

# Practical Machine Learning

## Exercise 3 Spring 2024

Lisa Stafford

February 20, 2024

# Abstract

Model selection is a common problem for machine learning architects where one selects the model that does the best job of generalizing and most appropriate for the given problem. Those selecting the model should choose the best balance between under- and over-fitting. When selecting a model, one must examine the value of different predictive methods to identify the one that best fits observable data. There are two main factors when choosing the most appropriate model in machine learning. One is the reason for choosing it, and the other is the model's performance. Reasons for choosing a model should be based on the available data set, and the problem task.

## Introduction

In order to learn and determine which models are most effective, we used a wine data set obtained from the University of Irvine [1]. Training is performed on each subset of red and white wine within the data set. Several models are trained implementing libraries directly from scikitlearn [4] and then training and testing performance scores for each model and wine subsets are taken. Model selection determinations can then be made based on the performance scores.

## Dataset Description

Our Wine Dataset is actually comprised of two different wine subsets of data. While both are related to variants of [1] "Portuguese Vhinho Verde wine" the data is actually split in two subsets - one red subset and one white subset. Each subset contains 11 features and 1 label for wine quality. All features are continuous float values. Within the model the "wine quality" labels are recognized as discrete categorical integer values ranging from 3 to 9 for the white wines, and 3 and 8 for the red wines. White wine contains a total of 4898 total data instances, while red wine has a total of 1598 data instances as shown in the following table images. Figure 1 (left) displays feature values, statistics, and total instances for white wine, while Figure 1 (right) displays feature values, statistics and total instances available for red wine.

```
RangeIndex: 4898 entries, 0 to 4897
Data columns (total 12 columns):
#   Column              Non-Null Count  Dtype
---  -
0   fixed acidity        4898 non-null   float64
1   volatile acidity     4898 non-null   float64
2   citric acid          4898 non-null   float64
3   residual sugar       4898 non-null   float64
4   chlorides            4898 non-null   float64
5   free sulfur dioxide  4898 non-null   float64
6   total sulfur dioxide 4898 non-null   float64
7   density              4898 non-null   float64
8   pH                  4898 non-null   float64
9   sulphates            4898 non-null   float64
10  alcohol              4898 non-null   float64
11  quality              4898 non-null   int64
dtypes: float64(11), int64(1)
memory usage: 459.3 KB
Categorical columns: []
Numerical columns: ['fixed acidity', 'volatile acidity', 'citric acid', 'residual sugar', 'chlorides', 'free sulfur dioxide', 'total sulfur dioxide', 'density', 'pH', 'sulphates', 'alcohol', 'quality']
```

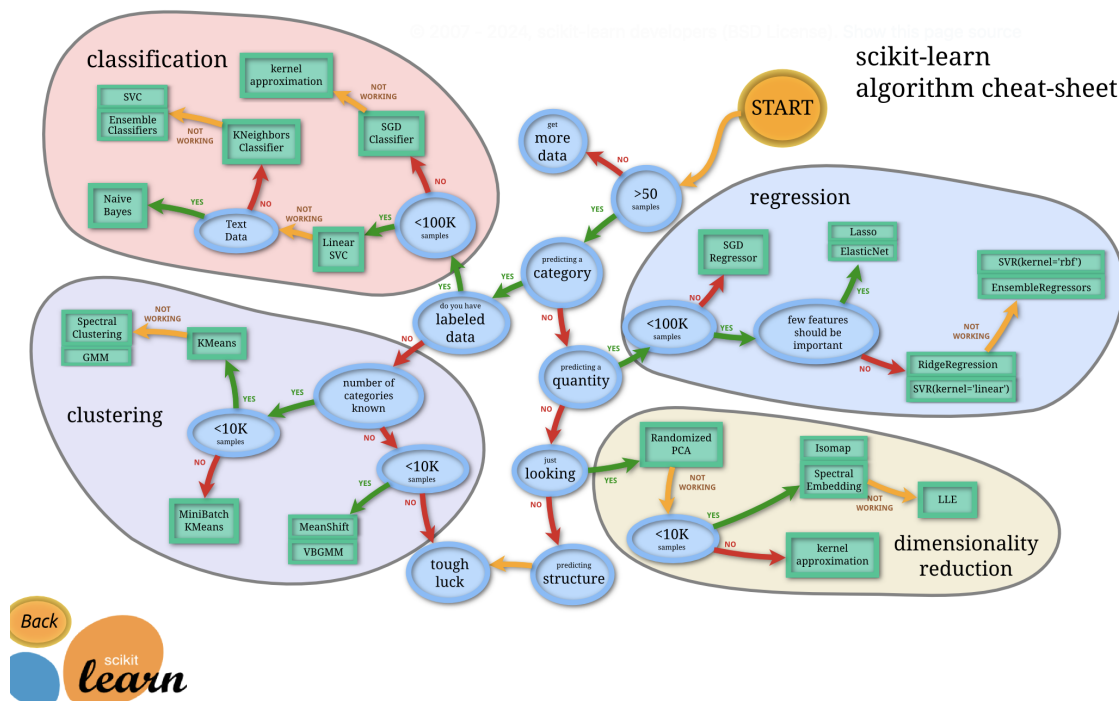
```
RangeIndex: 1599 entries, 0 to 1598
Data columns (total 12 columns):
#   Column              Non-Null Count  Dtype
---  -
0   fixed acidity        1599 non-null   float64
1   volatile acidity     1599 non-null   float64
2   citric acid          1599 non-null   float64
3   residual sugar       1599 non-null   float64
4   chlorides            1599 non-null   float64
5   free sulfur dioxide  1599 non-null   float64
6   total sulfur dioxide 1599 non-null   float64
7   density              1599 non-null   float64
8   pH                  1599 non-null   float64
9   sulphates            1599 non-null   float64
10  alcohol              1599 non-null   float64
11  quality              1599 non-null   int64
dtypes: float64(11), int64(1)
memory usage: 150.0 KB
Categorical columns: []
Numerical columns: ['fixed acidity', 'volatile acidity', 'citric acid', 'residual sugar', 'chlorides', 'free sulfur dioxide', 'total sulfur dioxide', 'density', 'pH', 'sulphates', 'alcohol', 'quality']
```

Data for both red and white wine datasets is somewhat "pre-cleaned". The feature values are all continuous and both red and white datasets contain no missing data values. However it should be noted that these features are not evenly distributed. The data values are somewhat bell-shaped with a larger number of data instances with labels closer to the mean. As described on the dataset website [1] both "dataset... classes are ordered and [unbalanced] so there are many more normal wines than excellent or poor ones.... [also], both datasets have 11 physiochemical features... and a sensory output label, ('quality')". Because Standard Scaling and Transformation and Log Transformation worked best as a data transformation method in Exercise 2, we use the same code to perform these two specific data cleaning methods on the data again.

## Experimental Setup

We are looking at the wine quality dataset, and trying to run a small number of machine learning algorithms on the dataset to determine the best machine learning algorithm for the dataset, where the "best" algorithm may be a set of algorithms.

We know that this particular dataset is supervised and the problem is a classification problem with a known number of categories. Therefore we know that the best ML models will be in that realm. To make an initial decision, this "cheat sheet" decision chart from sklearn [4] was followed:



A decision was made to evaluate a number of popular Classification ML Models/Algorithms, since it was hard to be sure which would be more effective. The machine learning algorithms that fit best with this particular dataset are as follows and therefore we will utilize these listed in our initial ML model selection:

1. Random Forest Classifier
2. Linear Support Vector Machine
3. Classic Support Vector Machine
4. K-Neighbors Classifier
5. Stochastic Gradient Descent Classifier
6. Neural Network using MultiLayer Perceptron

While neural networks are not listed in the above decision chart, they are frequently used when training ML models due to their flexibility for use in a wide variety of ML problems, on different data types, and generally high level of performance. Therefore, they were included in this exercise.

Training is performed on each subset of red and white wine within the data set and data is parsed into 80% for training, and 20% for testing. Several models are trained implementing libraries directly from scikitlearn [4]. While we have a number of different model options, we specify limited hyper-parameters for each model, and instead frequently proceed with training using the default hyper-parameters specified through scikitlearn in most cases.

Next, training and testing performance scores for each model and wine subsets are taken. Model selection determinations and evaluation is then to be made based on ten individual and unrelated train and test runs to get performance scores of each model. Performance values are averaged after 10 independent model train/test runs.

## Results

Machine learning performance results were never very highly performative, but resulted as follows. The highest machine learning test performance was found to be different models, depending on the wine subset. The first graphic displays the results of each model before performing standard scaling and log scaling on specific data features. The following shows performance of each model after performing data scaling. As is notable, data scaling has a significant effect on certain ML model performance.

Pre Scaling:

	Model_Type	White_Training_Score	White_Test_Score	Red_Training_Score	Red_Test_Score
0	Linear SVC	0.538030	0.491837	0.580141	0.628125
1	Polynomial SVC	0.465033	0.431633	0.509773	0.528125
2	K-Nearest Neighbor	1.000000	0.580612	1.000000	0.618750
3	SGD Classifier	0.456611	0.417347	0.499609	0.556250
4	MultiLayer Perceptron	0.548239	0.472449	0.613761	0.628125

Post Scaling:

	Model_Type	White_Training_Score	White_Test_Score	Red_Training_Score	Red_Test_Score
0	Linear SVC	0.537264	0.493878	0.587177	0.631250
1	Polynomial SVC	0.631445	0.522449	0.741204	0.628125
2	K-Nearest Neighbor	1.000000	0.641837	1.000000	0.675000
3	SGD Classifier	0.515722	0.485510	0.701955	0.600625
4	MultiLayer Perceptron	0.643543	0.539694	0.701955	0.600625

As shown in the charts, highest performance for the white wine subset was obtained using K-Nearest Neighbor for both pre and post scaling on data. KNN achieves 58% and 64% accuracy for performance on test data. The red wine data subset resulted in the same performance for both Linear Support Vector Machine and NN using a Multi Layer Perceptron with 62.8% accuracy for performance on test data when running on non-scaled data however KNN performs far better after data scaling with a 67.5% accuracy rate.

With the above changes to performance happening after certain methods of scaling, we try 2 different scaling methods and then one more run combining these two scaling methods and note resulting changes to performance for each model. When looking at the average, and best scores for each model type, the following were noted. On the left are best scores per runs, and on the right are average scores for each category after 5 runs per data scaling type:

	Model_Type	Wine_Type	Training_Score	Test_Score	Data_Scale
0	Random Forest	red	1.0	0.74375	standard
1	Random Forest	white	1.0	0.666327	log
2	Linear SVC	red	0.587177	0.63125	standard
3	Linear SVC	white	0.535733	0.494898	logstd
4	Polynomial SVC	red	0.741204	0.628125	standard
5	Polynomial SVC	white	0.631445	0.522449	standard
6	K-Nearest Neighbor	red	1.0	0.696875	log
7	K-Nearest Neighbor	white	1.0	0.641837	standard
8	SGD Classifier	red	0.582486	0.61875	log
9	SGD Classifier	white	0.531138	0.508163	logstd
10	MultiLayer Perceptron	red	0.59656	0.6375	standard
11	MultiLayer Perceptron	white	0.59367	0.539796	standard

	Model_Type	Wine_Type	Training_Score	Test_Score	Data_Scale
3	Random Forest	red	1.000000	0.724167	average
7	Random Forest	white	1.000000	0.653741	average
11	Linear SVC	red	0.590826	0.629167	average
15	Linear SVC	white	0.535988	0.494218	average
19	Polynomial SVC	red	0.698202	0.622917	average
23	Polynomial SVC	white	0.603709	0.504422	average
27	K-Nearest Neighbor	red	1.000000	0.681250	average
31	K-Nearest Neighbor	white	1.000000	0.627211	average
35	SGD Classifier	red	0.544436	0.546667	average
39	SGD Classifier	white	0.492700	0.463741	average
43	MultiLayer Perceptron	red	0.575137	0.583542	average
47	MultiLayer Perceptron	white	0.564608	0.512313	average

Certain models are more performant when data is scaled in specific ways but all models perform best when using some form of data scaling. Random Forest appears to be less sensitive to scaling but performed best with standard scaling. Linear SVC performs best with standard scaling on average. Polynomial SVC always performed best with standard scaling. KNN performed well regardless of scaling but did best with a log scaler on average. SGD performed best with a combination scaler, and MLP performed best with standard scaling. Average performance was best with the Random Forest method across all data scaling but KNN was a close second on all data.

The last process within training validation was to re-perform all model training using kfold cross validation on the sets into 5 partitions, resulting in an 80% training and 20% testing datasets. Results were very similar to those from previous runs though more "averaged out".

Best Results for these models was similar in performance to the more "average performance" results from the prior runs:

	<b>Model_Type</b>	<b>Wine_Type</b>	<b>Training_Score</b>	<b>Test_Score</b>	<b>Data_Scale</b>
<b>0</b>	Random Forest	red	1.0	0.7075	standard
<b>1</b>	Random Forest	white	1.0	0.673469	none
<b>2</b>	Random Forest KF	red	1.0	0.699171	log
<b>3</b>	Random Forest KF	white	1.0	0.687421	logstd
<b>4</b>	Linear SVC	red	0.576314	0.625	none
<b>5</b>	Linear SVC	white	0.534985	0.515102	standard
<b>6</b>	Linear SVC KF	red	0.597248	0.595374	log
<b>7</b>	Linear SVC KF	white	0.538587	0.533075	logstd
<b>8</b>	Polynomial SVC	red	0.753962	0.62	standard
<b>9</b>	Polynomial SVC	white	0.630003	0.515918	standard
<b>10</b>	Poly SVC KF	red	0.772514	0.609136	logstd
<b>11</b>	Poly SVC KF	white	0.627807	0.534301	standard
<b>12</b>	K-Nearest Neighbor	red	1.0	0.675	log
<b>13</b>	K-Nearest Neighbor	white	1.0	0.64	standard
<b>14</b>	KNN KF	red	1.0	0.669165	standard
<b>15</b>	KNN KF	white	1.0	0.662308	logstd
<b>16</b>	SGD Classifier	red	0.557965	0.6125	standard
<b>17</b>	SGD Classifier	white	0.536346	0.52	logstd
<b>18</b>	SGD KF	red	0.563322	0.568474	standard
<b>19</b>	SGD KF	white	0.519854	0.519192	logstd
<b>20</b>	MultiLayer Perceptron	red	0.609675	0.63	standard
<b>21</b>	MultiLayer Perceptron	white	0.564933	0.539592	standard
<b>22</b>	MLP KF	red	0.620856	0.603509	standard
<b>23</b>	MLP KF	white	0.575133	0.557168	standard

The one thing I could more easily take away from the runs with kfold cross validation is that the data scale methods used in these runs are likely going to be more consistently the best data scaling to result in the best performance for the underlying model.

## References

- [1] A. Asuncion, D. Newman, UCI Machine Learning Repository, University of California, Irvine (2007). Obtained from <https://archive-beta.ics.uci.edu/dataset/186/wine+quality>.

- [2] C. Harris, K. Millman, S. van der Walt, Array programming with NumPy. *Nature* 585, 357–362 (2020). DOI: 10.1038/s41586-020-2649-2. <https://numpy.org/doc/stable/reference/generated/numpy.isnan.html>
- [3] M. Waskom, (2021). seaborn: statistical data visualization. *Journal of Open Source Software*, 6(60), 3021, <https://doi.org/10.21105/joss.03021>.
- [4] Scikit-learn: Machine Learning in Python, Pedregosa et al., *JMLR* 12, pp. 2825-2830, 2011.