

## Problem 1

I created a data pipeline that was able to handle data from all four datasets. A function read in the dataset from CSV or scikit directly, and split into X and y. For the heart disease data, further preprocessing was required to one-hot encode several of the columns. Next, the X and y data were passed into my pipeline function, which performed a train/test split into x test, y test, x train, and y train. Next, missing values were imputed using the scikit SimpleImputer, as decision tree doesn't support missing values in the data. Next, a C4.5 decision tree, random forest, naive Bayes, and Support Vector Machine Model were created for all four data sets. After the models were each returned, I calculated the confusion matrices and associated metrics for each model, as well as ROC curves for the binary models (*cannot create ROC for non-binary classifiers*). Visualizations were saved and metrics saved to a .txt file. **Included figures and metrics in report are only for the wine dataset. All other data set figures and metrics are included in the attached zipped folder.**

model	sensitivity (recall)	specificity	balanced accuracy	precision	F1 Score
C4.5 Decision Tree	.947	.5	.333	.886	.975*
Naive Bayes	1	.525	.424	1	1*
Random Forest	1	.525	.424	1	1*
SVM	1	.513	.424	1	.975*

Table 1: Metrics for each model for the wine dataset

The calculated metrics indicate that several of the models had very high sensitivity, precision, and F1 score. The specificity is lower, and accordingly the balance accuracy is lower as well. All the models performed similarly. When we look at the confusion matrices for each model, we can see by the diagonal opacity that the models was classify most (or all) tuples correctly. Thusly, our true positive score was very high.

For the wine dataset, we do not have ROC curves as it is a multiclass (3-class) classification problem.

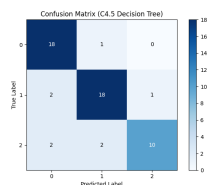


Figure 1: Wine dataset decision tree confusion matrix

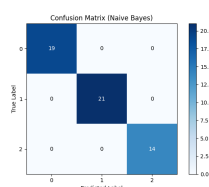


Figure 2: Wine dataset naive Bayes confusion matrix

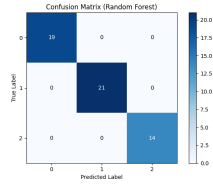


Figure 3: Wine dataset random forest confusion matrix

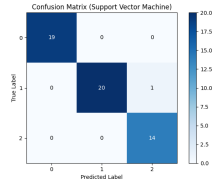


Figure 4: Wine dataset support vector machine confusion matrix

## HW8.py

---

```

1  #Disclaimer: parts or the whole of functions were generated using generative AI tools, then checked and modified for precision.
2
3  import os
4  import numpy as np
5  import pandas as pd
6  import matplotlib.pyplot as plt
7  from sklearn.model_selection import train_test_split
8  from sklearn.metrics import confusion_matrix, roc_curve, auc
9  from sklearn.tree import DecisionTreeClassifier
10 from sklearn.naive_bayes import GaussianNB
11 from sklearn.ensemble import RandomForestClassifier
12 from sklearn.svm import SVC
13 from sklearn.preprocessing import OneHotEncoder
14 from sklearn.impute import SimpleImputer
15 from sklearn.datasets import load_wine, load_breast_cancer, load_iris
16
17 def train_test_split_data(X, y, test_size=0.3, random_state=42):
18     """Split the dataset into training and testing sets."""
19     X_train, X_test, y_train, y_test = train_test_split(X, y, test_size=test_size, random_state=random_state)
20     return X_train, X_test, y_train, y_test
21
22 def c45_decision_tree(X_train, y_train, X_test):
23     """Train and evaluate a C4.5 decision tree model."""
24     model = DecisionTreeClassifier(criterion='entropy')
25     model.fit(X_train, y_train)
26     y_pred = model.predict(X_test)
27     return y_pred, model
28
29 def naive_bayes(X_train, y_train, X_test):
30     """Train and evaluate a naive Bayes model."""
31     model = GaussianNB()
32     model.fit(X_train, y_train)
33     y_pred = model.predict(X_test)
34     return y_pred, model
35
36 def random_forest(X_train, y_train, X_test):
37     """Train and evaluate a random forest model."""
38     model = RandomForestClassifier(n_estimators=100)
39     model.fit(X_train, y_train)

```

```

40     y_pred = model.predict(X_test)
41     return y_pred, model
42
43 def support_vector_machine(X_train, y_train, X_test):
44     """Train and evaluate a support vector machine model."""
45     model = SVC(kernel='linear')
46     model.fit(X_train, y_train)
47     y_pred = model.predict(X_test)
48     return y_pred, model
49
50 def calculate_metrics(y_true, y_pred):
51     """Calculate the confusion matrix and various classification metrics."""
52     cm = confusion_matrix(y_true, y_pred)
53     if cm.shape == (2, 2):
54         tn, fp, fn, tp = cm.ravel()
55         sensitivity = tp / (tp + fn)
56         specificity = tn / (tn + fp)
57         balanced_accuracy = (sensitivity + specificity) / 2
58         precision = tp / (tp + fp)
59         recall = sensitivity
60         f1_score = 2 * precision * recall / (precision + recall)
61         return cm, sensitivity, specificity, balanced_accuracy, precision, recall, f1_score
62     elif cm.shape == (3, 3):
63         # Calculate metrics for the three-class case
64         tp0 = cm[0, 0]
65         tp1 = cm[1, 1]
66         tp2 = cm[2, 2]
67         fp0 = np.sum(cm[1:, 0])
68         fp1 = np.sum(np.concatenate((cm[:, 1], cm[2:, 1])))
69         fp2 = np.sum(cm[:, 2])
70         fn0 = np.sum(cm[0, 1:])
71         fn1 = np.sum(np.concatenate((cm[0, :1], cm[0, 2:])))
72         fn2 = np.sum(cm[:, 2, 0])
73         tn0 = np.sum(cm[1:, 1:]) # includes class 1, class 2
74         tn1 = np.sum(np.concatenate((cm[:, 1, :1], cm[:, 1, 2:], cm[2:, :1], cm[2:, 2:])))
75         tn2 = np.sum(cm[:, 2, :2])
76         sensitivity0 = tp0 / (tp0 + fn0)
77         sensitivity1 = tp1 / (tp1 + fn1)
78         sensitivity2 = tp2 / (tp2 + fn2)
79         specificity0 = tn0 / (tn0 + fp0)
80         specificity1 = tn1 / (tn1 + fp1)
81         specificity2 = tn2 / (tn2 + fp2)
82         balanced_accuracy = (sensitivity0 + sensitivity1 + sensitivity2) / 3
83         precision0 = tp0 / (tp0 + fp0)
84         precision1 = tp1 / (tp1 + fp1)
85         precision2 = tp2 / (tp2 + fp2)
86         precision = (precision0 + precision1 + precision2) / 3
87         recall = (sensitivity0 + sensitivity1 + sensitivity2) / 3
88         f1_score = 2 * precision * recall / (precision + recall)
89         return cm, sensitivity0, sensitivity1, sensitivity2, specificity0, specificity1, specificity2, balanced_accuracy, precision, recall, f1_score
90     else:
91         print('Error: Confusion matrix has unexpected shape')
92         return None, None, None, None, None, None, None, None, None, None, None
93
94
95 def plot_confusion_matrix(cm, model_name, output_dir):
96     """Plot the confusion matrix."""
97     plt.figure()
98     plt.imshow(cm, interpolation='nearest', cmap=plt.cm.Blues)
99     plt.title(f'Confusion Matrix ({model_name})')
100    plt.colorbar()
101    tick_marks = np.arange(cm.shape[0])
102    plt.xticks(tick_marks, tick_marks)
103    plt.yticks(tick_marks, tick_marks)
104    plt.xlabel('Predicted Label')

```

```

105     plt.ylabel('True Label')
106     for i in range(cm.shape[0]):
107         for j in range(cm.shape[1]):
108             plt.text(j, i, format(cm[i, j], 'd'), ha="center", va="center", color="white" if cm[i, j] > cm.max() / 2 else "black")
109     plt.tight_layout()
110     plt.savefig(os.path.join(output_dir, 'confusion_matrix.png'))
111     plt.close()
112
113 def plot_roc_curve(y_true, y_score, model_name, output_dir):
114     """Plot the ROC curve."""
115     fpr, tpr, _ = roc_curve(y_true, y_score)
116     roc_auc = auc(fpr, tpr)
117
118     plt.figure()
119     plt.plot(fpr, tpr, color='darkorange', lw=2, label=f'ROC Curve (AUC = {roc_auc:.3f})')
120     plt.plot([0, 1], [0, 1], color='navy', lw=2, linestyle='--')
121     plt.xlim([0.0, 1.0])
122     plt.ylim([0.0, 1.05])
123     plt.xlabel('False Positive Rate')
124     plt.ylabel('True Positive Rate')
125     plt.title(f'ROC Curve ({model_name})')
126     plt.legend(loc="lower right")
127     plt.savefig(os.path.join(output_dir, 'roc_curve.png'))
128     plt.close()
129
130
131 def pipeline(X, y, dataset_name, output_dir='output'):
132     """Run the classification pipeline on the specified dataset."""
133     if not os.path.exists(output_dir):
134         os.makedirs(output_dir)
135
136     # Create a separate output folder for this dataset
137     dataset_output_dir = os.path.join(output_dir, dataset_name)
138     if not os.path.exists(dataset_output_dir):
139         os.makedirs(dataset_output_dir)
140
141     # Split the data into training and testing sets
142     X_train, X_test, y_train, y_test = train_test_split_data(X, y)
143
144     # Impute missing values
145     imputer = SimpleImputer(strategy='mean')
146     X_train = imputer.fit_transform(X_train)
147     X_test = imputer.transform(X_test)
148
149     # Train and evaluate the C4.5 decision tree model
150     y_pred, model = c45_decision_tree(X_train, y_train, X_test)
151     cm, *metrics = calculate_metrics(y_test, y_pred)
152     if cm is not None:
153         model_name = 'C4.5 Decision Tree'
154         model_output_dir = os.path.join(output_dir, dataset_name, model_name)
155         if not os.path.exists(model_output_dir):
156             os.makedirs(model_output_dir)
157         with open(os.path.join(model_output_dir, 'metrics.txt'), 'w') as f:
158             f.write(f'Confusion Matrix ({model_name}):\n{cm}\n\n')
159             f.write(f'Sensitivity: {metrics[0]:.3f}\n')
160             f.write(f'Specificity: {metrics[1]:.3f}\n')
161             f.write(f'Balanced Accuracy: {metrics[2]:.3f}\n')
162             f.write(f'Precision: {metrics[3]:.3f}\n')
163             f.write(f'Recall: {metrics[4]:.3f}\n')
164             f.write(f'F1 Score: {metrics[5]:.3f}\n')
165         y_score = model.predict_proba(X_test)[:, 1]
166         plot_confusion_matrix(cm, model_name, model_output_dir)
167         if cm.shape == (2, 2):
168             plot_roc_curve(y_test, y_score, model_name, model_output_dir)
169

```

```

170 # Train and evaluate the naive Bayes model
171 y_pred, model = naive_bayes(X_train, y_train, X_test)
172 cm, *metrics = calculate_metrics(y_test, y_pred)
173 if cm is not None:
174     model_name = 'Naive Bayes'
175     model_output_dir = os.path.join(output_dir, dataset_name, model_name)
176     if not os.path.exists(model_output_dir):
177         os.makedirs(model_output_dir)
178     with open(os.path.join(model_output_dir, 'metrics.txt'), 'w') as f:
179         f.write(f'Confusion Matrix ({model_name}): \n{cm}\n\n')
180         f.write(f'Sensitivity: {metrics[0]:.3f}\n')
181         f.write(f'Specificity: {metrics[1]:.3f}\n')
182         f.write(f'Balanced Accuracy: {metrics[2]:.3f}\n')
183         f.write(f'Precision: {metrics[3]:.3f}\n')
184         f.write(f'Recall: {metrics[4]:.3f}\n')
185         f.write(f'F1 Score: {metrics[5]:.3f}\n')
186     y_score = model.predict_proba(X_test)[: , 1]
187     plot_confusion_matrix(cm, model_name, model_output_dir)
188     if cm.shape == (2, 2):
189         plot_roc_curve(y_test, y_score, model_name, model_output_dir)
190
191 # Train and evaluate the random forest model
192 y_pred, model = random_forest(X_train, y_train, X_test)
193 cm, *metrics = calculate_metrics(y_test, y_pred)
194 if cm is not None:
195     model_name = 'Random Forest'
196     model_output_dir = os.path.join(output_dir, dataset_name, model_name)
197     if not os.path.exists(model_output_dir):
198         os.makedirs(model_output_dir)
199     with open(os.path.join(model_output_dir, 'metrics.txt'), 'w') as f:
200         f.write(f'Confusion Matrix ({model_name}): \n{cm}\n\n')
201         f.write(f'Sensitivity: {metrics[0]:.3f}\n')
202         f.write(f'Specificity: {metrics[1]:.3f}\n')
203         f.write(f'Balanced Accuracy: {metrics[2]:.3f}\n')
204         f.write(f'Precision: {metrics[3]:.3f}\n')
205         f.write(f'Recall: {metrics[4]:.3f}\n')
206         f.write(f'F1 Score: {metrics[5]:.3f}\n')
207     y_score = model.predict_proba(X_test)[: , 1]
208     plot_confusion_matrix(cm, model_name, model_output_dir)
209     if cm.shape == (2, 2):
210         plot_roc_curve(y_test, y_score, model_name, model_output_dir)
211
212 # Train and evaluate the support vector machine model
213 y_pred, model = support_vector_machine(X_train, y_train, X_test)
214 cm, *metrics = calculate_metrics(y_test, y_pred)
215 if cm is not None:
216     model_name = 'Support Vector Machine'
217     model_output_dir = os.path.join(output_dir, dataset_name, model_name)
218     if not os.path.exists(model_output_dir):
219         os.makedirs(model_output_dir)
220     with open(os.path.join(model_output_dir, 'metrics.txt'), 'w') as f:
221         f.write(f'Confusion Matrix ({model_name}): \n{cm}\n\n')
222         f.write(f'Sensitivity: {metrics[0]:.3f}\n')
223         f.write(f'Specificity: {metrics[1]:.3f}\n')
224         f.write(f'Balanced Accuracy: {metrics[2]:.3f}\n')
225         f.write(f'Precision: {metrics[3]:.3f}\n')
226         f.write(f'Recall: {metrics[4]:.3f}\n')
227         f.write(f'F1 Score: {metrics[5]:.3f}\n')
228     y_score = model.decision_function(X_test)
229     plot_confusion_matrix(cm, model_name, model_output_dir)
230     if cm.shape == (2, 2):
231         plot_roc_curve(y_test, y_score, model_name, model_output_dir)
232
233 def load_data(filename):
234     """Load the heart disease dataset."""

```

```

235     df = pd.read_csv(filename, usecols=["age", "sex", "cp", "trestbps", "chol", "fbs", "restecg", "thalch", "exang", "oldpeak",
236
237     # Replace the target variable values
238     df["num"] = df["num"].replace({1: 0, 2: 1, 3: 1, 4: 1})
239
240     # One-hot encode categorical variables
241     df = pd.get_dummies(df, columns=["sex", "cp", "fbs", "restecg", "exang", "slope", "ca", "thal"])
242
243     # Split the data into X and y
244     X = df.drop("num", axis=1)
245     y = df["num"]
246
247     return X, y
248
249 def main():
250     sklearn_datasets = [('wine', load_wine()), ('breast_cancer', load_breast_cancer()), ('iris', load_iris())]
251     for dataset_name, dataset in sklearn_datasets:
252         X = dataset.data
253         y = dataset.target
254         pipeline(X, y, dataset_name)
255     X, y = load_data('heart_disease_uci.csv')
256     pipeline(X, y, 'heart_disease_uci')
257
258
259 if __name__ == '__main__':
260     main()

```

---