

ESE417 Final Project

By Josh Hart, Ritika Tejwani, and Sam Schwartz

Table of Contents

INTRODUCTION.....	1
METHODS.....	1
Preprocessing.....	1
Principal Component Analysis.....	2
Feature Selection.....	2
Logistic Regression Model.....	3
Random Forest Classifier Model.....	4
Support Vector Machine Model.....	4
RESULTS AND ANALYSIS.....	4
Logistic Regression Model.....	4
Random Forest Classifier Model.....	5
Support Vector Machine Model.....	6
CONCLUSIONS.....	7
Division of Work.....	7
References.....	7
APPENDIX.....	8

INTRODUCTION

Machine learning is a field of computer science that focuses on utilizing algorithms and datasets to make predictions about data a computer would otherwise be unable to without being explicitly told. In this project we designed, implemented, and tested three different machine learning models on a high-dimension, non-linearly separable dataset. The red wine dataset has 12 attributes, 11 of which are input variables: fixed acidity, volatile acidity, citric acid, residual sugar, chlorides, free sulfur dioxide, total sulfur dioxide, density, pH, sulfates, and alcohol; one of which is the target variable: quality. We decided to apply a logistic regression, random forest classifier, and support vector machine algorithm on the red wine dataset due to their varying levels of complexity and unique approaches to classification with the goal of optimal classification performance on this dataset in mind. For each approach, we also implemented various measures to improve performance such as data scaling and exhaustive hyperparameter evaluation. Ultimately, we came to the conclusion that a tuned random forest classifier model clearly performs the strongest on this dataset.

METHODS

After importing the dataset using pandas' `.read_csv()` method, we first looked at the dataset itself. Based on the printed dataset, we see that there are 1,599 data instances and 12 attributes. The majority of the features are continuous variables, with the target variable (quality) being the only categorical variable.

Preprocessing

When approaching the red wine dataset, we first implemented preprocessing methods to make the dataset feasible to train a model on. First, for each feature, we checked if there were any null instances in the data, and if so, replaced them with a value of 0 using the `.fillna(0)` method. Having instances of null in the dataset adversely affects a model's performance and overall accuracy by creating bias because the null values are treated as a value. We ensured that every column (feature) had 0 instances of null values.

# of Null Instances for each feature:		<class 'pandas.core.frame.DataFrame'> RangeIndex: 1599 entries, 0 to 1598 Data columns (total 12 columns):			
		#	Column	Non-Null Count	Dtype
fixed acidity	0	0	fixed acidity	1599 non-null	float64
volatile acidity	0	1	volatile acidity	1599 non-null	float64
citric acid	0	2	citric acid	1599 non-null	float64
residual sugar	0	3	residual sugar	1599 non-null	float64
chlorides	0	4	chlorides	1599 non-null	float64
free sulfur dioxide	0	5	free sulfur dioxide	1599 non-null	float64
total sulfur dioxide	0	6	total sulfur dioxide	1599 non-null	float64
density	0	7	density	1599 non-null	float64
pH	0	8	pH	1599 non-null	float64
sulphates	0	9	sulphates	1599 non-null	float64
alcohol	0	10	alcohol	1599 non-null	float64
quality	0	11	quality	1599 non-null	int64
dtype: int64		dtypes: float64(11), int64(1) memory usage: 150.0 KB			

Next, we created histograms of each feature to visualize the distribution of values for each variable. We see again that all of the variables are continuous except for the target variable, quality. Quality is rated on a scale from 0-10 but the only classes that appear in the dataset are 3-8. Fixed acidity, volatile acidity,

citric acid, residual sugar, chlorides, free sulfur dioxide, total sulfur dioxide, sulfates, and alcohol are right-skewed distributions while density, pH, and quality follow approximately normal distributions. (Reference appendix for histograms)

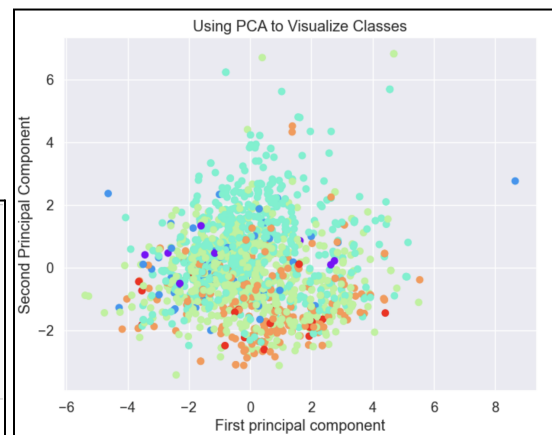
We then segmented the dataset into input variables (X) and a target variable (y). X includes the 11 input variables (fixed acidity, volatile acidity, citric acid, residual sugar, chlorides, free sulfur dioxide, total sulfur dioxide, sulfates, alcohol, density, and pH) while y is quality.

Principal Component Analysis

Next, we implemented Principal Component Analysis (PCA), which allowed us to visualize the different quality classes (values 3-8) in two-dimensional space. PCA is a form of linear dimensionality reduction by using singular value decomposition to bring the data to lower dimensional space for data visualization purposes. The resulting figure below highlights that this dataset is not linearly separable and involves multi-class classification.

PCA also gave us information about explained variance and the corresponding ratios. Using `pca.explained_variance_ratio_`, we found that the first principal component is only responsible for 28.17% of the variance while the second principal component explains 17.51% of the variance. Together, these two components explain approximately 45% of the variance in the outcomes.

```
pca_data = wineData.copy()
X_pca = pca_data.loc[:, 'fixed acidity':'alcohol']
y_pca = pca_data['quality']
X_pca.tail()
X_pca = StandardScaler().fit_transform(X_pca)
pca = PCA(n_components = 2) #2D
X_pca = pca.fit_transform(X_pca) #fitting the PCA
X_pca.shape
(1599, 2)
```



explained variance ratio: [0.2372106 0.12296171]

Feature Selection

Now that we understood the target variable distribution, we wanted to investigate the input variables through a feature selection method. We decided to use a univariate feature selection method (SelectKBest) because we wanted to see the relationships between each input variable and the target variable individually to further evaluate and choose the k highest scoring features. We used the chi-squared statistical test to determine which features have the best relationships to the quality of wine. We found that total sulfur dioxide had the best relationships with wine quality by far, with a score of 2755.56, and free sulfur dioxide and alcohol followed with 161.94 and 46.42, respectively. Because chlorides, pH, and

density had very low scores, we decided to drop them from the dataset and use the highest scoring eight features.

```
X_f = wineData.loc[:, 'fixed acidity':'alcohol']
y_f = wineData['quality']

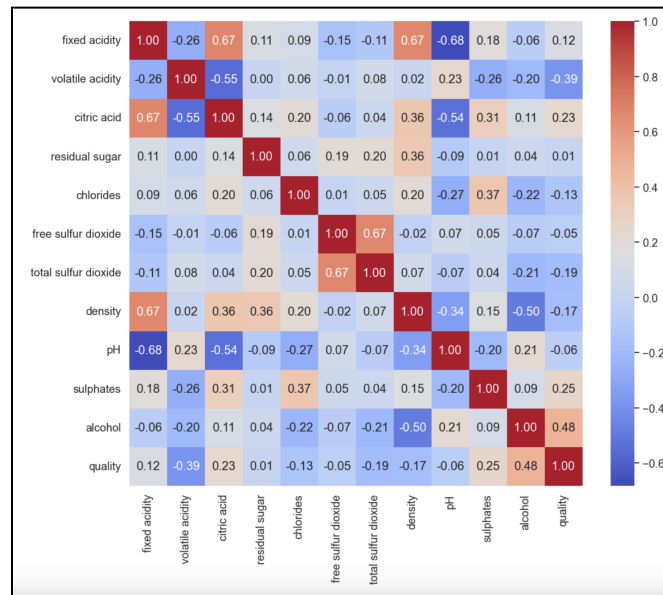
bestFeatures = SelectKBest(score_func=chi2, k=11)
bestFeaturesFit = bestFeatures.fit(X_f, y_f)
scores = pd.DataFrame(bestFeaturesFit.scores_)
predictor_variables = pd.DataFrame(X_f.columns) |

#concatenate scores with predictor names
predScores = pd.concat([predictor_variables, scores], axis=1)
predScores.columns = ['Predictor Variable', 'Score']

print("Features ranked from highest score to lowest")
print(predScores.nlargest(12, 'Score'))
```

	Predictor Variable	Score
6	total sulfur dioxide	2755.557984
5	free sulfur dioxide	161.936036
10	alcohol	46.429892
1	volatile acidity	15.580289
2	citric acid	13.025665
0	fixed acidity	11.260652
9	sulphates	4.558488
3	residual sugar	4.123295
4	chlorides	0.752426
8	pH	0.154655
7	density	0.000230

To further evaluate the relationships between each feature variable and the target variable, we used a correlation matrix and heatmap. Seeing that none of the input features had a significant correlation with wine quality, we decided that feature selection was not the appropriate method to approach this dataset as it would lower model accuracy and performance, so we decided to use all 11 features in the dataset for training the algorithms. (We later tested the models on the feature selection filtered dataset and found that each model had significantly worse accuracy scores, thus supporting our proposition that feature selection should not be used for this analysis.)



After splitting the data into testing and training datasets (70% training, 30% testing) and scaling the data, we started training and testing 3 models: logistic regression, random forest classifier, and support vector machine. For each model, we evaluated the performance using information from the corresponding classification report and confusion matrix.

Logistic Regression Model

We decided to use a logistic regression model for this dataset because we thought the underlying logistic function would be an appropriate predictor to apply on the various continuous variables we have as inputs. With this model, we decided to investigate how training on scaled data changes model performance as compared to unscaled data. We used the `LogisticRegression()` class from the `sklearn` package, fit it on the scaled and unscaled training datasets, predicted the `y` values given the respective test data, and evaluated performance of the model in each case through outputting an accuracy score, classification report, and confusion matrix. (Reference appendix for code)

Random Forest Classifier Model

We also created a Random Forest Classifier model. We chose this model in particular because Random Forest Classification because the performance is less influenced by outliers than other algorithms, it provides more information into the relationships between individual features and the target variable, and performs well on multinomial classification problems. We used the `RandomForestClassifier()` class from the `sklearn` package to implement this model and we focused on evaluating the impact of hyperparameter tuning on model performance. For the unoptimized model, we set `n_estimators` to 10 as that is the default for previous versions of `sklearn` examples. We then used `GridSearchCV` to tune the hyperparameters, which does an exhaustive search over the parameter values to find the optimal values. This is feasible because the dataset is not very large and thus not very computationally expensive for an exhaustive search. After finding the optimal values of `max_depth` and `n_estimators`, we created an instance of the `RandomForestClassifier()` class with these values, fit them to the training dataset, predicted values for the test data, and output accuracy and performance measures for evaluation. (Reference appendix for code)

Support Vector Machine Model

Lastly, we chose to create a Support Vector Machine (SVM) model because it is optimal for classification problems and it can handle non-linearly separable data well. In addition, while logistic regression involves a probabilistic approach to creating a separating hyperplane, SVM uses a statistical approach, which might be better fitted for this type of dataset. We used the `SVC()` class from the `sklearn` package to create the SVM model and compared the accuracy and performance of SVM models with different types of kernels: `rbf`, `sigmoid`, and `poly`. We also investigated the effect of the regularization parameter `C` value on model performance. Seeing that the model with the `rbf` kernel performed the best on this dataset in the previous part, we kept the kernel constant at `rbf` and created several instances of `SVC`, fitted the model on the training data, predicted values on the test data, and evaluated performance for various values of `C` = [1, 10, 100, 1000]. We then plotted the model accuracy score against the different values of `C` to visualize the effect of the parameter on model performance. Finally, we used `GridSearchCV` to tune the hyperparameters using an exhaustive search over parameter values. After finding the optimal values of kernel and `C`, we created instances of `SVC()` with these parameters, fit them to the training dataset, predicted values for the test data, and output accuracy and performance measures for evaluation. (Reference appendix for code)

RESULTS AND ANALYSIS

Logistic Regression Model

The logistic regression model that was fitted and tested on unscaled data had an accuracy score of 59.17% while the model for scaled data produced an accuracy score of 61.25%. This 2.08% increase in the model accuracy score demonstrates the importance of scaling the data before training or testing as it helps improve model performance.

Accuracy on test data without data scaling 0.5916666666666667 Accuracy on test data with data scaling 0.6125

But overall, these are not very high performance scores and would not be ideal for analyzing a dataset. This may reflect the nature of logistic regression models, which tend to perform better on binary classification problems as compared to multinomial classification problems.

We also looked at the differences in the weights (coefficients) and intercepts of the two logistic regression models. Generally speaking, if the intercept is positive, the likelihood of the corresponding class outcome will be greater than 50%. For the unscaled dataset model, the intercepts corresponding to classes 4 and 5 were positive, while for the scaled dataset model, the intercepts corresponding to classes 5, 6, and 7 were positive. This contributes to the proposition that scaling data highly influences the model outcomes.

Logistic Regression (not scaled data) Constant: [-0.00654774 0.07035266 0.48444819 -0.09470671 -0.38663586 -0.06691054]
--

Logistic Regression (scaled data) Constant: [-3.10996993 -0.16978106 2.70458974 2.86322264 0.59741842 -2.88547981]

Random Forest Classifier Model

We first created a Random Forest Classifier model (before hypertuning the parameters) and this model had an accuracy score of approximately 65.83%. Using an exhaustive grid search method, we found that the optimal max_depth level was 16 and the optimal n_estimators value was 200. Using these parameters, we increased the model accuracy score by just over 4.5%. This demonstrates the importance of hypertuning parameters to optimal values in order to improve model performance.

{'max_depth': 16, 'n_estimators': 200}
--

Accuracy on unoptimized RFC 0.6583333333333333 Accuracy on optimized RFC 0.7041666666666667
--

Overall, the Random Forest Classifier model had the highest accuracy score of any model we tested with an optimized model accuracy score of approximately 70.42%. This reflects the nature of the Random

Forest Classification method in being optimal for multinomial classification problems and handling datasets with continuous variables well.

Support Vector Machine Model

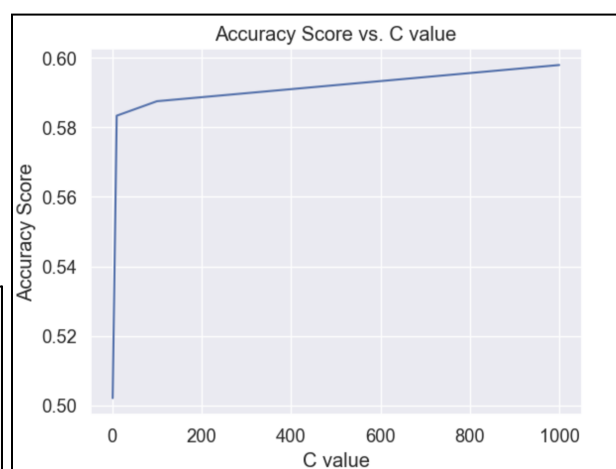
When comparing different kernel types for SVM, we found that for this dataset, the model with an rbf kernel had the highest accuracy score of 50%, while the model with a poly kernel followed with an accuracy score of 48% and then model with a sigmoid kernel with an accuracy score of 44%. The model with an rbf kernel performed the best because the kernel is optimal for non-linear datasets with multiple classes. The rbf kernel is able to map points to higher-dimensional spaces that are easier to separate and thus predict classifications.

rbf kernel classification report:					sigmoid kernel classification report:					poly kernel classification report:				
	precision	recall	f1-score	support		precision	recall	f1-score	support		precision	recall	f1-score	support
3	0.00	0.00	0.00	3	3	0.00	0.00	0.00	3	3	0.00	0.00	0.00	3
4	0.00	0.00	0.00	11	4	0.00	0.00	0.00	11	4	0.00	0.00	0.00	11
5	0.68	0.43	0.53	213	5	0.48	0.54	0.51	213	5	0.76	0.31	0.44	213
6	0.43	0.82	0.56	181	6	0.42	0.53	0.47	181	6	0.42	0.92	0.58	181
7	1.00	0.02	0.03	65	7	0.00	0.00	0.00	65	7	1.00	0.02	0.03	65
8	0.00	0.00	0.00	7	8	0.00	0.00	0.00	7	8	0.00	0.00	0.00	7
accuracy			0.50	480	accuracy			0.44	480	accuracy			0.48	480
macro avg	0.35	0.21	0.19	480	macro avg	0.15	0.18	0.16	480	macro avg	0.36	0.21	0.17	480
weighted avg	0.60	0.50	0.45	480	weighted avg	0.37	0.44	0.40	480	weighted avg	0.63	0.48	0.42	480

We also evaluated the impact of the regularization parameter C in SVM models, which is related to the margin of the hyperplane. In SVM models, there is a tradeoff between finding the largest minimum margin and achieving the lowest testing error. A lower value of C indicates a higher margin and as the value of C increases, the margin decreases. Typically, the smaller the margin, the higher the accuracy. As we can see in the resulting graph below, as the value of C increased, the accuracy score of the model increased, reaching approximately 59.79% when C=1000. But we must also consider that C can't be too big because it will cause a large testing error error (overfitted model not generalizable) and C can't be too small because it will cause a large training error (underfitted model). This highlights the importance of optimizing the value of C through hyperparameter tuning.

```
#Testing the effect of C values on SVM accuracy
C_values = [1, 10, 100, 1000]
Cval_scores = []

for c in C_values:
    wineSVM = SVC(C=c, kernel='rbf')
    wineSVM.fit(X_train, y_train)
    y_pred = wineSVM.predict(X_test)
    Cval_scores.append(accuracy_score(y_test, y_pred))
```



After establishing the importance of hyperparameter tuning, we used an exhaustive grid search method to find that the optimal kernel was rbf and the optimal C value was 1000. This produced a model with an

accuracy score of 59.79%. Overall, this accuracy is not very high, which likely reflects the variation in the dataset or misleading feature information that makes it difficult to predict the various classes.

Accuracy on optimized SVM: 0.5979166666666667				
	precision	recall	f1-score	support
3	0.00	0.00	0.00	0
4	0.00	0.00	0.00	0
5	0.83	0.67	0.74	264
6	0.57	0.51	0.54	204
7	0.11	0.58	0.18	12
8	0.00	0.00	0.00	0
accuracy			0.60	480
macro avg	0.25	0.29	0.24	480
weighted avg	0.70	0.60	0.64	480

CONCLUSIONS

Ultimately, the Random Forest Classification model with hyperparameter tuning had the best performance on the test dataset with an accuracy score of approximately 70.4%. The untuned random forest classifier came in second with 65.8% accuracy, followed by optimized and unoptimized logistic regression and then lastly SVM. Although these accuracy levels may seem quite low, the complex and (at times) seemingly contradictory nature of the dataset implies that these models are performing quite well, especially the random forest classifier. We can conclude that for complex, high-dimension, non-linearly separable datasets that the random forest classifier (and generally speaking, more complex, robust, and computationally expensive classification algorithms) will perform the strongest.

Division of Work

Sam: Logistic regression model and random forest classifier model code, introduction and conclusion sections of paper

Josh: Feature selection code, classification report and confusion matrix figures in code, methods section of paper

Ritika: Preprocessing and PCA code, support vector machine model code, results and analysis section of paper

References

- PS5_Heart Disease Example
- <https://gursev-pirge.medium.com/performance-comparison-of-multi-class-classification-algorithms-606e8ba4e0ee>
- <https://machinelearningmastery.com/feature-selection-with-real-and-categorical-data/>
- https://scikit-learn.org/stable/modules/feature_selection.html
- https://scikit-learn.org/stable/modules/generated/sklearn.model_selection.GridSearchCV.html
- <https://quantifyinghealth.com/interpret-logistic-regression-coefficients/#:~:text=If%20the%20interpret%20has%20a,outcome%20will%20be%20exactly%200.5.>
- <https://datascience.stackexchange.com/questions/54751/when-to-use-random-forest>

APPENDIX

Accuracy on optimized SVM: 0.597916666666667

	precision	recall	f1-score	support
3	0.00	0.00	0.00	0
4	0.00	0.00	0.00	0
5	0.83	0.67	0.74	264
6	0.37	0.51	0.54	204
7	0.11	0.58	0.18	12
8	0.00	0.00	0.00	0
accuracy			0.60	480
macro avg	0.25	0.29	0.24	480
weighted avg	0.70	0.60	0.64	480

