



*University of Essex*

**Department of Mathematical Sciences**

---

COMBINED COURSEWORK REASSESSMENT FOR  
CE802

COMPARATIVE STUDY:  
INVESTIGATING THE EFFECTIVENESS  
OF MACHINE LEARNING FOR  
TARGETED EMAIL MARKETING

**WORD COUNT: 792**

**2202026**

---

August 23, 2023  
Colchester

---

# Contents

<b>1</b>	<b>Introduction</b>	<b>3</b>
<b>2</b>	<b>Data Exploration and Pre-processing</b>	<b>4</b>
<b>3</b>	<b>Results and Analysis</b>	<b>7</b>
3.1	Interpretation . . . . .	8
<b>4</b>	<b>Conclusions</b>	<b>9</b>
<b>5</b>	<b>Appendix</b>	<b>10</b>

---

## Introduction

This research delved into the evaluation of various machine learning methodologies aimed at predicting email success in generating sales. Leveraging historical data supplied by the company, this study explored the effectiveness of algorithms including Naive Bayes, Decision Trees, Logistic Regression, and Support Vector Machines (SVM). The overarching objective was to pinpoint the optimal strategy for enhancing email campaign precision and bolstering overall conversion rates.

## Data Exploration and Pre-processing

The data set required little to no pre-processing techniques as it consisted of numerical variables and had no missing values.

The plots below were done to identify false predictors. Most of the plots have a pattern which means that most of our independent variables have a relationship with the dependent variable which is the class variable. The only plot that doesn't have a pattern is the F4 plot.

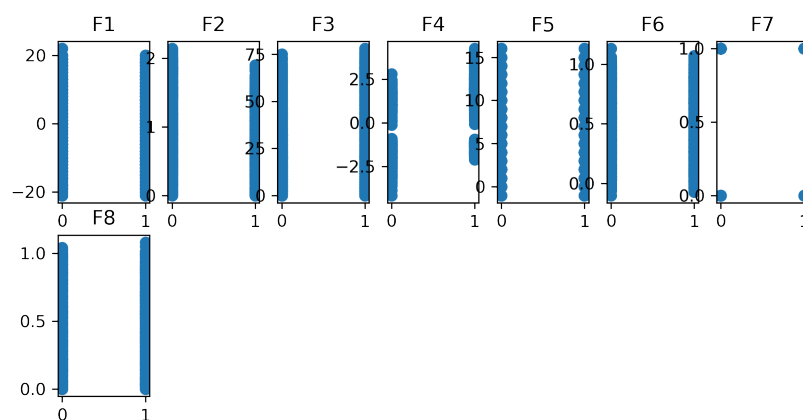


Figure 2.1: Scatter Plot of Features and Labels

The patterns shown in 2.1 also showcase that the features are not identical as they all have their own unique relationship with the labels.

Scaling of each input features was done to ensure there is fair comparison between them.

The figure below showcases how the features used had no significant outliers.

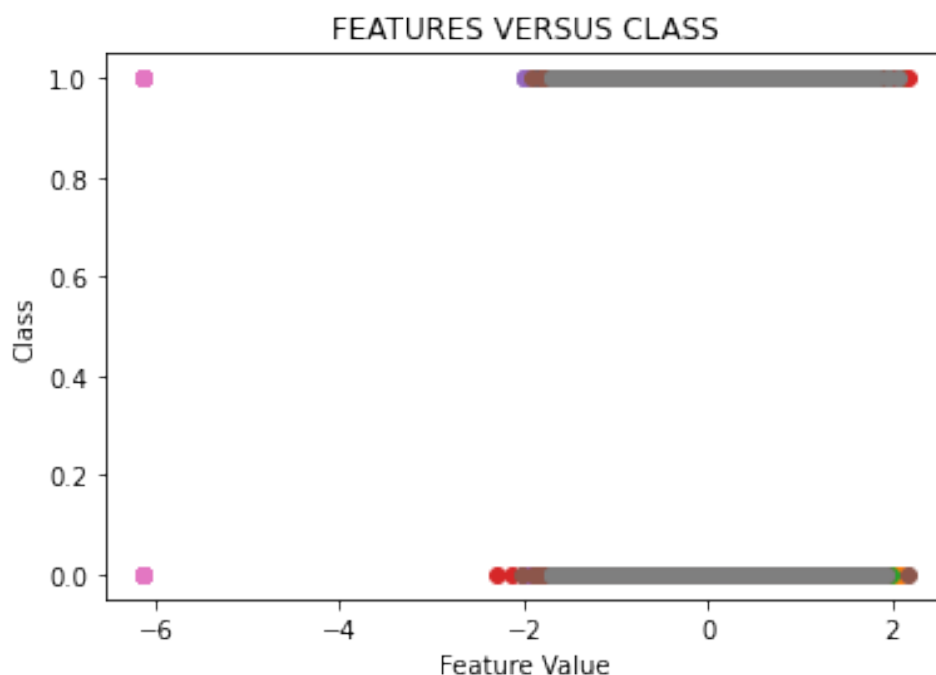


Figure 2.2: Features versus Class

The Scikit-learn library in Python was utilized to implement and evaluate the machine learning models. The dataset consisted of 1000 examples, each with 8 features and a binary label indicating whether the corresponding email led to a sale. The data set was randomly split into a training set (70%) and a testing set (30%) for model evaluation. For each algorithm, the model was trained on the training set and its performance was evaluated on the testing set.

The distribution of the classes is slightly evenly split as depicted in the chart below.

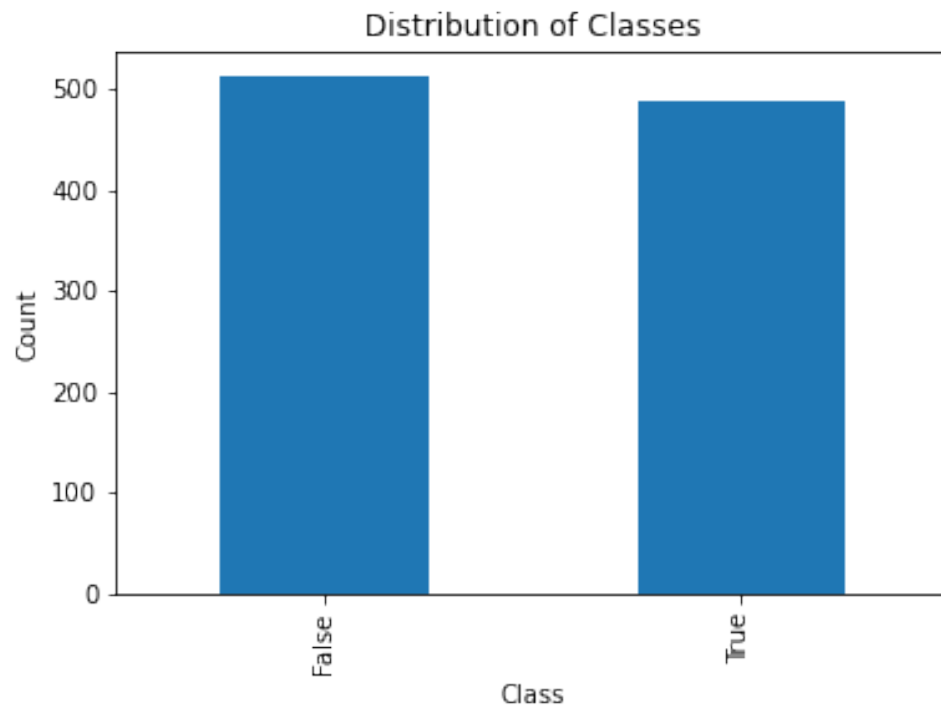


Figure 2.3: Distribution of the classes

## Results and Analysis

The performance of each algorithm was measured using the following key performance metrics accuracy, precision, recall, F1-score and Cohen's Kappa statistic.

The performance of these performance metrics on each classification technique is depicted in the chart below.

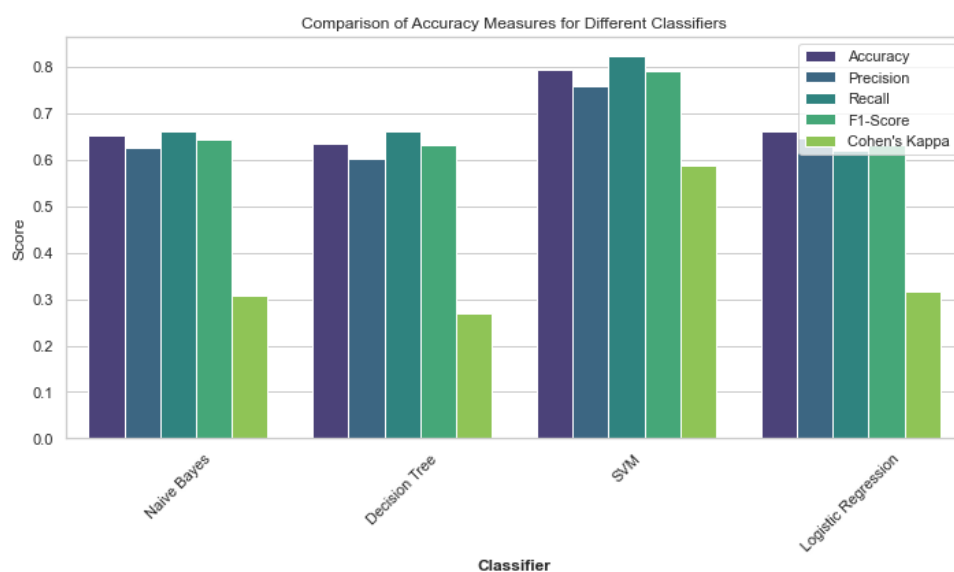


Figure 3.1: Comparison of Accuracy Measures for Different Classifiers

### 3.1 Interpretation

The performance metrics on 3.1 provide insights into the effectiveness of different classifiers in a classification task. The "Accuracy" measure signifies the proportion of correct predictions overall. "Precision" illustrates the ratio of true positive predictions to all instances predicted as positive, emphasizing the precision of positive predictions. "Recall" quantifies the ratio of true positives to actual positive instances, indicating the classifier's capability to identify all relevant cases. The "F1-Score" strikes a balance between precision and recall, offering a unified value that considers both false positives and false negatives. Additionally, "Cohen's Kappa" statistic evaluates the agreement between the classifier's predictions and the actual classes, while considering random chance.

Among the classifiers examined, the Support Vector Machine (SVM) stands out with notable performance. SVM achieves an accuracy of 0.793, indicating a high proportion of correct predictions. Its precision value of 0.760 showcases accurate positive predictions, while its recall value of 0.824 indicates its proficiency in capturing actual positive instances. The F1-Score, at 0.791, reflects the balance between precision and recall. Furthermore, the Cohen's Kappa score of 0.587 signifies substantial agreement beyond random chance.

In comparison, the Naive Bayes classifier achieves an accuracy of 0.653 with a precision of 0.627, recall of 0.662, F1-Score of 0.644, and Cohen's Kappa of 0.307. The Decision Tree classifier records an accuracy of 0.633, precision of 0.603, recall of 0.662, F1-Score of 0.631, and Cohen's Kappa of 0.268. The Logistic Regression classifier yields an accuracy of 0.660, precision of 0.647, recall of 0.620, F1-Score of 0.633, and Cohen's Kappa of 0.317. These metrics collectively highlight the SVM's superiority in this classification task compared to the other classifiers.



---

## Conclusions

The comparative study highlighted the strengths and weaknesses of different machine learning techniques for predicting email campaign effectiveness. Overall, the study suggests that leveraging Support Vector Machines could significantly enhance the company's targeted email marketing strategy and campaign conversion rates.

Support Vector Machines (SVMs) offer significant advantages for targeted email marketing. They excel in handling both linear and complex non-linear relationships in customer data, making them adept at capturing intricate behaviors and preferences. SVMs' ability to manage high-dimensional data, tackle imbalanced datasets, and resist the influence of outliers is particularly valuable in email marketing scenarios. Their interpretability through support vectors aids in understanding influential features, while parameter tuning controls overfitting and ensures generalization. SVMs accommodate smaller datasets common in marketing campaigns and maintain consistency in predictions. Their adaptability to changing data suits evolving marketing strategies, and diverse kernel functions offer flexibility to model different customer behaviors. While parameter tuning and computational complexity should be considered, SVMs stand as a powerful tool for enhancing targeted email marketing efforts.

The code for this study is available in the Jupyter notebook provided in the appendix section.

---

## Appendix

```
1 #INVESTIGATE THE PERFORMANCE OF A NUMBER OF MACHINE LEARNING PROCEDURES ON
   THIS DATASET USING PYTHON/SCIKIT-LEARN.
2
3 #DATA PRE-PROCESSING
4
5 ##### Importing necessary libraries
6
7 from sklearn.tree import DecisionTreeClassifier
8 from sklearn.model_selection import train_test_split
9 from sklearn.metrics import classification_report
10 import pandas as pd
11 from sklearn.preprocessing import StandardScaler
12
13 data = pd.read_csv('CE802_Ass_Resit.csv')
14 data
15
16 ###Checking the data types
17 data.dtypes
18
19
20 #####Transferring the data to two numpy arrays, one containing the variables
   (features) and another the labels (ground truth).
21
22 import numpy as np
```

```
23 x = data.loc[:, data.columns != 'C'].to_numpy() #features
24 y = data.loc[:, data.columns == 'C'].to_numpy() #labels
25
26 print(data.shape)
27 print(x.shape, y.shape)
28
29 ##### Identifying false predictors
30 import matplotlib.pyplot as plt
31
32 plt.figure(figsize=(8,6))
33 for i in range(x.shape[1]):
34     plt.subplot(3, 7, i+1)
35     plt.scatter(y,x[:,i])
36     plt.xticks([0, 1])
37     plt.title(data.columns[i])
38     # Save the plot as an image file
39     plt.savefig('false_predictors.png', dpi=300)
40
41 np.sum(np.isnan(x), 0)
42 np.sum(np.isnan(y), 0)
43
44 # From the code above we can see that none of features and the labels had
    missing values.
45
46 ### Convert the label's column vector to a 1d array
47 y = y.ravel()
48
49 ###Scaling the data
50 scaler = StandardScaler()
51 X = scaler.fit_transform(x)
52
53 ##Identifying Outliers
54 # Create the scatter plot
55 for i in range(8):
56     plt.scatter(X[:, i], y)
57
58 # Add a title and axis labels
59 plt.title('FEATURES VERSUS CLASS')
60 plt.xlabel('Feature Value')
```

```
61 plt.ylabel('Class')
62 plt.savefig('featuresversusclass.png', dpi=300)
63 plt.show()
64
65 import pandas as pd
66
67 # Convert the NumPy array y into a pandas Series
68 y_series = pd.Series(y)
69
70 # Use value_counts() to count unique values and plot a bar chart
71 y_series.value_counts().plot.bar()
72
73 # Add labels and title
74 plt.xlabel('Class')
75 plt.ylabel('Count')
76 plt.title('Distribution of Classes')
77 plt.show()
78
79 ##Split the data into training and test sets
80 X_train, X_val, y_train, y_val = train_test_split(X, y, test_size=0.3,
        random_state=42)
81
82 # A NAIVE BAYES CLASSIFIER
83
84 from sklearn.naive_bayes import GaussianNB
85 from sklearn.model_selection import GridSearchCV
86 from sklearn.metrics import accuracy_score, precision_score, recall_score,
        f1_score, cohen_kappa_score
87
88 # Create a Naive Bayes classifier
89 clf = GaussianNB()
90
91 # Naive Bayes doesn't have hyperparameters to tune with GridSearchCV
92
93 # Fit the classifier to the training data
94 clf.fit(X_train, y_train)
95
96 # Make predictions using the trained model
97 naive_bayes_predictions = clf.predict(X_val)
```

```
98
99 # Calculate accuracy
100 naive_bayes_accuracy = accuracy_score(y_val, naive_bayes_predictions)
101 print("Validation accuracy:", naive_bayes_accuracy)
102
103 # Calculate precision
104 naive_bayes_precision = precision_score(y_val, naive_bayes_predictions)
105 print("Precision:", naive_bayes_precision)
106
107 # Calculate recall
108 naive_bayes_recall = recall_score(y_val, naive_bayes_predictions)
109 print("Recall:", naive_bayes_recall)
110
111 # Calculate F1-score
112 naive_bayes_f1 = f1_score(y_val, naive_bayes_predictions)
113 print("F1-score:", naive_bayes_f1)
114
115 # Calculate the Cohen's Kappa statistic
116 naive_bayes_kappa = cohen_kappa_score(y_val, naive_bayes_predictions)
117 print("Cohen's Kappa:", naive_bayes_kappa)
118
119 #DECISION TREES
120
121 from sklearn.tree import DecisionTreeClassifier
122 from sklearn.model_selection import GridSearchCV
123
124 # Create a Decision Tree classifier
125 clf = DecisionTreeClassifier(random_state=0)
126
127 # Define the parameter grid for grid search
128 param_grid = {
129     'criterion': ['gini', 'entropy'],
130     'max_depth': [None, 5, 10, 20],
131     'min_samples_split': [2, 5, 10],
132     'min_samples_leaf': [1, 2, 4]
133 }
134
135 # Create the GridSearchCV object
136 grid_search = GridSearchCV(clf, param_grid, cv=5, scoring='accuracy')
```

```
137
138 # Fit the grid search to the training data
139 grid_search.fit(X_train, y_train)
140
141 # Get the best parameters and the corresponding model
142 best_params = grid_search.best_params_
143 best_model = grid_search.best_estimator_
144
145 # Make predictions using the best model
146 decision_tree_predictions = best_model.predict(X_val)
147 print("Best parameters:", best_params)
148
149 # Calculate accuracy
150 decision_tree_accuracy = accuracy_score(y_val, decision_tree_predictions)
151 print("Validation accuracy:", decision_tree_accuracy)
152
153 # Calculate precision
154 decision_tree_precision = precision_score(y_val, decision_tree_predictions)
155 print("Precision:", decision_tree_precision)
156
157 # Calculate recall
158 decision_tree_recall = recall_score(y_val, decision_tree_predictions)
159 print("Recall:", decision_tree_recall)
160
161 # Calculate F1-score
162 decision_tree_f1 = f1_score(y_val, decision_tree_predictions)
163 print("F1-score:", decision_tree_f1)
164
165 # Calculate the Cohen's Kappa statistic
166 decision_tree_kappa = cohen_kappa_score(y_val, decision_tree_predictions)
167 print("Cohen's Kappa:", decision_tree_kappa)
168
169 ## SUPPORT VECTOR MACHINES
170 from sklearn.svm import SVC
171 from sklearn.model_selection import GridSearchCV
172 from sklearn.metrics import accuracy_score
173
174 # Create an SVM classifier
175 clf = SVC(random_state=0)
```

```
176
177 # Define the parameter grid for grid search
178 param_grid = {
179     'C': [0.1, 1, 10],
180     'kernel': ['linear', 'poly', 'rbf', 'sigmoid'],
181     'degree': [2, 3, 4],
182     'gamma': ['scale', 'auto']
183 }
184
185 # Create the GridSearchCV object
186 grid_search = GridSearchCV(clf, param_grid, cv=5, scoring='accuracy')
187
188 # Fit the grid search to the training data
189 grid_search.fit(X_train, y_train)
190
191 # Get the best parameters and the corresponding model
192 best_params = grid_search.best_params_
193 best_model = grid_search.best_estimator_
194
195 # Make predictions using the best model
196 svm_predictions = best_model.predict(X_val)
197 print("Best parameters:", best_params)
198
199 # Calculate accuracy
200 svm_accuracy = accuracy_score(y_val, svm_predictions)
201 print("Validation accuracy:", svm_accuracy)
202
203 # Calculate precision
204 svm_precision = precision_score(y_val, svm_predictions)
205 print("Precision:", svm_precision)
206
207 # Calculate recall
208 svm_recall = recall_score(y_val, svm_predictions)
209 print("Recall:", svm_recall)
210
211 # Calculate F1-score
212 svm_f1 = f1_score(y_val, svm_predictions)
213 print("F1-score:", svm_f1)
214
```

```
215 # Calculate the Cohen's Kappa statistic
216 svm_kappa = cohen_kappa_score(y_val, svm_predictions)
217 print("Cohen's Kappa:", svm_kappa)
218
219 #LOGISTIC REGRESSIONS
220
221 from sklearn.linear_model import LogisticRegression
222 from sklearn.model_selection import GridSearchCV
223 from sklearn.metrics import accuracy_score
224
225 # Create a Logistic Regression classifier
226 clf = LogisticRegression(max_iter=10000)
227
228 # Define the parameter grid for grid search
229 param_grid = {
230     'penalty': ['l1', 'l2'],
231     'C': [0.1, 1, 10],
232     'solver': ['newton-cg', 'lbfgs', 'liblinear', 'sag', 'saga']
233 }
234
235 # Create the GridSearchCV object
236 grid_search = GridSearchCV(clf, param_grid, cv=5, scoring='accuracy')
237
238 # Fit the grid search to the training data
239 grid_search.fit(X_train, y_train)
240
241 # Get the best parameters and the corresponding model
242 best_params = grid_search.best_params_
243 best_model = grid_search.best_estimator_
244
245 # Make predictions using the best model
246 logistic_regression_predictions = best_model.predict(X_val)
247 print("Best parameters:", best_params)
248
249 # Calculate accuracy
250 logistic_regression_accuracy = accuracy_score(y_val,
        logistic_regression_predictions)
251 print("Validation accuracy:", logistic_regression_accuracy)
252
```



```
253 # Calculate precision
254 logistic_regression_precision = precision_score(y_val,
        logistic_regression_predictions)
255 print("Precision:", logistic_regression_precision)
256
257 # Calculate recall
258 logistic_regression_recall = recall_score(y_val,
        logistic_regression_predictions)
259 print("Recall:", logistic_regression_recall)
260
261 # Calculate F1-score
262 logistic_regression_f1 = f1_score(y_val, logistic_regression_predictions)
263 print("F1-score:", logistic_regression_f1)
264
265 # Calculate the Cohen's Kappa statistic
266 logistic_regression_kappa = cohen_kappa_score(y_val,
        logistic_regression_predictions)
267 print("Cohen's Kappa:", logistic_regression_kappa)
268
269 import pandas as pd
270
271 # Create a dictionary to store the accuracy measures
272 accuracy_measures = {
273     "Classifier": ["Naive Bayes", "Decision Tree", "SVM", "Logistic
        Regression"],
274     "Accuracy": [naive_bayes_accuracy, decision_tree_accuracy, svm_accuracy
        , logistic_regression_accuracy],
275     "Precision": [naive_bayes_precision, decision_tree_precision,
        svm_precision, logistic_regression_precision],
276     "Recall": [naive_bayes_recall, decision_tree_recall, svm_recall,
        logistic_regression_recall],
277     "F1-Score": [naive_bayes_f1, decision_tree_f1, svm_f1,
        logistic_regression_f1],
278     "Cohen's Kappa": [naive_bayes_kappa, decision_tree_kappa, svm_kappa,
        logistic_regression_kappa]
279 }
280
281 # Create a DataFrame from the dictionary
282 accuracy_df = pd.DataFrame(accuracy_measures)
```

```
283
284 # Print the DataFrame
285 print(accuracy_df)
286
287 import pandas as pd
288 import seaborn as sns
289 import matplotlib.pyplot as plt
290
291 # Assuming you have a DataFrame named accuracy_df
292
293 # Set the style of the plot
294 sns.set(style="whitegrid")
295
296 # Melt the DataFrame to create a suitable format for the plot
297 melted_df = accuracy_df.melt(id_vars='Classifier', var_name='Metric',
    value_name='Score')
298
299 # Create the grouped bar plot using seaborn
300 plt.figure(figsize=(10, 6))
301 sns.barplot(x='Classifier', y='Score', hue='Metric', data=melted_df,
    palette='viridis')
302
303 # Add labels and title
304 plt.xlabel('Classifier', fontweight='bold')
305 plt.ylabel('Score')
306 plt.title('Comparison of Accuracy Measures for Different Classifiers')
307
308 # Rotate x-axis labels for better visibility
309 plt.xticks(rotation=45)
310
311 # Show the legend
312 plt.legend()
313
314 # Show the plot
315 plt.tight_layout()
316 plt.show()
317
318 ##support vector machines performed the best
319
```

```
320 import numpy as np
321 from sklearn.metrics import confusion_matrix
322 from sklearn.model_selection import train_test_split
323 from sklearn.svm import SVC
324
325 # Split the data into training and validation sets
326 X_train, X_val, y_train, y_val = train_test_split(X, y, test_size=0.2,
327     random_state=0)
328
329 # Create and train an SVM classifier with the best parameters
330 clf = SVC(C=10, kernel='rbf', degree=2, gamma='scale', probability=True)
331 clf.fit(X_train, y_train)
332
333 # Make predictions on the validation set
334 predictions = clf.predict(X_val)
335
336 import seaborn as sns
337 import matplotlib.pyplot as plt
338
339 # Calculate the confusion matrix
340 cm = confusion_matrix(y_val, predictions)
341
342 # Visualize the confusion matrix using a heatmap
343 plt.figure(figsize=(8, 6))
344 sns.heatmap(cm, annot=True, fmt="d", cmap="Blues", cbar=False)
345 plt.title("Confusion Matrix")
346 plt.xlabel("Predicted Labels")
347 plt.ylabel("True Labels")
348 plt.show()
```

Listing 5.1: Codes