Among the options that we are left with is undersampling or oversampling but for now we will just stick with the data as is and see how our model fares.

We basically do nothing regarding this.

## Feature correlation

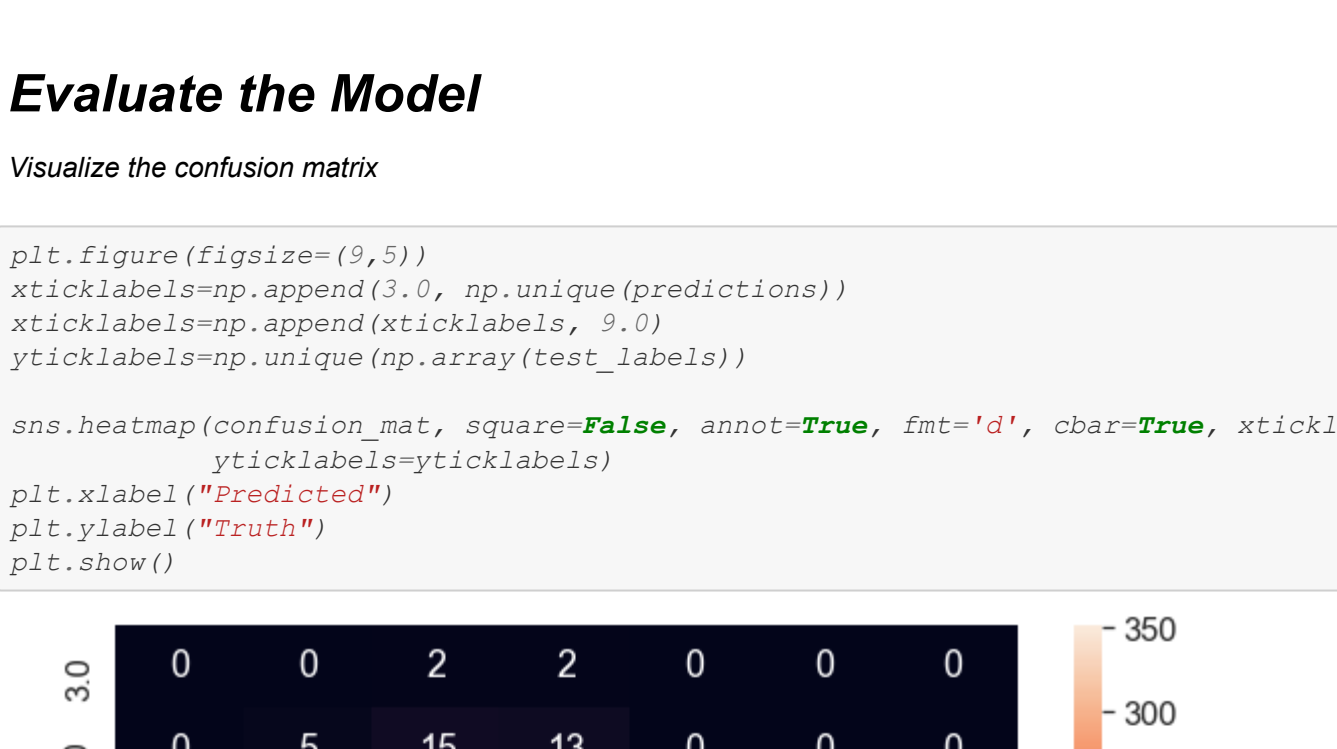
A high correlation between features is almost always not feasible for the ML algorithms. Therefore we always aim to feature engineer correlated features.

## Correlation with the quality variable

Let's see important variables in predicting the quality by checking the feature's correlation with it.

```
In [12]: # Drop the quality and plot with pandas
features.drop("quality", axis=1).corrwith(features['quality']).plot.bar(figsize=(10,5),
                                title="Correlation with the class variable", rot=45, fontsize=15)

plt.show()
```



So, taking the correlation into account, we can see that alcohol is the most important feature in our feature vector for predicting the quality of the wine. Note that even variables that are negatively correlated with it (like density) are indeed important because they could be contributors of a lower quality.

Citric acid and free sulfur dioxide seem to have the lowest correlation.

alcohol and density seem to have the highest correlation.

## Correlation Matrix: Correlation of features with each-other.

Variables to utilize when using sns heatmap

- vmask:** Achoring the color map to values. If not specified, it will infer from the data
- linewidths:** The widths of the lines dividing the cells
- cbar\_kwds:** in the below case we specify the height of the stripe showing the values
- annot:** allow annotations inside the boxes
- fmt:** formatting the annotations

```
In [13]: sns.set(style='white', font_scale=1.5)

# Compute the correlation matrix
# The correlation matrix is an array where the main diagonal separates two identical triangles
corrMatrix = features.corr()

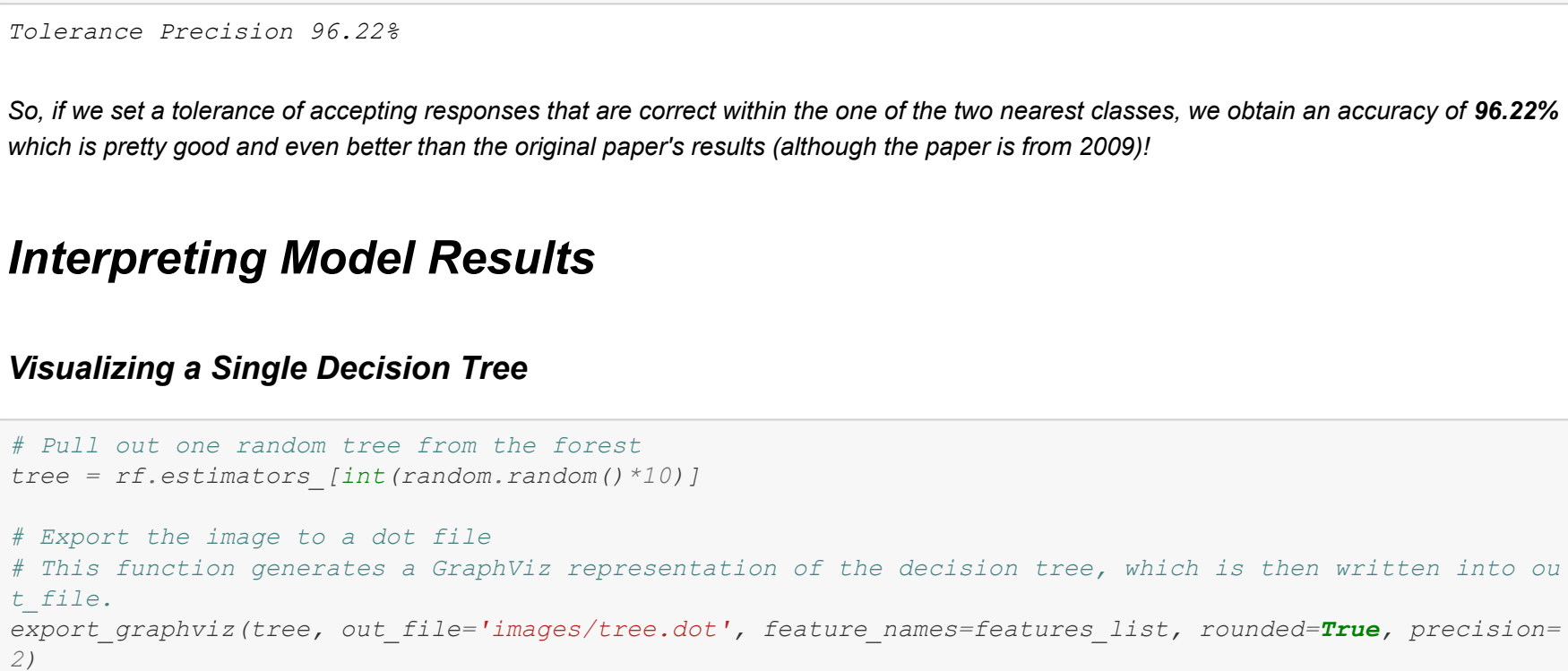
# Create a mask for the upper triangle so that we can ignore it later when building the heatmap
# When we pass this mask to the heatmap function it will generate no data for the upper triangle
mask = np.zeros_like(corrMatrix, dtype=np.bool)
# Get the indices of the upper-triangle of arr: triu_indices
mask[np.triu_indices_from(mask)] = True

# Generate a custom diverging colormap
# Colormap for the different values of the correlation matrix
cmap = sns.diverging_palette(220, 10, as_cmap=True)

plt.figure(figsize=(3, 7))
plt.title("Correlation Matrix of Features")

# Draw the heatmap
sns.heatmap(corrMatrix, square=True, mask=mask, cmap=cmap, center=0, linewidths=2.0, cbar_kwds={"shrink": 0.6})

plt.show()
```



The features seem to be somewhat correlated with each-other.

- alcohol and density are negatively correlated to a rather significant amount.
- density and residual sugar are positively correlated to a significant amount as well

density and alcohol are also highly correlated with the quality. An increase of alcohol increases the quality but an increase in density decreases it. Therefore it doesn't seem such a good idea to merge them in one single feature.

## Prepare the data for the machine learning model

- One-Hot Encoding
- Train and Features
- Class and Features
- Scaling/Standardisation if needed

One-Hot encoding: We need not perform this for our current dataset

## Features and classes

```
In [14]: # The predicted classes
labels = np.array(features['quality'])

# Features
features_list = list(features.columns)
features_list.remove('quality')

# Features values
features_values = features.drop('quality', axis=1)
```

## Train & Test Data

Use stratify because our classes are not balanced. With stratify we can preserve the class imbalance into the train, test sets. This will ensure a better model.

```
In [15]: train_features, test_features, train_labels, test_labels = train_test_split(features_values, labels,
                                        train_size=0.2, random_state=3, stratify=labels)

print("Shape of our train and test samples")
print("Shape of train features:", train_features.shape)
print("Shape of train labels:", train_labels.shape)
print("Shape of test features:", test_features.shape)
print("Shape of test labels:", test_labels.shape)

Shape of our train and test samples
Shape of train features: (3917, 11)
Shape of train labels: (3917,)
Shape of test features: (980, 11)
Shape of test labels: (980,)
```

## Train the Model

Yes! With the state-of-the-art scikit-learn algorithm, training our model looks this the following!

```
In [16]: # Instantiate Model
# Set a random_state in order to get consistent results
rf = RandomForestClassifier(n_estimators=100, random_state=3)

# Train the model
rf.fit(train_features, train_labels);
```

## Make Predictions on Test Data

One way to do this is using the confusion matrix and then calculating the Micro-Precision. We use Micro- and not Macro- precision because Macro-Precision is not flexible to class-imbalance problem which we do have in this case.

MicroP =  $\frac{TP}{TP+FP}$

```
In [17]: predictions = rf.predict(test_features)
confusion_mat = confusion_matrix(test_labels, predictions)
microP = sum(np.diag(confusion_mat)) / sum(sum(confusion_mat))

print("Micro precision {:.8%}.format(round(microP*100, 2)))

# Calculate with library
print("Precision {:.8%}.format(round(100*accuracy_score(test_labels, predictions), 2)))

Micro precision 68.06%
Precision 68.06%
```

## Evaluate the Model

Visualize the confusion matrix

```
In [18]: plt.figure(figsize=(5,5))
yticklabels=np.append(0, np.unique(predictions))
xticklabels=np.append(0, np.unique(predictions))
yticklabels=np.array(test_labels)
yticklabels=np.array(test_labels)
sns.heatmap(confusion_mat, square=False, annot=True, fmt='d', cbar=True, xticklabels=xticklabels,
            yticklabels=yticklabels)
plt.xlabel("Predicted")
plt.ylabel("Truth")
plt.show()
```



It looks like our model is not doing such a good job since its accuracy is 68.06%.

It completely misses out on classes with quality 9 and 3. This makes sense because we do not have so many samples for both of them. For class 9 specifically we had 1 training and 1 testing samples, which is terrible. It also misclassifies some values from the other classes.

As already investigated a little bit earlier a reason for this might be due the class-imbalance problem.

However, the paper reports that, when admitting only the correct classified classes, the overall accuracy was 62.4%. So already our model is faring better than the SVM and NN version's of the paper. The performance is then substantially increased when they set a tolerance of accepting responses that are correct within the one of the two nearest classes, obtaining an accuracy of 89.0%.

Let us calculate this kind of tolerance accuracy as described by the paper and observe the results of the model. Incidentally, this is the reason why we also manually calculated the micro precision using the confusion matrix earlier.

```
In [20]: sumTp = sum(sum(confusion_mat))
sumFP = 0
for i in range(len(confusion_mat)):
    for j in range(len(confusion_mat)):
        # element in main diagonal
        if i == j:
            sumTp += confusion_mat[i][j]
        # element around main diagonal
        elif (i - j) < 4 & (i < 3):
            sumTp += confusion_mat[i][j]
        elif (j - i) < 4 & (i > 0):
            sumTp += confusion_mat[i][j]

microP_tol = sumTp / sumTp + sumFP
print("Tolerance Precision {:.8%}.format(round(microP_tol*100, 2)))

Tolerance Precision 96.22%
```

So, if we set a tolerance of accepting responses that are correct within the one of the two nearest classes, we obtain an accuracy of 96.22% which is pretty good and even better than the original paper's results (although the paper is from 2009!).

## Interpreting Model Results

### Visualizing a Single Decision Tree

```
In [21]: # Pull out one random tree from the forest
tree = rf.estimators_[int(random.random()*10)]

# Export the image to a dot file
# This function generates a graphviz representation of the decision tree, which is then written into a file.
export_graphviz(tree, out_file='images/tree.dot', feature_names=features_list, rounded=True, precision=2)

# Use dot file to create a graph
(graph, ) = pydot.graph_from_dot_file('images/tree.dot')

# Write graph to a png file
graph.write_png('images/tree.png')
```

```
In [22]: print("The depth of the tree is", tree.get_max_depth())

The depth of the tree is 26
```

### Visualizing a smaller Decision Tree

```
In [23]: mini_rf = RandomForestClassifier(n_estimators=100, max_depth=3, random_state=3)
mini_rf.fit(train_features, train_labels)
predictions2 = mini_rf.predict(test_features)

print("The accuracy of the mini tree is {:.8%}.format(round(100*accuracy_score(test_labels, predictions2), 2)))

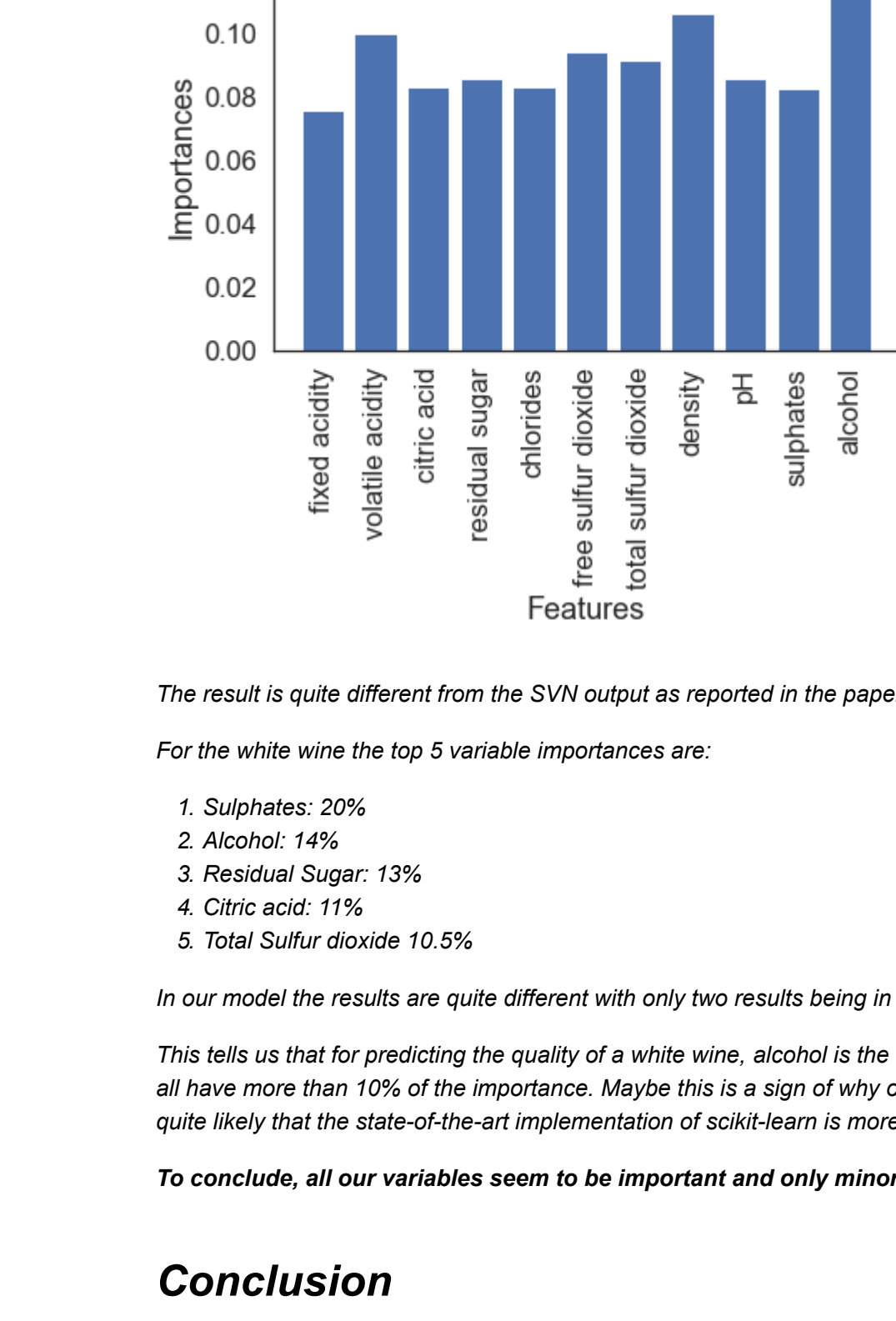
The accuracy of the mini tree is 53.67%
```

Understandably the precision is decreasing. Let's visualize the tree so that we can better understand what's going on. Incidentally decision trees have the advantage of being readily interpretable at a time when many ML algorithms appear to be black boxes.

```
In [24]: mini_tree = mini_rf.estimators_[1]
export_graphviz(mini_tree, out_file='images/mini_tree
```



```
[27]: # Visualize variable importances
plt.bar(features_list, importances)
plt.xticks(rotation='vertical')
plt.xlabel("Features"); plt.ylabel("Importances"); plt.title("Variable Importances"); plt.show()
```



The result is quite different from the SVN output as reported in the paper.

For the white wine the top 5 variable importances are:

- 1. Sulphates: 20%
- 2. Alcohol: 14%
- 3. Residual Sugar: 13%
- 4. Citric acid: 11%
- 5. Total Sulfur dioxide 10.5%

In our model the results are quite different with only two results being in the top five for both models.

This tells us that for predicting the quality of a white wine, alcohol is the most important feature, followed by volatile acidity and density which all have more than 10% of the importance. Maybe this is a sign of why our model is doing better. After all the paper is from 2009 and it is quite likely that the state-of-the-art implementation of scikit-learn is more robust than the algorithm used by the paper.

To conclude, all our variables seem to be important and only minor differences exist between them.

## Conclusion

The model that we build can classify with 96.22% the quality of the wine (with a tolerance of 1 point in both directions). This is much better than the 89% accuracy of the paper (from 2009).

Next we will take a look at possible ways to adjust and improve the model.