

Final Project: Predict Apple Quality

Ziwei Su, Wenhao Gu

1 Introduction

Evaluation and prediction of fruit quality are critical for ensuring a steady supply of high-quality fruits to meet market demands. The Apple Quality dataset, available on Kaggle Apple Quality, presents an interesting and valuable resource for predicting the quality of apples based on diverse characteristics. The primary objective of this project is to conduct a comparative analysis of various predictive models to achieve high predictive accuracy in determining apple quality while concurrently evaluating the significance of various apple attributes in predicting apple quality.

2 Data Description

The dataset comprises 4,000 rows and 9 columns, with each row corresponding to the distinct attributes of an individual apple. These attributes include:

- **A_id**: Unique identifier for each fruit.
- **Size**: Size of the fruit.
- **Weight**: Weight of the fruit.
- **Sweetness**: Degree of sweetness of the fruit.
- **Crunchiness**: Texture indicating the crunchiness of the fruit.
- **Juiciness**: Level of juiciness of the fruit.
- **Ripeness**: Stage of ripeness of the fruit.
- **Acidity**: Acidity level of the fruit.
- **Quality**: Overall quality of the fruit. It takes value in {good, bad}.

The following table gives the first five rows of the Apple Quality dataset.

A_id	Size	Weight	Sweetness	Crunchiness	Juiciness	Ripeness	Acidity	Quality
0	-3.970049	-2.512336	5.346330	-1.012009	1.844900	0.329840	-0.491590	good
1	-1.195217	-2.839257	3.664059	1.588232	0.853286	0.867530	-0.722809	good
2	-0.292024	-1.351282	-1.738429	-0.342616	2.838636	-0.038033	2.621636	bad
3	-0.657196	-2.271627	1.324874	-0.097875	3.637970	-3.413761	0.790723	good
4	1.364217	-1.296612	-0.384658	-0.553006	3.030874	-1.303849	0.501984	good

Table 1: Apple Quality Dataset Sample

Figure 1 summarizes the dataset. We choose Size, Weight, Sweetness, Crunchiness, Juiciness, Ripeness, and Acidity as the predictors, and Quality is the binary response. We exclude A_id because the identifier is irrelevant to our analysis. We transform Acidity to float, normalize all the predictors, and transform the binary response to {0,1}. Figure 2 displays the correlation coefficients among the 7 predictors, showing little collinearity among them. This can also be observed from the scatter matrix in Figure 3.

```

<class 'pandas.core.frame.DataFrame'>
RangeIndex: 4000 entries, 0 to 3999
Data columns (total 8 columns):
#   Column      Non-Null Count  Dtype
---  -
0   Size        4000 non-null   float64
1   Weight       4000 non-null   float64
2   Sweetness    4000 non-null   float64
3   Crunchiness  4000 non-null   float64
4   Juiciness    4000 non-null   float64
5   Ripeness     4000 non-null   float64
6   Acidity      4000 non-null   object
7   Quality      4000 non-null   object
dtypes: float64(6), object(2)
memory usage: 250.1+ KB

```

Figure 1: Summary of the Apple Quality Dataset

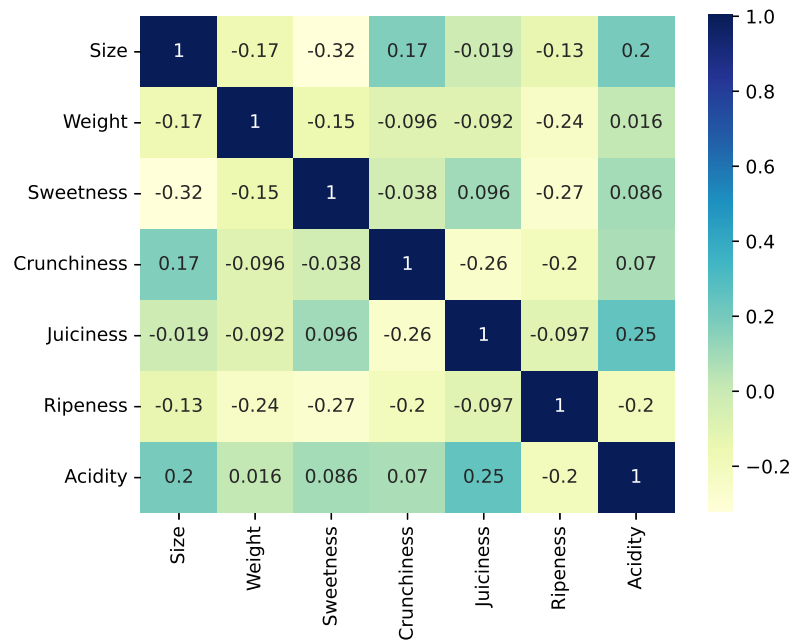


Figure 2: Correlations of the predictors

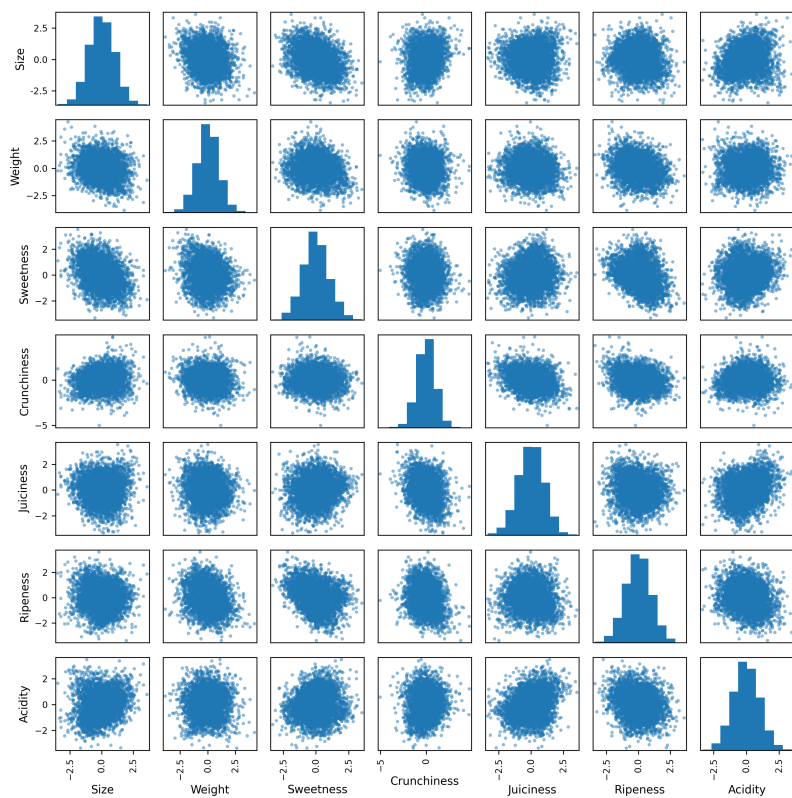


Figure 3: Scatter matrix of the predictors

3 Model Fitting and Analysis

In this section, we will be deploying the following eight models to the Apple Quality dataset for classification: Logistic Regression (LogR), Naive Bayes (NB), Linear Discriminant Analysis (LDA), Support Vector Machine (SVM), K-Nearest Neighbors (KNN), Neural Network (NN), Classification Tree (CT), Light Gradient Boosting Machine (LGBM), and Random Forest (RF). We measure the predictive power of each model with the mean cross-validation (CV) error with multiple replications of 10-fold CV, summarized in Table 2. In general, we use the permutation-based variable importance

Table 2: Mean CV accuracy of models.

Model	CV accuracy
Logistic Regression (LogR)	0.748
Naive Bayes (NB)	0.749
Linear Discriminant Analysis (LDA)	0.748
Support Vector Machine (SVM)	0.915
K-Nearest Neighbors (KNN)	0.900
Neural Network (NN)	0.946
Classification Tree (CT)	0.819
Light Gradient Boosting Machine (LGBM)	0.900
Random Forest (RF)	0.891

measure (VIM) to quantify the importance of the predictors. We use other approaches for the predictor importance such as the Accumulated Local Effects (ALE) plot for the logistic regression model and the impurity-based VIM for tree models.

3.1 Logistic Regression (LogR)

In this section, we use the logistic regression model with *sklearn.linear_model.LogisticRegression*. The hyperparameters to tune are C , *penalty*, and *l1_ratio*, where C is the inverse of regularization strength, *penalty* is the type of regularization penalty, and *l1_ratio* is a mixing parameter taking values between 0 and 1. Possible regularization methods include $l1$ (LASSO), $l2$ (Ridge), and elastic net, a mixture of $l1$ and $l2$ regularization. $l1_ratio = 0$ is equivalent to using $l2$ penalty, while setting $l1_ratio = 1$ is equivalent to using $l1$ penalty. For $0 < l1_ratio < 1$, the penalty is a combination of $l1$ and $l2$. We perform grid search with 10-fold cross-validation repeated 10 times on the following hyperparameter grid:

$$C \in \{0.01, 0.05, 0.1, 0.5, 1, 5, 10\}, \text{penalty} \in \{l1, l2, \text{elasticnet}, \text{none}\}, l1_ratio \in \{0.1, 0.3, 0.5, 0.7, 0.9\}.$$

The best hyperparameter configuration is $C = 0.05$, $l1_ratio = 0.7$, *penalty* = *elasticnet*, and the corresponding model has a CV mean accuracy of around 0.748.

We consider the importance of predictors via the coefficients of each predictor in the model, the permutation VIM, and the ALE plot against the log odds. The feature coefficients are as follows:

- Size: The coefficient for size is 1.1277.
- Sweetness: The coefficient for sweetness is 0.9920.
- Juiciness: The coefficient for juiciness is 0.7664.
- Acidity: The coefficient for acidity is -0.5371 .
- Weight: The coefficient for weight is 0.3507.
- Ripeness: The coefficient for ripeness is -0.2454 .
- Crunchiness: The coefficient for crunchiness is 0.0096.

The permutation VIM is given in Figure 4, and the ALE plot is given in Figure 5. We can tell from the coefficients, the permutation VIM, and the ALE that Size appears to be the most important predictor, followed by Sweetness, and Crunchiness appears to be the least important.

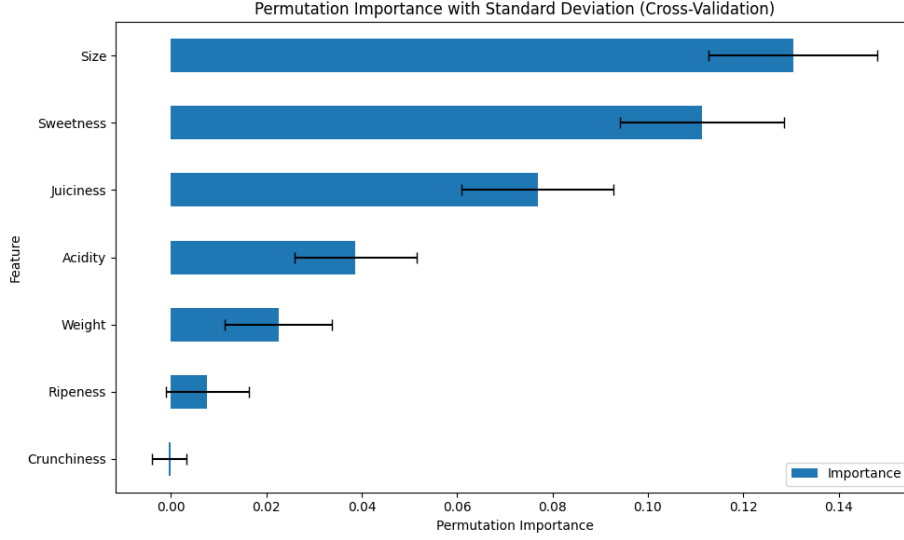


Figure 4: Permutation VIM of the best LogR model

3.2 Naive Bayes (NB)

In this section, we use the Naive Bayes model with *sklearn.naive_bayes.GaussianNB*. The only hyperparameter to tune is the amount of smoothing (*var_smoothing*), which is added to the variance of each feature during fitting to alleviate overfitting. We perform grid search with 10-fold cross-validation on the following hyperparameter grid:

$$var_smoothing \in \{1e-2, 5e-2, 1e-1, 5e-1, 1, 5, 10\}.$$

The best hyperparameter value is *var_smoothing* = 0.01, and the corresponding model achieves a CV mean accuracy of around 0.749. We consider the importance of predictors via the permutation VIM, given in Figure 6. According to the permutation VIM, the order of predictor importance from most to least is: Juiciness, Size, Sweetness, Ripeness, Weight, Crunchiness, and Acidity. The first six predictors vary little in their VIM scores, while Acidity appears to be very unimportant in comparison.

3.3 Linear Discriminant Analysis (LDA)

In this section, we use the Linear Discriminant Analysis model with *sklearn.discriminant_analysis*. There are no hyperparameters to tune, as when we set the parameter *solver* = 'lsqr' (least squares solution), the shrinkage parameter can be automatically selected using the Ledoit-Wolf lemma. The model achieves a CV mean accuracy of around 0.7479. We consider the importance of predictors via the permutation VIM, given in Figure 7. According to the permutation VIM, the most important predictors in descending order are: Size, Sweetness, and Juiciness, while the least important predictor is Crunchiness.

3.4 Support Vector Machine (SVM)

In this section we use the support vector machine (SVM) algorithm with *sklearn.svm.SVC*. The parameters are *C* and *kernel*, where *C* is the regularization parameter and *kernel* is the kernel for the SVM classifier. But we choose *kernel* = 'rbf', meaning that we are using Gaussian kernel, because 'rbf' is efficient to compute. we use K-fold cross validation to do a three stage grid search over *C*. First we search

$$C \in \{0.1, 1, 10, 50, 100, 1000\}$$

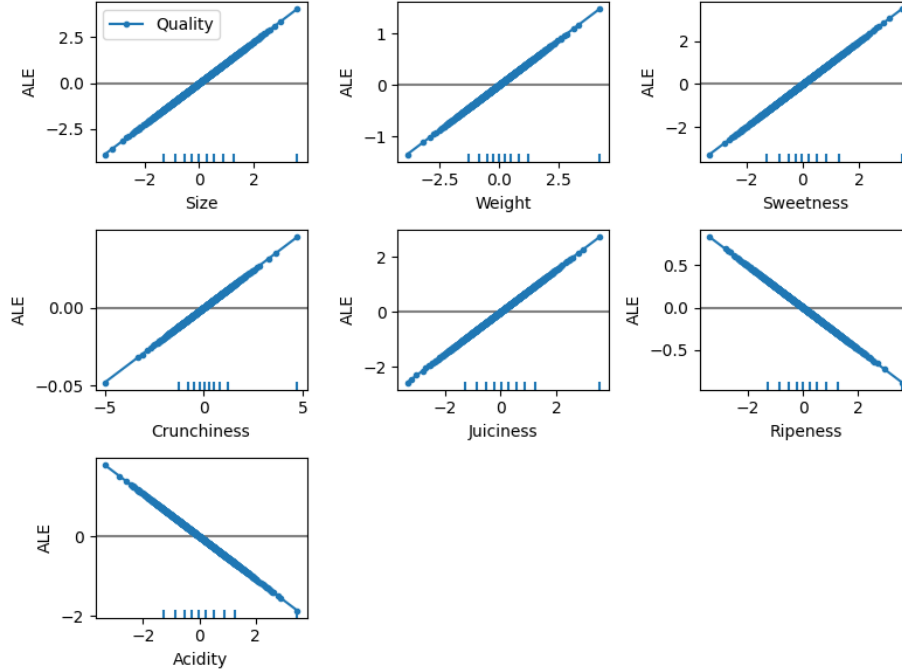


Figure 5: ALE of the best LogR model

and found that the best C in this grid is 50. We then search

$$C \in \{10, 20, 30, \dots, 90, 100\}$$

and found that the best parameter is $C = 40$. Finally, we search

$$C \in \{30, 31, 32, \dots, 49, 50\}$$

and the best parameter is $C = 41$ with the best CV accuracy being 0.9151. For the model with the selected best parameter, we calculate the CV permutation VIM importance of the predictors and visualize them in 8.

From the permutation VIM importance plot, the most important three predictors are Ripeness, Sweetness and Size, and the least important predictor is Crunchiness.

3.5 K-nearest Neighbors (KNN)

In this section we use the K-Nearest Neighbors (KNN) algorithm to fit the prediction model. The parameter for the distance-based KNN classifier is $n_neighbors$, and we use K-fold cross validation to do a grid search over $n_neighbors \in \{1, 2, \dots, 30\}$. The resulting cross validation accuracy versus $n_neighbors$ is shown in Figure 9.

The best parameter is $n_neighbors = 11$, and the CV accuracy is 0.90045. For the model with the selected best parameter, we calculate the CV permutation VIM importance of the predictors, which are shown in Figure 10.

From the permutation VIM importance plot, the most important three predictors are Ripeness, Size and Sweetness, and the least important two predictors are Crunchiness and Weight.

3.6 Neural Network (NN)

In this section, we use a neural network model with `sklearn.neural_network.MLPClassifier`. The activation function is the logistic sigmoid function. The hyperparameters to tune are the number of

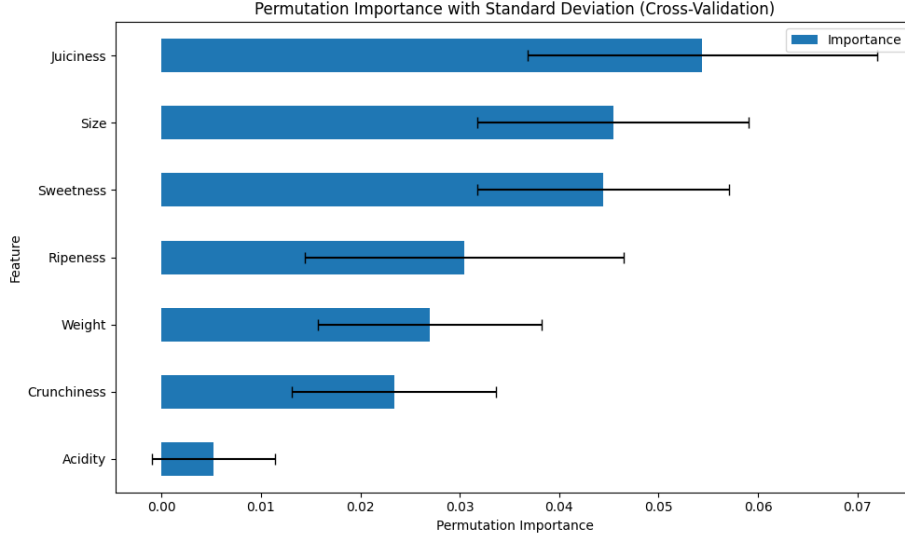


Figure 6: Permutation VIM of the best Naive Bayes model

neurons in the hidden layer (*hidden_layer_sizes*), the l_2 regularization parameter α , and the initial learning rate (*learning_rate_init*). We perform a grid search with 10-fold cross-validation repeated 1 time on the following hyperparameter grid:

$$hidden_layer_sizes \in \{(50,), (100,)\}, \alpha \in \{0.0001, 0.001, 0.01\}, learning_rate_init \in \{0.01, 0.1\}.$$

The neural network is very costly to tune with cross-validation. Due to time constraints, only 1 replication of 10-fold cross-validation was performed. The best hyperparameter configuration is $\alpha = 0.01$, *hidden_layer_sizes* = (50,), *learning_rate_init* = 0.01, and the corresponding model has a CV mean accuracy of around 0.946. We consider the importance of predictors via the permutation VIM, given in Figure 11. The permutation VIM shows that the most important predictor is Ripeness, followed closely by Sweetness and Size. The next tier of predictors, in order of decreasing importance, are Juiciness, Acidity, Weight and Crunchiness, which differ little in importance from each other. Overall, the difference in predictor importance for all predictors is not very large.

3.7 Classification Tree (CT)

In this section, we use the classification tree model *sklearn.tree.DecisionTreeClassifier*. The hyperparameters are the maximum depth of the tree (*max_depth*), the minimum number of samples required to be at a leaf node (*min_samples_split*), the minimum number of samples required to be at a leaf node (*min_samples_leaf*), and the complexity parameter used for Minimal Cost-Complexity Pruning (*ccp_alpha*). We perform a four stage grid search with 10-fold cross-validation repeated. We first search over

$$\begin{aligned} max_depth &\in \{5, 10, 15, 20\} \\ min_samples_split &\in \{2, 3, 5, 10\} \\ min_samples_leaf &\in \{2, 3, 5, 10\} \\ ccp_alpha &\in \{0.0, 0.001, 0.01, 0.1\} \end{aligned}$$

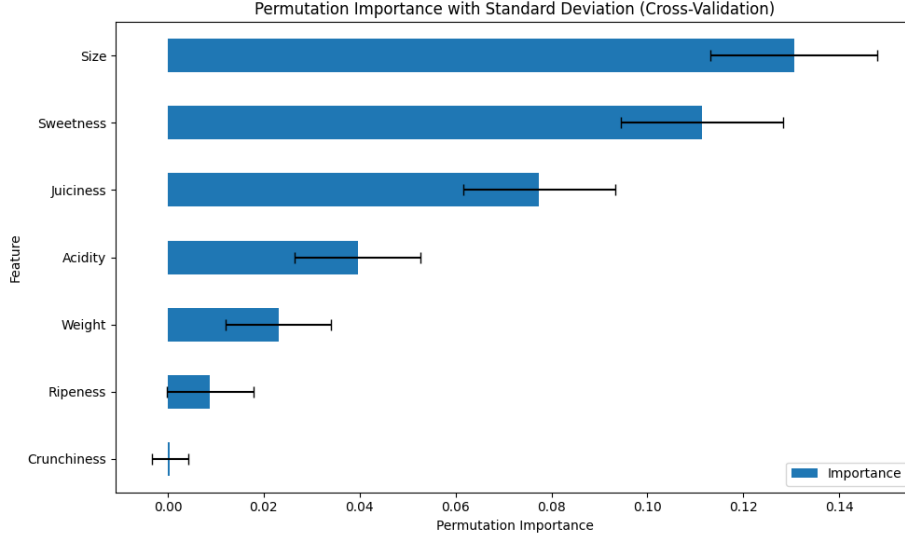


Figure 7: Permutation VIM of the LDA model

The best hyperparameter for this set of parameters is $max_depth = 15$, $min_samples_split = 5$, $min_samples_leaf = 2$, $ccp_alpha = 0.001$. Next we search over

$$\begin{aligned}
 max_depth &\in \{13, 15, 17, 19\} \\
 min_samples_split &\in \{4, 5, 7, 9\} \\
 min_samples_leaf &\in \{2, 3, 4\} \\
 ccp_alpha &\in \{0.0, 0.0001, 0.0005, 0.001\}
 \end{aligned}$$

and the best parameter in the 2nd search is $max_depth = 15$, $min_samples_split = 7$, $min_samples_leaf = 4$, $ccp_alpha = 0.0005$. In the 3rd search we examine

$$\begin{aligned}
 max_depth &\in \{14, 15, 16\} \\
 min_samples_split &\in \{6, 7, 8\} \\
 min_samples_leaf &\in \{3, 4\} \\
 ccp_alpha &\in \{0.0002, 0.0003, 0.0005, 0.0007, 0.0009\}
 \end{aligned}$$

and the best parameter in the 3rd search is $max_depth = 16$, $min_samples_split = 8$, $min_samples_leaf = 4$, $ccp_alpha = 0.0007$. In the last search, we test

$$\begin{aligned}
 max_depth &\in \{15, 16, 17\} \\
 min_samples_split &\in \{7, 8, 9\} \\
 min_samples_leaf &\in \{3, 4, 5\} \\
 ccp_alpha &\in \{0.0006, 0.0007, 0.0008\}
 \end{aligned}$$

and the best model is $max_depth = 15$, $min_samples_split = 7$, $min_samples_leaf = 5$, $ccp_alpha = 0.0007$. The corresponding best CV accuracy is 0.8195. We show the feature importance generate by the `sklearn.tree.DecisionTreeClassifier.feature_importances_` attribute in Figure 12

this shows that the most important predictor is Ripeness, and the least important predictor is Crunchiness. We next examine the CV permutation in Figure 13

this shows that the most important predictor is Ripeness, and the least two important predictors are Crunchiness and Weight, which is the same as in the results shown in the `feature_importance_` function. However, the importance ranking of the other predictors are not the same.

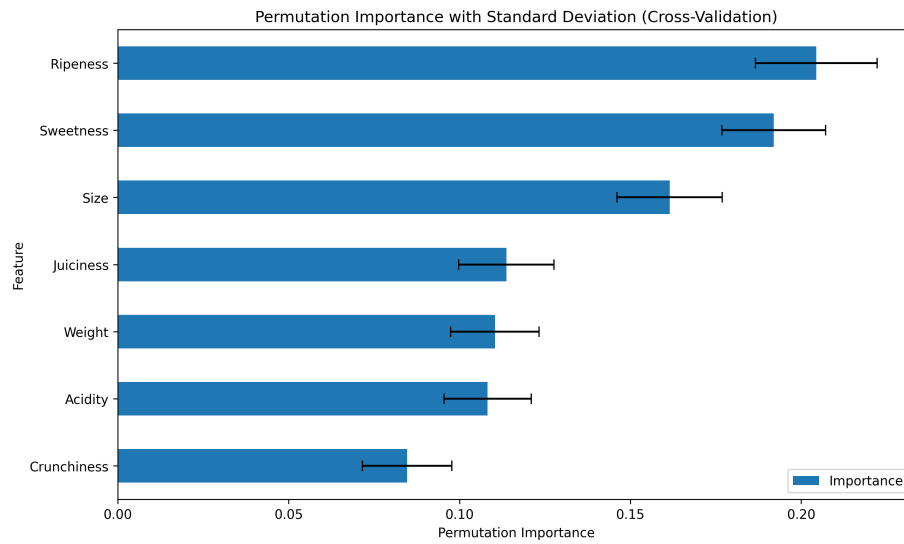


Figure 8: Permutation VIM of the best KNN model

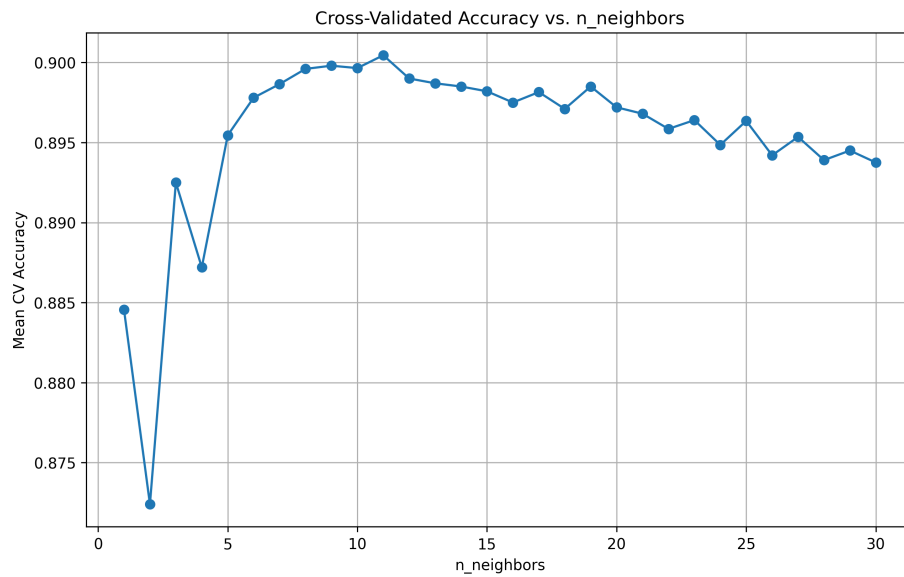


Figure 9: CV accuracy v.s. $n_neighbors$

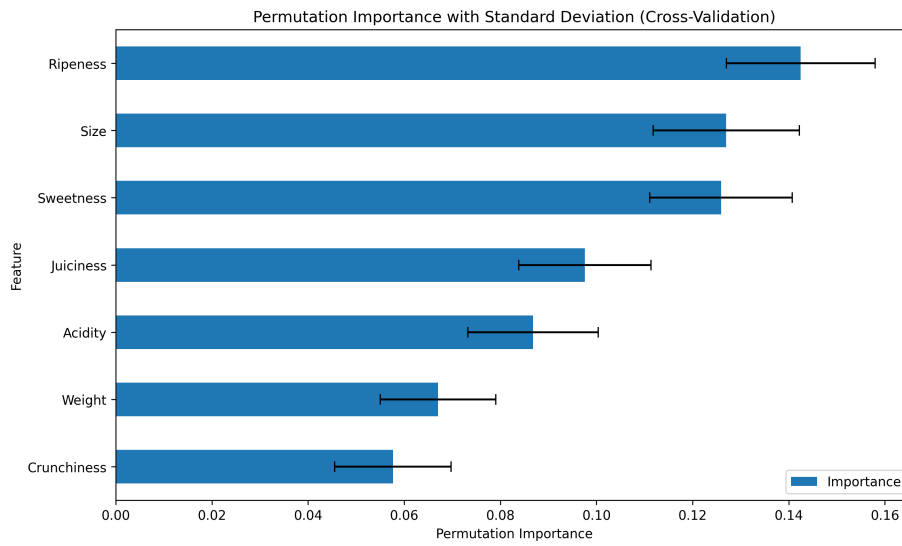


Figure 10: Permutation VIM of the best KNN model

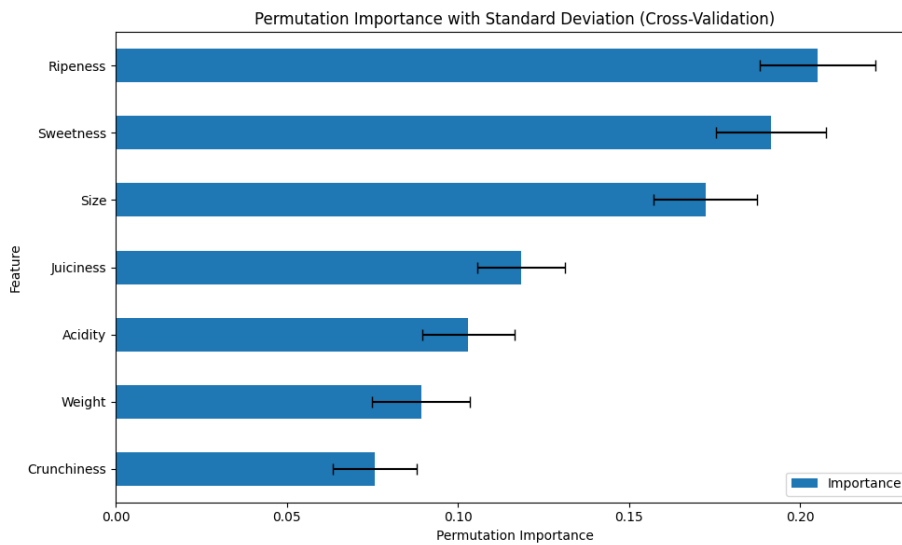


Figure 11: Permutation VIM of the best neural network model

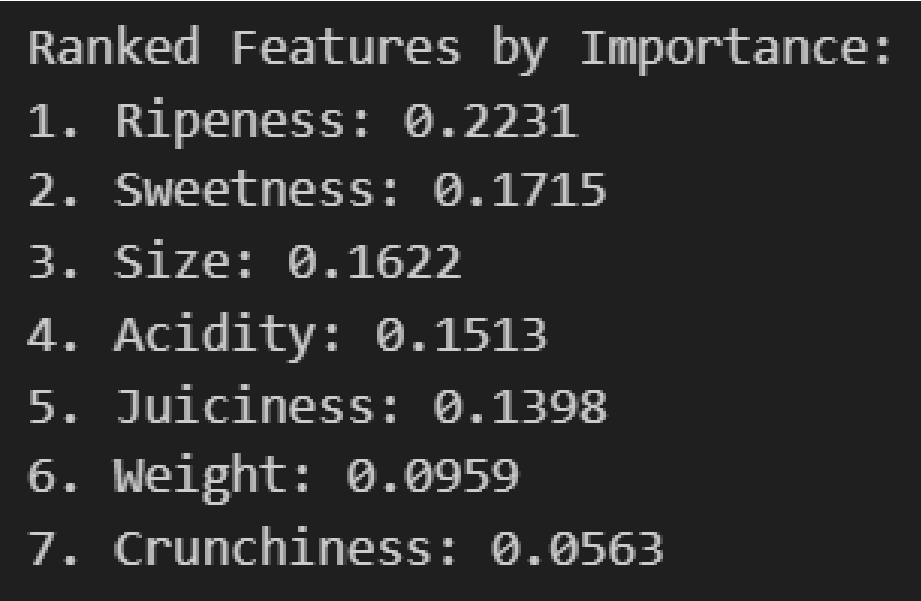


Figure 12: Feature importance of the best classification tree

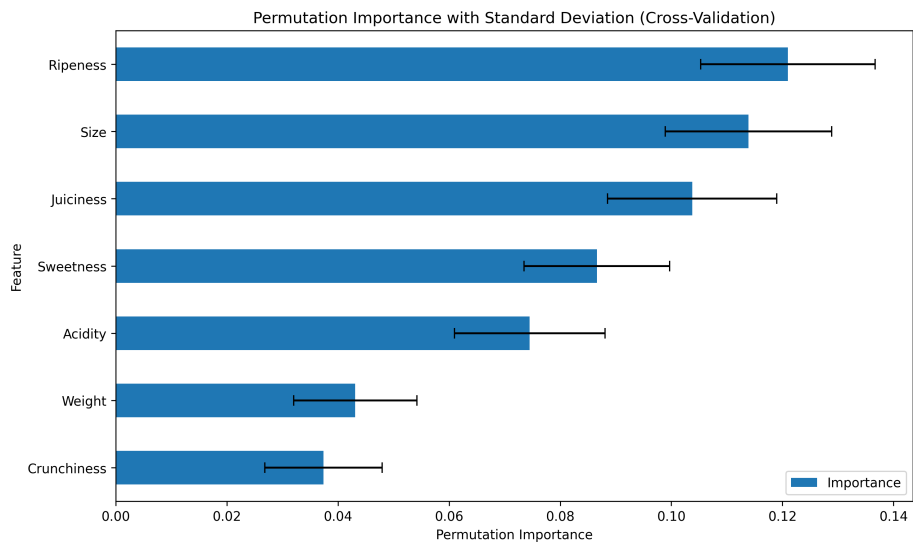


Figure 13: Permutation VIM of the best classification tree

3.8 Light Gradient Boosting Machine (LGBM)

In this section, we use the Light Gradient Boosting Machine *lightgbm.LGBMClassifier*. The hyperparameters are *learning_rate*, *n_estimators*, *max_depth* and *num_leaves*. We do a four stage grid search with 10-fold cross-validation repeated. We first search over

$$\begin{aligned} \text{learning_rate} &\in \{0.01, 0.1, 0.5\} \\ \text{n_estimators} &\in \{200, 300\} \\ \text{max_depth} &\in \{10, 20, 30\} \\ \text{num_leaves} &\in \{15, 31, 63\} \end{aligned}$$

The best hyperparameter for this set of parameters is *learning_rate* = 0.1, *n_estimators* = 300, *max_depth* = 30, *num_leaves* = 63. Next we search over

$$\begin{aligned} \text{learning_rate} &\in \{0.05, 0.1, 0.2\} \\ \text{n_estimators} &\in \{300, 400\} \\ \text{max_depth} &\in \{30, 40\} \\ \text{num_leaves} &\in \{63, 96\} \end{aligned}$$

and the best parameter in the 2nd search is *learning_rate* = 0.2, *n_estimators* = 300, *max_depth* = 30, *num_leaves* = 96. In the 3rd search we examine

$$\begin{aligned} \text{learning_rate} &\in \{0.18, 0.2, 0.22\} \\ \text{n_estimators} &\in \{300, 400\} \\ \text{max_depth} &\in \{30\} \\ \text{num_leaves} &\in \{96\} \end{aligned}$$

and the best parameter in the 3rd search is *learning_rate* = 0.22, *n_estimators* = 400, *max_depth* = 30, *num_leaves* = 96. In the last search, we test

$$\begin{aligned} \text{learning_rate} &\in \{0.22, 0.24, 0.26, 0.28, 0.3\} \\ \text{n_estimators} &\in \{300, 400, 500\} \\ \text{max_depth} &\in \{30\} \\ \text{num_leaves} &\in \{96\} \end{aligned}$$

and the best model is *learning_rate* = 0.22, *n_estimators* = 500, *max_depth* = 30, *num_leaves* = 96. The corresponding best CV accuracy is 0.9003. We show the feature importance generate by the *lightgbm.plot_importance* function in Figure 14

this shows that the most important predictor is Ripeness, and the least important predictor is Crunchiness. We next examine the CV permutation in Figure 15

this shows that the most three important predictors are Ripeness, Size and Sweetness, and the least two important predictors are Crunchiness and Weight. However, the importance ranking of the other predictors are not the same.

3.9 Random Forest (RF)

In this section, we use the random forest model with *sklearn.ensemble.RandomForestClassifier*. The hyperparameters to tune are the number of trees (*n_estimators*), the maximum depth of the tree (*max_depth*), the minimum number of samples required to split an internal node (*min_samples_split*), and the minimum number of samples required to be at a leaf node (*min_samples_leaf*). We perform a grid search with 10-fold cross-validation repeated 1 time on the following hyperparameter grid:

$$\text{n_estimators} \in \{250, 300, 350\}, \text{max_depth} \in \{20, 25, 30\},$$

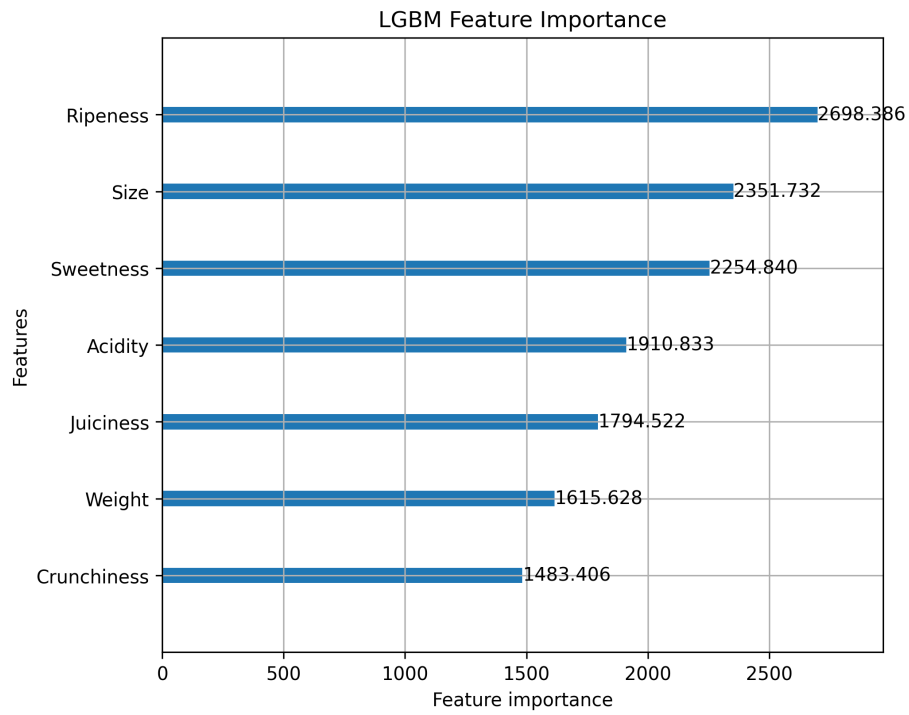


Figure 14: Feature importance of the best classification tree

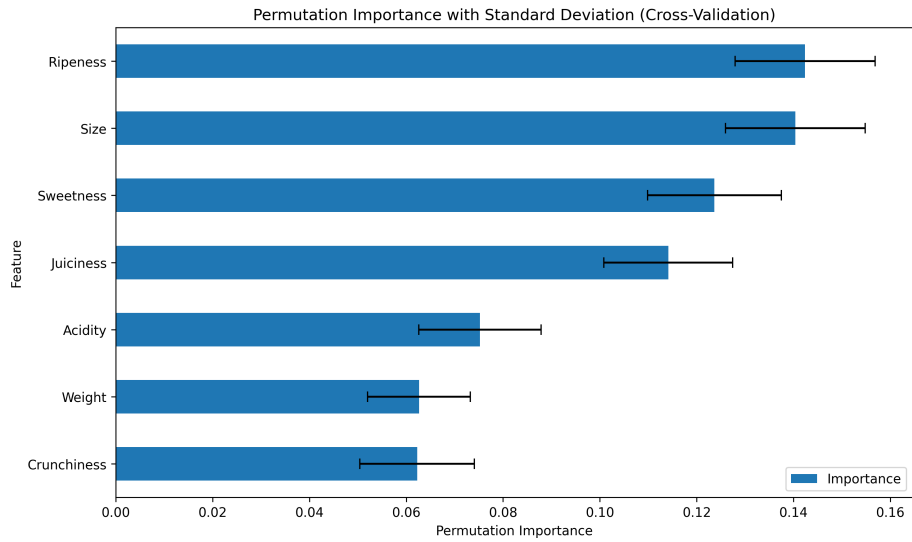


Figure 15: Permutation VIM of the best classification tree

$min_samples_split \in \{2, 5, 10\}, min_samples_leaf \in \{1, 2, 5\}.$

The best hyperparameter configuration is $max_depth = 25$, $min_samples_leaf = 1$, $min_samples_split = 2$, $n_estimators = 300$, and the corresponding model has a CV mean accuracy of around 0.891. The predictor importance is given by both the impurity-based measure (in Figure 16) and the permutation VIM (in Figure 17). Both measures agree that Ripeness, Size, Juiciness are the three most

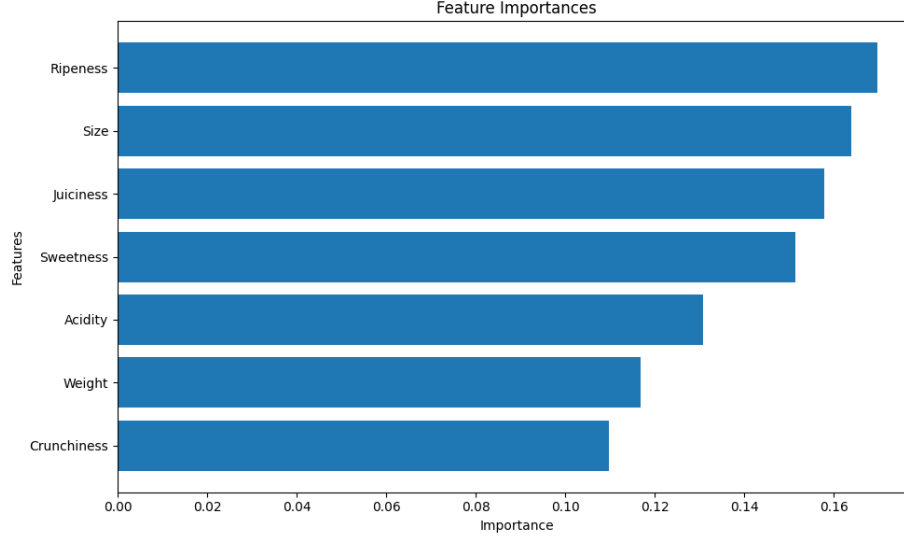


Figure 16: Impurity-based predictor importance of the best random forest model

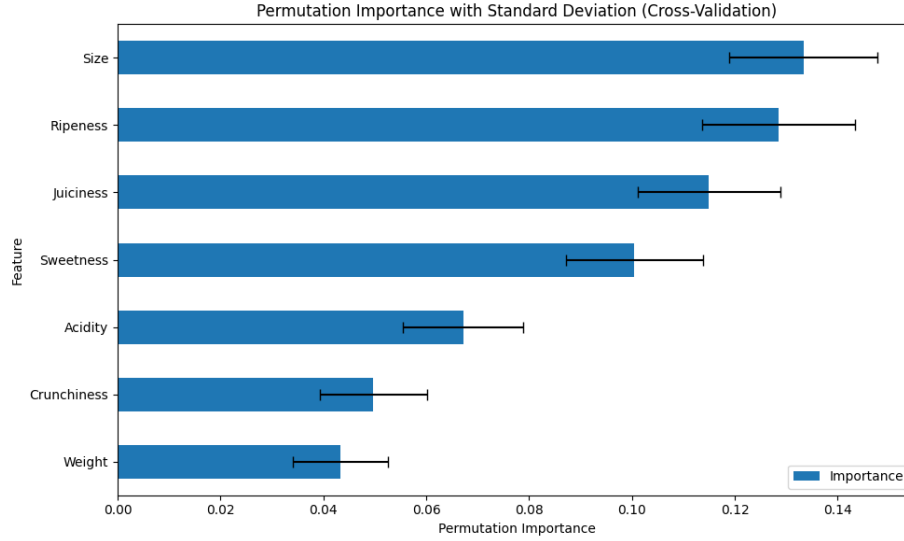


Figure 17: Permutation VIM of the best random forest model

important predictors, and Acidity, Crunchiness and Weight are the three least important predictors, but the difference between importance measures for different predictors seem small.

4 Conclusion

From Table 2, we conclude that the best model in terms of predictive power is Neural Network (NN) with a CV accuracy of around 0.946, followed by Support Vector Machine (SVM), Light Gradient Boosting Tree (LGBM), K-Nearest Neighbors (KNN) and Random Forest (RF), all with accuracy around 0.9. Simpler linear models such as Logistic Regression (LR), Naive Bayes (NB), and Linear Discriminant Analysis (LDA) only achieve CV accuracy of around 0.75.

The superior performance of NN and the relatively bad performance of linear models are expected, but the relatively good performance of KNN is unexpected for us. We speculate that the relatively good performance of KNN is due to the small number of predictors (7) and the small size of training data (4000). KNN is suitable for low-dimensional feature space due to less sparsity and a small training set due to its locality. Therefore, we recommend using NN for the best predictive power and KNN for a good trade-off between predictive performance and model interpretability.

In terms of the feature importance, different models yield different rankings. Ripeness and Size appear as the top two predictors, rank consistently among the top 3 across different models. Conversely, Crunchiness and Weights are at the lower end of the ranking, frequently appearing in the bottom 2. However, the relative rankings of the remaining predictors remain unclear to us at this point.

A Code

```
1 import pandas as pd
2 import numpy as np
3 from sklearn.preprocessing import StandardScaler
4
5 # data overview
6 df = pd.read_csv("apple_quality.csv")
7 df = df.loc[0:df.shape[0]-2,:]
8 df.head()
9
10 # remove id
11 df.drop('A_id', axis=1, inplace=True)
12 df.info()
13
14 # data preprocessing
15 df['Acidity'] = df['Acidity'].astype(float) # convert to float
16
17 # data normalization
18 # get the number of columns in the dataframe
19 num_columns = df.shape[1]
20
21 # create scaler
22 scaler = StandardScaler()
23
24 # select all columns except the last one
25 columns_to_normalize = df.columns[:num_columns-1]
26 print(columns_to_normalize)
27 # normalize the selected columns
28 df[columns_to_normalize] = scaler.fit_transform(df[columns_to_normalize])
29
30 # convert response to categorical
31 last_column = df.columns[-1]
32 df[last_column] = df[last_column].astype('category')
33 df.info()
34
35 from sklearn.preprocessing import LabelEncoder
36
37 # divide the data and encode the response
38 X = df[df.columns[:num_columns-1]]
39 y = df[df.columns[-1]]
40 label_encoder = LabelEncoder()
41 label_encoder.fit(['bad', 'good'])
```

```

42 y_encoded = label_encoder.transform(y)
43
44 #####
45 # Scatter plots
46 #####
47
48 import matplotlib.pyplot as plt
49 import seaborn as sns
50
51 axes = pd.plotting.scatter_matrix(X, alpha=0.5, figsize=(10, 10))
52 plt.tight_layout()
53 plt.savefig('scatter_matrix.png', dpi=300)
54
55 #####
56 # show correlation heatmaps
57 #####
58
59 dataplot = sns.heatmap(X.corr(), cmap="YlGnBu", annot=True)
60 fig = dataplot.get_figure()
61 fig.savefig("correlations.eps", format="eps", bbox_inches="tight")
62
63 # displaying heatmap
64 plt.show()
65
66
67 #####
68 # Model 1: Logistic Regression
69 #####
70
71 from sklearn.linear_model import LogisticRegression
72 from sklearn.model_selection import GridSearchCV, RepeatedKFold
73
74 # create model
75 model = LogisticRegression(max_iter=1000, random_state=0, solver='saga')
76
77 # define the hyperparameter grid
78 param_grid = {'C': [0.01, 0.05, 0.1, 0.5, 1, 5, 10],
79               'penalty': ['l1', 'l2', 'elasticnet', 'none'],
80               'l1_ratio': [0.1, 0.3, 0.5, 0.7, 0.9]}
81
82 # create cv object
83 cv = RepeatedKFold(n_splits=10, n_repeats=10, random_state=0)
84
85 # create grid-search cv object
86 grid_search = GridSearchCV(estimator=model, param_grid=param_grid, cv=cv, n_jobs=-1,
87                             scoring='accuracy')
88
89 # perform grid-search
90 grid_search.fit(X, y_encoded)
91
92 # get the best hyperparameters
93 best_params = grid_search.best_params_
94 best_score = grid_search.best_score_
95 best_model = grid_search.best_estimator_
96
97 # best hyperparameters
98 best_params
99
100 # best CV accuracy
101 best_score
102
103 from sklearn.metrics import accuracy_score
104
105 y_pred = best_model.predict(X)
106 # best model training error
107 accuracy_score(y_encoded, y_pred)
108
109 # get the coefficients of the trained model

```



```

109 coefficients = best_model.coef_[0]
110
111 # get the feature names
112 feature_names = X.columns
113
114 # create a dictionary to store the coefficients
115 feature_coefs = dict(zip(feature_names, coefficients))
116
117 # sort the feature coefficients by their absolute values in descending order
118 sorted_feature_importances = sorted(feature_coefs.items(), key=lambda x: abs(x[1]),
    reverse=True)
119
120 # print the feature coefficients
121 print("Feature Coefficients:")
122 for feature, coefficient in sorted_feature_importances:
123     print(f"{feature}: Coefficient = {coefficient:.4f}")
124
125 # permuted VIM
126 from sklearn.inspection import permutation_importance
127 import matplotlib.pyplot as plt
128
129 # initialize variables to store the permutation importances and standard deviations
130 all_importances = []
131 all_std_devs = []
132
133 # perform cross-validation
134 for train_index, test_index in cv.split(X):
135     X_train, X_test = X.iloc[train_index], X.iloc[test_index]
136     y_train, y_test = y_encoded[train_index], y_encoded[test_index]
137
138     # Fit the model on the training data
139     model.fit(X_train, y_train)
140
141     # Compute permuted variable importance measure for the test set
142     result = permutation_importance(model, X_test, y_test, scoring='accuracy',
    n_repeats=10, random_state=0, n_jobs=-1)
143
144     # Store the importances and standard deviations for each fold
145     all_importances.append(result.importances_mean)
146     all_std_devs.append(result.importances_std)
147
148 # compute the average importances and standard deviations across all folds
149 avg_importances = np.mean(all_importances, axis=0)
150 avg_std_devs = np.mean(all_std_devs, axis=0)
151
152 # create a DataFrame to store the VIM and standard deviations
153 feature_importance = pd.DataFrame({
154     'Feature': X.columns,
155     'Importance': avg_importances,
156     'Standard Deviation': avg_std_devs
157 })
158
159 # sort the feature importances in descending order
160 feature_importance = feature_importance.sort_values('Importance', ascending=True)
161
162 # plot the feature importances with standard deviations
163 ax = feature_importance.plot(x='Feature', y='Importance', kind='barh', figsize=(10,
    6), xerr='Standard Deviation', capsize=4)
164 ax.set_xlabel('Permutation Importance')
165 ax.set_title('Permutation Importance with Standard Deviation (Cross-Validation)')
166 plt.tight_layout()
167 plt.show()
168
169 # construct ALE plots
170 from alibi.explainers import ALE, plot_ale
171 logit_fun_lr = best_model.decision_function
172 proba_fun_lr = best_model.predict_proba
173 logit_ale_lr = ALE(logit_fun_lr, feature_names=feature_names, target_names=[str(y.

```

```

        name)])
174 proba_ale_lr = ALE(proba_fun_lr, feature_names=feature_names, target_names=[str(y.
        name)])
175 logit_exp_lr = logit_ale_lr.explain(X.values)
176 proba_exp_lr = proba_ale_lr.explain(X.values)
177 plot_ale(logit_exp_lr, n_cols=3, fig_kw={'figwidth': 8, 'figheight': 6}, sharey=None)
178
179 #####
180 # Model 2: Naive Bayes
181 #####
182
183 from sklearn.naive_bayes import GaussianNB
184 from sklearn.model_selection import GridSearchCV, RepeatedKFold
185
186 # create model
187 model = GaussianNB()
188
189 # define the hyperparameter grid
190 param_grid = {'var_smoothing':[1e-2, 5e-2, 1e-1, 5e-1, 1, 5, 10]}
191
192 # create cv object
193 cv = RepeatedKFold(n_splits=10, n_repeats=10, random_state=0)
194
195 # create grid-search cv object
196 grid_search = GridSearchCV(estimator=model, param_grid=param_grid, cv=cv, scoring='
    accuracy')
197
198 # perform grid-search
199 grid_search.fit(X, y_encoded)
200
201 # get the best hyperparameters
202 best_params = grid_search.best_params_
203 best_score = grid_search.best_score_
204 best_model = grid_search.best_estimator_
205
206 # best hyperparameters
207 best_params
208
209 # best CV accuracy
210 best_score
211
212 from sklearn.metrics import accuracy_score
213
214 y_pred = best_model.predict(X)
215 # best model training error
216 accuracy_score(y_encoded, y_pred)
217
218 # permuted VIM
219 from sklearn.inspection import permutation_importance
220 import matplotlib.pyplot as plt
221
222 # initialize variables to store the permutation importances and standard deviations
223 all_importances = []
224 all_std_devs = []
225
226 # perform cross-validation
227 for train_index, test_index in cv.split(X):
228     X_train, X_test = X.iloc[train_index], X.iloc[test_index]
229     y_train, y_test = y_encoded[train_index], y_encoded[test_index]
230
231     # Fit the model on the training data
232     model.fit(X_train, y_train)
233
234     # Compute permuted variable importance measure for the test set
235     result = permutation_importance(model, X_test, y_test, scoring='accuracy',
        n_repeats=10, random_state=0, n_jobs=-1)
236
237     # Store the importances and standard deviations for each fold

```

```

238     all_importances.append(result.importances_mean)
239     all_std_devs.append(result.importances_std)
240
241 # compute the average importances and standard deviations across all folds
242 avg_importances = np.mean(all_importances, axis=0)
243 avg_std_devs = np.mean(all_std_devs, axis=0)
244
245 # create a DataFrame to store the VIM and standard deviations
246 feature_importance = pd.DataFrame({
247     'Feature': X.columns,
248     'Importance': avg_importances,
249     'Standard Deviation': avg_std_devs
250 })
251
252 # sort the feature importances in descending order
253 feature_importance = feature_importance.sort_values('Importance', ascending=True)
254
255 # plot the feature importances with standard deviations
256 ax = feature_importance.plot(x='Feature', y='Importance', kind='barh', figsize=(10,
257     6), xerr='Standard Deviation', capsize=4)
258 ax.set_xlabel('Permutation Importance')
259 ax.set_title('Permutation Importance with Standard Deviation (Cross-Validation)')
260 plt.tight_layout()
261 plt.show()
262 #####
263 # Model 3: LDA
264 #####
265
266 from sklearn.discriminant_analysis import LinearDiscriminantAnalysis
267 from sklearn.model_selection import cross_val_score, RepeatedKFold
268
269 # create model
270 model = LinearDiscriminantAnalysis(shrinkage='auto', solver='lsqr')
271
272 # create cv object
273 cv = RepeatedKFold(n_splits=10, n_repeats=10, random_state=0)
274
275 # get the cv score
276 cv_score = cross_val_score(model, X, y_encoded, cv=cv, scoring='accuracy')
277
278 # best CV accuracy
279 cv_score.mean()
280
281 # permuted VIM
282 from sklearn.inspection import permutation_importance
283 import matplotlib.pyplot as plt
284
285 # initialize variables to store the permutation importances and standard deviations
286 all_importances = []
287 all_std_devs = []
288
289 # perform cross-validation
290 for train_index, test_index in cv.split(X):
291     X_train, X_test = X.iloc[train_index], X.iloc[test_index]
292     y_train, y_test = y_encoded[train_index], y_encoded[test_index]
293
294     # Fit the model on the training data
295     model.fit(X_train, y_train)
296
297     # Compute permuted variable importance measure for the test set
298     result = permutation_importance(model, X_test, y_test, scoring='accuracy',
299         n_repeats=10, random_state=0, n_jobs=-1)
300
301     # Store the importances and standard deviations for each fold
302     all_importances.append(result.importances_mean)
303     all_std_devs.append(result.importances_std)
304

```

```

304 # compute the average importances and standard deviations across all folds
305 avg_importances = np.mean(all_importances, axis=0)
306 avg_std_devs = np.mean(all_std_devs, axis=0)
307
308 # create a DataFrame to store the VIM and standard deviations
309 feature_importance = pd.DataFrame({
310     'Feature': X.columns,
311     'Importance': avg_importances,
312     'Standard Deviation': avg_std_devs
313 })
314
315 # sort the feature importances in descending order
316 feature_importance = feature_importance.sort_values('Importance', ascending=True)
317
318 # plot the feature importances with standard deviations
319 ax = feature_importance.plot(x='Feature', y='Importance', kind='barh', figsize=(10,
320     6), xerr='Standard Deviation', capsize=4)
321 ax.set_xlabel('Permutation Importance')
322 ax.set_title('Permutation Importance with Standard Deviation (Cross-Validation)')
323 plt.tight_layout()
324 plt.show()
325
326 #####
327 # Model 4: SVM
328 #####
329
330 from sklearn.model_selection import GridSearchCV, RepeatedKFold
331 from sklearn.neighbors import KNeighborsClassifier
332
333 # set model
334 knn = KNeighborsClassifier()
335
336 # search parameters
337 param_grid = {'n_neighbors': range(1, 31)}
338
339 # grid search to find the best parameter with K-fold CV.
340 kfold = RepeatedKFold(n_splits=10, n_repeats=5, random_state=0)
341 grid_search = GridSearchCV(knn, param_grid, cv=kfold, scoring='accuracy', verbose=1)
342
343 # Fit the model
344 grid_search.fit(X, y_encoded)
345
346 # best model
347 best_knn_params = grid_search.best_params_
348 print("Best Parameters: ", best_knn_params)
349 print("Best Cross-Validated Score", grid_search.best_score_)
350
351 # plot parameter tuning curve
352 cv_results_df = pd.DataFrame(grid_search.cv_results_)
353 n_neighbors = cv_results_df['param_n_neighbors']
354 CV_acc = cv_results_df['mean_test_score']
355
356 plt.figure(figsize=(10, 6))
357 plt.plot(n_neighbors, CV_acc, marker='o')
358 plt.title('Cross-Validated Accuracy vs. n_neighbors')
359 plt.xlabel('n_neighbors')
360 plt.ylabel('Mean CV Accuracy')
361 plt.grid(True)
362 plt.savefig("knn-tuning.png", format="png", bbox_inches="tight", dpi=300)
363 plt.show()
364
365 ##### VIM importance
366 from sklearn.inspection import permutation_importance
367
368 # Initialize variables to store the permutation importances and standard deviations
369 all_importances = []
370 all_std_devs = []
371

```

```

371 # Perform cross-validation
372 for train_index, test_index in kfold.split(X):
373     X_train, X_test = X.iloc[train_index], X.iloc[test_index]
374     y_train, y_test = y_encoded[train_index], y_encoded[test_index]
375
376     # Fit the model on the training data
377     knn.fit(X_train, y_train)
378
379     # Compute permuted variable importance measure for the test set
380     result = permutation_importance(knn, X_test, y_test, scoring='accuracy',
381                                     n_repeats=10, random_state=0, n_jobs=-1)
382
383     # Store the importances and standard deviations for each fold
384     all_importances.append(result.importances_mean)
385     all_std_devs.append(result.importances_std)
386
387 # Compute the average importances and standard deviations across all folds
388 avg_importances = np.mean(all_importances, axis=0)
389 avg_std_devs = np.mean(all_std_devs, axis=0)
390
391 # create a DataFrame to store the VIM and standard deviations
392 feature_importance = pd.DataFrame({
393     'Feature': X.columns,
394     'Importance': avg_importances,
395     'Standard Deviation': avg_std_devs
396 })
397
398 # sort the feature importances in descending order
399 feature_importance = feature_importance.sort_values('Importance', ascending=True)
400
401 # plot the feature importances with standard deviations
402 ax = feature_importance.plot(x='Feature', y='Importance', kind='barh', figsize=(10,
403                                     6), xerr='Standard Deviation', capsize=4)
404 ax.set_xlabel('Permutation Importance')
405 ax.set_title('Permutation Importance with Standard Deviation (Cross-Validation)')
406 plt.tight_layout()
407 plt.savefig("knn-vim.png", format="png", bbox_inches="tight", dpi=300)
408 plt.show()
409
410 #####
411 # Model 5: KNN
412 #####
413
414 from sklearn.model_selection import GridSearchCV, RepeatedKFold
415 from sklearn.svm import SVC
416
417 # set model
418 svc = SVC()
419
420 #####
421 # 1st search
422 #####
423
424 # search parameters
425 param_grid = {'C': [0.1, 1, 10, 50, 100, 1000], 'kernel': ['rbf']}
426
427 # grid search to find the best parameter with K-fold CV.
428 kfold = RepeatedKFold(n_splits=10, n_repeats=5, random_state=0)
429 grid_search = GridSearchCV(svc, param_grid, cv=kfold, scoring='accuracy', verbose=1)
430
431 # Fit the model
432 grid_search.fit(X, y_encoded)
433
434 # best model
435 best_svc_params = grid_search.best_params_
436 print("Best Parameters: ", best_svc_params)
437 print("Best Cross-Validated Score", grid_search.best_score_)
438

```

```

437 #####
438 # 2nd search
439 #####
440
441 # search parameters
442 param_grid = {'C': [10, 20, 30, 40, 50, 60, 70, 80, 90, 100], 'kernel': ['rbf']}
443
444 # grid search to find the best parameter with K-fold CV.
445 kfold = RepeatedKFold(n_splits=10, n_repeats=5, random_state=0)
446 grid_search = GridSearchCV(svc, param_grid, cv=kfold, scoring='accuracy', verbose=1)
447
448 # Fit the model
449 grid_search.fit(X, y_encoded)
450
451 # best model
452 best_svc_params = grid_search.best_params_
453 print("Best Parameters: ", best_svc_params)
454 print("Best Cross-Validated Score", grid_search.best_score_)
455
456 #####
457 # 3rd search
458 #####
459
460 # search parameters
461 param_grid = {'C': range(30,51), 'kernel': ['rbf']}
462
463 # grid search to find the best parameter with K-fold CV.
464 kfold = RepeatedKFold(n_splits=10, n_repeats=5, random_state=0)
465 grid_search = GridSearchCV(svc, param_grid, cv=kfold, scoring='accuracy', verbose=1)
466
467 # Fit the model
468 grid_search.fit(X, y_encoded)
469
470 # best model
471 best_svc_params = grid_search.best_params_
472 print("Best Parameters: ", best_svc_params)
473 print("Best Cross-Validated Score", grid_search.best_score_)
474
475
476 ##### VIM importance
477
478 svc = SVC(**best_svc_params)
479
480 # Initialize variables to store the permutation importances and standard deviations
481 all_importances = []
482 all_std_devs = []
483
484 # Perform cross-validation
485 for train_index, test_index in kfold.split(X):
486     X_train, X_test = X.iloc[train_index], X.iloc[test_index]
487     y_train, y_test = y_encoded[train_index], y_encoded[test_index]
488
489     # Fit the model on the training data
490     svc.fit(X_train, y_train)
491
492     # Compute permuted variable importance measure for the test set
493     result = permutation_importance(svc, X_test, y_test, scoring='accuracy',
494                                     n_repeats=10, random_state=0, n_jobs=-1)
495
496     # Store the importances and standard deviations for each fold
497     all_importances.append(result.importances_mean)
498     all_std_devs.append(result.importances_std)
499
500 # Compute the average importances and standard deviations across all folds
501 avg_importances = np.mean(all_importances, axis=0)
502 avg_std_devs = np.mean(all_std_devs, axis=0)
503
504 # create a DataFrame to store the VIM and standard deviations

```

```

504 feature_importance = pd.DataFrame({
505     'Feature': X.columns,
506     'Importance': avg_importances,
507     'Standard Deviation': avg_std_devs
508 })
509
510 # sort the feature importances in descending order
511 feature_importance = feature_importance.sort_values('Importance', ascending=True)
512
513 # plot the feature importances with standard deviations
514 ax = feature_importance.plot(x='Feature', y='Importance', kind='barh', figsize=(10,
515     6), xerr='Standard Deviation', capsize=4)
516 ax.set_xlabel('Permutation Importance')
517 ax.set_title('Permutation Importance with Standard Deviation (Cross-Validation)')
518 plt.tight_layout()
519 plt.savefig("svc-vim.png", format="png", bbox_inches="tight", dpi=300)
520 plt.show()
521
522 #####
523 # Model 6: neural network
524 #####
525
526 from sklearn.neural_network import MLPClassifier
527 from sklearn.model_selection import GridSearchCV, RepeatedKFold
528
529 # create model
530 model = MLPClassifier(random_state=0, activation='logistic', solver='adam', max_iter
531     =1000, learning_rate='adaptive')
532
533 # define the hyperparameter grid
534 param_grid = {'hidden_layer_sizes':[(50,), (100,)],
535     'alpha':[0.0001, 0.001, 0.01],
536     'learning_rate_init':[0.01, 0.1]}
537
538 # create cv object
539 cv = RepeatedKFold(n_splits=10, n_repeats=1, random_state=0)
540
541 # create grid-search cv object
542 grid_search = GridSearchCV(estimator=model, param_grid=param_grid, cv=cv, scoring='
543     accuracy')
544
545 # perform grid-search
546 grid_search.fit(X, y_encoded)
547
548 # get the best hyperparameters
549 best_params = grid_search.best_params_
550 best_score = grid_search.best_score_
551 best_model = grid_search.best_estimator_
552
553 # best hyperparameters
554 best_params
555
556 # best CV accuracy
557 best_score
558
559 from sklearn.metrics import accuracy_score
560
561 y_pred = best_model.predict(X)
562 # best model training error
563 accuracy_score(y_encoded, y_pred)
564
565 # permuted VIM
566 from sklearn.inspection import permutation_importance
567 import matplotlib.pyplot as plt
568
569 # initialize variables to store the permutation importances and standard deviations
570 all_importances = []
571 all_std_devs = []

```

```

569
570 # perform cross-validation
571 for train_index, test_index in cv.split(X):
572     X_train, X_test = X.iloc[train_index], X.iloc[test_index]
573     y_train, y_test = y_encoded[train_index], y_encoded[test_index]
574
575     # Fit the model on the training data
576     model.fit(X_train, y_train)
577
578     # Compute permuted variable importance measure for the test set
579     result = permutation_importance(model, X_test, y_test, scoring='accuracy',
580                                   n_repeats=10, random_state=0, n_jobs=-1)
581
582     # Store the importances and standard deviations for each fold
583     all_importances.append(result.importances_mean)
584     all_std_devs.append(result.importances_std)
585
586 # compute the average importances and standard deviations across all folds
587 avg_importances = np.mean(all_importances, axis=0)
588 avg_std_devs = np.mean(all_std_devs, axis=0)
589
590 # create a DataFrame to store the VIM and standard deviations
591 feature_importance = pd.DataFrame({
592     'Feature': X.columns,
593     'Importance': avg_importances,
594     'Standard Deviation': avg_std_devs
595 })
596
597 # sort the feature importances in descending order
598 feature_importance = feature_importance.sort_values('Importance', ascending=True)
599
600 # plot the feature importances with standard deviations
601 ax = feature_importance.plot(x='Feature', y='Importance', kind='barh', figsize=(10,
602     6), xerr='Standard Deviation', capsize=4)
603 ax.set_xlabel('Permutation Importance')
604 ax.set_title('Permutation Importance with Standard Deviation (Cross-Validation)')
605 plt.tight_layout()
606 plt.show()
607
608 #####
609 # Model 7: Classification Tree
610 #####
611
612 from sklearn.tree import DecisionTreeClassifier
613 from sklearn.model_selection import GridSearchCV, RepeatedKFold
614
615 # set model
616 clt = DecisionTreeClassifier()
617
618 #####
619 # 1st search
620 #####
621
622 # search parameters
623 param_grid = {'max_depth': [5,10,15,20], 'min_samples_split': [2,3,5,10], '
624     min_samples_leaf': [2,3,5,10], 'ccp_alpha': [0.0,0.001,0.01,0.1]}
625 kfold = RepeatedKFold(n_splits=10, n_repeats=5, random_state=0)
626 grid_search = GridSearchCV(clt, param_grid, cv=kfold, scoring='accuracy', verbose=1)
627 grid_search.fit(X, y)
628
629 # best model
630 best_clt_params = grid_search.best_params_
631 print("Best Parameters: ", best_clt_params)
632 print("Best Cross-Validated Score", grid_search.best_score_)
633
634 #####
635 # 2nd search
636 #####

```



```

634
635 # search parameters
636 param_grid = {'max_depth': [13,15,17,19], 'min_samples_split': [4,5,7,9], '
        min_samples_leaf': [2,3,4], 'ccp_alpha': [0.0,0.0001,0.0005,0.001]}
637 kfold = RepeatedKFold(n_splits=10, n_repeats=5, random_state=0)
638 grid_search = GridSearchCV(clt, param_grid, cv=kfold, scoring='accuracy', verbose=1)
639 grid_search.fit(X, y)
640
641 # best model
642 best_clt_params = grid_search.best_params_
643 print("Best Parameters: ", best_clt_params)
644 print("Best Cross-Validated Score", grid_search.best_score_)
645
646 #####
647 # 3rd search
648 #####
649
650 # search parameters
651 param_grid = {'max_depth': [14,15,16], 'min_samples_split': [6,7,8], '
        min_samples_leaf': [3,4], 'ccp_alpha': [0.0002,0.0003,0.0005,0.0007,0.0009]}
652 kfold = RepeatedKFold(n_splits=10, n_repeats=5, random_state=0)
653 grid_search = GridSearchCV(clt, param_grid, cv=kfold, scoring='accuracy', verbose=1)
654 grid_search.fit(X, y)
655
656 # best model
657 best_clt_params = grid_search.best_params_
658 print("Best Parameters: ", best_clt_params)
659 print("Best Cross-Validated Score", grid_search.best_score_)
660
661 #####
662 # 4th search
663 #####
664
665 # search parameters
666 param_grid = {'max_depth': [15,16,17], 'min_samples_split': [7,8,9], '
        min_samples_leaf': [3,4,5], 'ccp_alpha': [0.0006,0.0007,0.0008]}
667 kfold = RepeatedKFold(n_splits=10, n_repeats=5, random_state=0)
668 grid_search = GridSearchCV(clt, param_grid, cv=kfold, scoring='accuracy', verbose=1)
669 grid_search.fit(X, y)
670
671 # best model
672 best_clt_params = grid_search.best_params_
673 print("Best Parameters: ", best_clt_params)
674 print("Best Cross-Validated Score", grid_search.best_score_)
675
676
677 ##### Built-in feature_importance_ function
678 clt = DecisionTreeClassifier(**best_clt_params)
679 clt.fit(X,y)
680
681 feature_importances = clt.feature_importances_
682 feature_names = X.columns
683
684 # print feature_importance_
685 feature_importance_dict = dict(zip(feature_names, feature_importances))
686 ranked_features = sorted(feature_importance_dict.items(), key=lambda x: x[1], reverse
        =True)
687 print("Ranked Features by Importance:")
688 for rank, (feature, importance) in enumerate(ranked_features, start=1):
689     print(f"{rank}. {feature}: {importance:.4f}")
690
691
692 ##### VIM importance
693 clt = DecisionTreeClassifier(**best_clt_params)
694
695 # Initialize variables to store the permutation importances and standard deviations
696 all_importances = []
697 all_std_devs = []

```

```

698
699 # Perform cross-validation
700 for train_index, test_index in kfold.split(X):
701     X_train, X_test = X.iloc[train_index], X.iloc[test_index]
702     y_train, y_test = y_encoded[train_index], y_encoded[test_index]
703
704     # Fit the model on the training data
705     clt.fit(X_train, y_train)
706
707     # Compute permuted variable importance measure for the test set
708     result = permutation_importance(clt, X_test, y_test, scoring='accuracy',
709                                     n_repeats=10, random_state=0, n_jobs=-1)
710
711     # Store the importances and standard deviations for each fold
712     all_importances.append(result.importances_mean)
713     all_std_devs.append(result.importances_std)
714
715 # Compute the average importances and standard deviations across all folds
716 avg_importances = np.mean(all_importances, axis=0)
717 avg_std_devs = np.mean(all_std_devs, axis=0)
718
719 # create a DataFrame to store the VIM and standard deviations
720 feature_importance = pd.DataFrame({
721     'Feature': X.columns,
722     'Importance': avg_importances,
723     'Standard Deviation': avg_std_devs
724 })
725
726 # sort the feature importances in descending order
727 feature_importance = feature_importance.sort_values('Importance', ascending=True)
728
729 # plot the feature importances with standard deviations
730 ax = feature_importance.plot(x='Feature', y='Importance', kind='barh', figsize=(10,
731     6), xerr='Standard Deviation', capsize=4)
732 ax.set_xlabel('Permutation Importance')
733 ax.set_title('Permutation Importance with Standard Deviation (Cross-Validation)')
734 plt.tight_layout()
735 plt.savefig("ct-vim.png", format="png", bbox_inches="tight", dpi=300)
736 plt.show()
737
738 #####
739 # Model 8: Light GBM
740 #####
741
742 import lightgbm as lgb
743 from sklearn.model_selection import GridSearchCV, RepeatedKFold
744 import warnings
745
746 # set model
747 lgbc = lgb.LGBMClassifier(num_threads=8, verbose=-1)
748
749 # warnings.filterwarnings("ignore", category=UserWarning, module="lightgbm")
750
751 #####
752 # 1st search
753 #####
754
755 param_grid = {
756     'learning_rate': [0.01, 0.1, 0.5],
757     'n_estimators': [200, 300],
758     'max_depth': [10, 20, 30],
759     'num_leaves': [15, 31, 63],
760 }
761
762 kfold = RepeatedKFold(n_splits=10, n_repeats=5, random_state=0)
763 grid_search = GridSearchCV(lgbc, param_grid, cv=kfold, scoring='accuracy', verbose=1)
764 grid_search.fit(X, y)

```

```

764
765 # best model
766 best_lgbc_params = grid_search.best_params_
767 print("Best Parameters: ", best_lgbc_params)
768 print("Best Cross-Validated Score", grid_search.best_score_)
769
770 #####
771 # 2nd search
772 #####
773
774 param_grid = {
775     'learning_rate': [0.05, 0.1, 0.2],
776     'n_estimators': [300, 400],
777     'max_depth': [30, 40],
778     'num_leaves': [63, 96],
779 }
780
781
782 kfold = RepeatedKFold(n_splits=10, n_repeats=5, random_state=0)
783 grid_search = GridSearchCV(lgbc, param_grid, cv=kfold, scoring='accuracy', verbose=1)
784 grid_search.fit(X, y)
785
786 # best model
787 best_lgbc_params = grid_search.best_params_
788 print("Best Parameters: ", best_lgbc_params)
789 print("Best Cross-Validated Score", grid_search.best_score_)
790
791 #####
792 # 3rd search
793 #####
794
795 param_grid = {
796     'learning_rate': [0.18, 0.2, 0.22],
797     'n_estimators': [300, 400],
798     'max_depth': [30],
799     'num_leaves': [96],
800 }
801
802
803 kfold = RepeatedKFold(n_splits=10, n_repeats=5, random_state=0)
804 grid_search = GridSearchCV(lgbc, param_grid, cv=kfold, scoring='accuracy', verbose=1)
805 grid_search.fit(X, y)
806
807 # best model
808 best_lgbc_params = grid_search.best_params_
809 print("Best Parameters: ", best_lgbc_params)
810 print("Best Cross-Validated Score", grid_search.best_score_)
811
812 #####
813 # 4th search
814 #####
815
816 param_grid = {
817     'learning_rate': [0.22, 0.24, 0.26, 0.28, 0.3],
818     'n_estimators': [300, 400, 500],
819     'max_depth': [30],
820     'num_leaves': [96],
821 }
822
823 kfold = RepeatedKFold(n_splits=10, n_repeats=5, random_state=0)
824 grid_search = GridSearchCV(lgbc, param_grid, cv=kfold, scoring='accuracy', verbose=1)
825 grid_search.fit(X, y)
826
827 # best model
828 best_lgbc_params = grid_search.best_params_
829 print("Best Parameters: ", best_lgbc_params)
830 print("Best Cross-Validated Score", grid_search.best_score_)
831

```

```

832
833 ##### built-in feature importance
834 lgbc = lgb.LGBMClassifier(**best_lgbc_params, num_threads=8, verbose=-1)
835 lgbc.fit(X,y)
836
837 lgb.plot_importance(lgbc, importance_type="gain", figsize=(7,6), title="LGBM Feature
      Importance")
838 plt.savefig("lgbm-importance.png", format="png", bbox_inches="tight", dpi=300)
839 plt.show()
840
841
842 ##### VIM importance
843 lgbc = lgb.LGBMClassifier(**best_lgbc_params, num_threads=8, verbose=-1)
844
845 # Initialize variables to store the permutation importances and standard deviations
846 all_importances = []
847 all_std_devs = []
848
849 # Perform cross-validation
850 for train_index, test_index in kfold.split(X):
851     X_train, X_test = X.iloc[train_index], X.iloc[test_index]
852     y_train, y_test = y_encoded[train_index], y_encoded[test_index]
853
854     # Fit the model on the training data
855     lgbc.fit(X_train, y_train)
856
857     # Compute permuted variable importance measure for the test set
858     result = permutation_importance(lgbc, X_test, y_test, scoring='accuracy',
      n_repeats=10, random_state=0, n_jobs=-1)
859
860     # Store the importances and standard deviations for each fold
861     all_importances.append(result.importances_mean)
862     all_std_devs.append(result.importances_std)
863
864 # Compute the average importances and standard deviations across all folds
865 avg_importances = np.mean(all_importances, axis=0)
866 avg_std_devs = np.mean(all_std_devs, axis=0)
867
868 # create a DataFrame to store the VIM and standard deviations
869 feature_importance = pd.DataFrame({
870     'Feature': X.columns,
871     'Importance': avg_importances,
872     'Standard Deviation': avg_std_devs
873 })
874
875 # sort the feature importances in descending order
876 feature_importance = feature_importance.sort_values('Importance', ascending=True)
877
878 # plot the feature importances with standard deviations
879 ax = feature_importance.plot(x='Feature', y='Importance', kind='barh', figsize=(10,
      6), xerr='Standard Deviation', capsize=4)
880 ax.set_xlabel('Permutation Importance')
881 ax.set_title('Permutation Importance with Standard Deviation (Cross-Validation)')
882 plt.tight_layout()
883 plt.savefig("lgbm-vim.png", format="png", bbox_inches="tight", dpi=300)
884 plt.show()
885
886 #####
887 # Model 9: Random Forest
888 #####
889
890 from sklearn.ensemble import RandomForestClassifier
891 from sklearn.model_selection import GridSearchCV, RepeatedKFold
892
893 # create model
894 model = RandomForestClassifier(random_state=0, n_jobs=-1)
895
896 # define the hyperparameter grid

```

```

897 param_grid = {'n_estimators':[250, 300, 350],
898               'max_depth':[20, 25, 30],
899               'min_samples_split':[2, 5, 10],
900               'min_samples_leaf':[1, 2, 5]}
901
902 # create cv object
903 cv = RepeatedKFold(n_splits=10, n_repeats=1, random_state=0)
904
905 # create grid-search cv object
906 grid_search = GridSearchCV(estimator=model, param_grid=param_grid, cv=cv, scoring='
    accuracy')
907
908 # perform grid-search
909 grid_search.fit(X, y_encoded)
910
911 # get the best hyperparameters
912 best_params = grid_search.best_params_
913 best_score = grid_search.best_score_
914 best_model = grid_search.best_estimator_
915
916 # best hyperparameters
917 best_params
918
919 # best CV accuracy
920 best_score
921
922 from sklearn.metrics import accuracy_score
923
924 y_pred = best_model.predict(X)
925 # best model training error
926 accuracy_score(y_encoded, y_pred)
927
928 # get the feature importances from the best model
929 importances = best_model.feature_importances_
930
931 # sort the feature importances in descending order
932 indices = np.argsort(importances)[::-1]
933
934 # rearrange the feature names based on the sorted importances
935 feature_names_new = [X.columns[i] for i in indices]
936
937 # create a bar plot
938 plt.figure(figsize=(10, 6))
939 plt.title("Feature Importances")
940 plt.barh(range(len(importances)), importances[indices])
941 plt.yticks(range(len(importances)), feature_names_new)
942 plt.xlabel("Importance")
943 plt.ylabel("Features")
944 plt.gca().invert_yaxis()
945 plt.tight_layout()
946 plt.show()
947
948 # initialize variables to store the permutation importances and standard deviations
949 all_importances = []
950 all_std_devs = []
951
952 # perform cross-validation
953 for train_index, test_index in cv.split(X):
954     X_train, X_test = X.iloc[train_index], X.iloc[test_index]
955     y_train, y_test = y_encoded[train_index], y_encoded[test_index]
956
957     # Fit the model on the training data
958     model.fit(X_train, y_train)
959
960     # Compute permuted variable importance measure for the test set
961     result = permutation_importance(model, X_test, y_test, scoring='accuracy',
    n_repeats=10, random_state=0, n_jobs=-1)
962

```

```

963     # Store the importances and standard deviations for each fold
964     all_importances.append(result.importances_mean)
965     all_std_devs.append(result.importances_std)
966
967 # compute the average importances and standard deviations across all folds
968 avg_importances = np.mean(all_importances, axis=0)
969 avg_std_devs = np.mean(all_std_devs, axis=0)
970
971 # create a DataFrame to store the VIM and standard deviations
972 feature_importance = pd.DataFrame({
973     'Feature': X.columns,
974     'Importance': avg_importances,
975     'Standard Deviation': avg_std_devs
976 })
977
978 # sort the feature importances in descending order
979 feature_importance = feature_importance.sort_values('Importance', ascending=True)
980
981 # plot the feature importances with standard deviations
982 ax = feature_importance.plot(x='Feature', y='Importance', kind='barh', figsize=(10,
983     6), xerr='Standard Deviation', capsize=4)
984 ax.set_xlabel('Permutation Importance')
985 ax.set_title('Permutation Importance with Standard Deviation (Cross-Validation)')
986 plt.tight_layout()
987 plt.show()

```