
Applying Machine Learning Algorithms



Submitted to:

Mr. Nauman Shafi

Submitted by:

Faisal Ilyas

2022-CS-63

Department of Computer Science
University of Engineering and Technology
Lahore, Pakistan

0.1 Code

```
1 # Importing necessary libraries
2 import pandas as pd
3 from sklearn.datasets import load_iris, fetch_california_housing
4 from sklearn.preprocessing import StandardScaler
5 from sklearn.model_selection import train_test_split
6 from sklearn.neighbors import KNeighborsClassifier
7 from sklearn.svm import SVC
8 from sklearn.ensemble import RandomForestClassifier
9 from sklearn.linear_model import LinearRegression
10 from sklearn.tree import DecisionTreeRegressor
11 from sklearn.metrics import accuracy_score, mean_squared_error, r2_score
12
13 # Task 1: Data Preprocessing
14
15 # Question 1: Load and Explore the Datasets
16 # Loading the Iris dataset for classification task
17 iris = load_iris()
18 iris_df = pd.DataFrame(data=iris.data, columns=iris.feature_names)
19 iris_df['Species'] = iris.target
20
21 # Loading the California Housing dataset for regression task
22 california = fetch_california_housing()
23 california_df = pd.DataFrame(data=california.data, columns=california.feature_names)
24 california_df['MedianHouseValue'] = california.target
25
26 # Exploring the datasets
27 # Displaying first few rows of Iris dataset
28 print("Iris Dataset:")
29 print(iris_df.head())
30
31 # Checking for missing values in Iris dataset
32 print("\nMissing values in Iris dataset:")
33 print(iris_df.isnull().sum())
34
35 # Displaying first few rows of California Housing dataset
36 print("\nCalifornia Housing Dataset:")
37 print(california_df.head())
38
39 # Checking for missing values in California Housing dataset
40 print("\nMissing values in California Housing dataset:")
41 print(california_df.isnull().sum())
```

Figure 1: Code

```
43 # Question 2: Feature Scaling and Splitting Data
44 # Feature scaling for both datasets using StandardScaler
45 scaler_iris = StandardScaler()
46 iris_scaled = scaler_iris.fit_transform(iris_df.drop(columns=['Species']))
47
48 scaler_california = StandardScaler()
49 california_scaled = scaler_california.fit_transform(california_df.drop(columns=['MedianHouseValue']))
50
51 # Splitting both datasets into training and testing sets (80% training, 20% testing)
52 X_train_iris, X_test_iris, y_train_iris, y_test_iris = train_test_split(
53     iris_scaled, iris_df['Species'], test_size=0.2, random_state=42
54 )
55
56 X_train_california, X_test_california, y_train_california, y_test_california = train_test_split(
57     california_scaled, california_df['MedianHouseValue'], test_size=0.2, random_state=42
58 )
59
60 # Task 2: Classification Algorithms
61
62 # Question 3: Apply k-Nearest Neighbors (k-NN) for Classification
63 # Training k-NN model for different values of k and evaluating accuracy
64 for k in range(1, 4):
65     knn = KNeighborsClassifier(n_neighbors=k)
66     knn.fit(X_train_iris, y_train_iris)
67     y_pred_iris = knn.predict(X_test_iris)
68     accuracy = accuracy_score(y_test_iris, y_pred_iris)
69     print(f'k = {k}, k-NN Accuracy = {accuracy}')
70
71 # Question 4: Apply Support Vector Machine (SVM) for Classification
72 # Training SVM model with different kernels and evaluating accuracy
73 kernels = ['linear', 'poly', 'rbf', 'sigmoid']
74 for kernel in kernels:
75     svm = SVC(kernel=kernel)
76     svm.fit(X_train_iris, y_train_iris)
77     y_pred_iris = svm.predict(X_test_iris)
78     accuracy = accuracy_score(y_test_iris, y_pred_iris)
79     print(f'SVM (kernel={kernel}) Accuracy = {accuracy}')
```

Figure 2: Code

```

81 # Question 5: Apply Random Forest Classifier
82 # Training Random Forest Classifier with 10 estimators and evaluating accuracy
83 rf = RandomForestClassifier(n_estimators=10)
84 rf.fit(X_train_iris, y_train_iris)
85 y_pred_iris = rf.predict(X_test_iris)
86 accuracy = accuracy_score(y_test_iris, y_pred_iris)
87 print(f'Random Forest (10 estimators) Accuracy = {accuracy}%')
88
89 # Task 3: Regression Algorithms
90
91 # Question 6: Apply Linear Regression for Regression
92 # Training Linear Regression model on California Housing dataset and evaluating performance
93 lr = LinearRegression()
94 lr.fit(X_train_california, y_train_california)
95 y_pred_california = lr.predict(X_test_california)
96 mse_lr = mean_squared_error(y_test_california, y_pred_california)
97 r2_lr = r2_score(y_test_california, y_pred_california)
98 print(f'Linear Regression MSE = {mse_lr}, R² = {r2_lr}')
99
100 # Question 7: Apply Decision Tree Regression
101 # Training Decision Tree Regressor and comparing with Linear Regression using MSE and R² score
102 dt = DecisionTreeRegressor()
103 dt.fit(X_train_california, y_train_california)
104 y_pred_california = dt.predict(X_test_california)
105 mse_dt = mean_squared_error(y_test_california, y_pred_california)
106 r2_dt = r2_score(y_test_california, y_pred_california)
107 print(f'Decision Tree MSE = {mse_dt}, R² = {r2_dt}')
108
109 # Task 4: Model Evaluation and Comparison
110
111 # Question 8: Evaluate Classification Models Using Classification Metrics
112 # Comparing k-NN, SVM, and Random Forest models based on accuracy
113 print(f'Final k-NN Accuracy: {accuracy_score(y_test_iris, knn.predict(X_test_iris))}%')
114 print(f'Final SVM Accuracy: {accuracy_score(y_test_iris, svm.predict(X_test_iris))}%')
115 print(f'Final Random Forest Accuracy: {accuracy_score(y_test_iris, rf.predict(X_test_iris))}%')
116
117 # Question 9: Evaluate Regression Models Using Regression Metrics
118 # Comparing Linear Regression and Decision Tree models using MSE and R² score
119 print(f'Final Linear Regression MSE: {mse_lr}, R²: {r2_lr}')
120 print(f'Final Decision Tree MSE: {mse_dt}, R²: {r2_dt}')
121

```

Figure 3: Code

```

122 # Task 5: Conclusion
123
124 # Question 10: Compare and Summarize the Findings
125 # Summary of classification models:
126 # - k-NN achieved perfect accuracy for k=2 and k=3.
127 # - SVM with the RBF kernel also reached perfect accuracy, while the linear and polynomial kernels performed slightly lower.
128 # - Random Forest classifier achieved perfect accuracy, making it one of the most reliable models.
129 #
130 # Summary of regression models:
131 # - Linear Regression resulted in a moderate prediction error with an MSE of 0.556.
132 # - Decision Tree Regressor slightly outperformed Linear Regression with a lower MSE of 0.505.
133

```

Figure 4: Code

1 Output

```
Iris Dataset:
  sepal length (cm)  sepal width (cm)  petal length (cm)  petal width (cm)  Species
0                5.1                3.5                1.4                0.2        0
1                4.9                3.0                1.4                0.2        0
2                4.7                3.2                1.3                0.2        0
3                4.6                3.1                1.5                0.2        0
4                5.0                3.6                1.4                0.2        0

Missing values in Iris dataset:
sepal length (cm)    0
sepal width (cm)     0
petal length (cm)    0
petal width (cm)     0
Species              0
dtype: int64

California Housing Dataset:
   MedInc  HouseAge  AveRooms  AveBedrms  Population  AveOccup  Latitude  Longitude  MedianHouseValue
0   8.3252    41.0   6.984127   1.023810     322.0   2.555556     37.88    -122.23         4.526
1   8.3014    21.0   6.238137   0.971880    2401.0   2.109842     37.86    -122.22         3.585
2   7.2574    52.0   8.288136   1.073446     496.0   2.802260     37.85    -122.24         3.521
3   5.6431    52.0   5.817352   1.073059     558.0   2.547945     37.85    -122.25         3.413
4   3.8462    52.0   6.281853   1.081081     565.0   2.181467     37.85    -122.25         3.422

Missing values in California Housing dataset:
MedInc          0
HouseAge        0
AveRooms        0
AveBedrms       0
Population      0
AveOccup        0
Latitude        0
Longitude       0
MedianHouseValue 0
dtype: int64
k = 1, k-NN Accuracy = 0.9666666666666667
k = 2, k-NN Accuracy = 1.0
k = 3, k-NN Accuracy = 1.0
SVM (kernel=linear) Accuracy = 0.9666666666666667
SVM (kernel=poly) Accuracy = 0.9666666666666667
SVM (kernel=rbf) Accuracy = 1.0
SVM (kernel=sigmoid) Accuracy = 0.9
Random Forest (10 estimators) Accuracy = 1.0
Linear Regression MSE = 0.5558915986952444, R² = 0.5757877060324508
Decision Tree MSE = 0.4983876455466812, R² = 0.6196701534999252
Final k-NN Accuracy: 1.0
Final SVM Accuracy: 0.9
Final Random Forest Accuracy: 1.0
Final Linear Regression MSE: 0.5558915986952444, R²: 0.5757877060324508
Final Decision Tree MSE: 0.4983876455466812, R²: 0.6196701534999252
🔗 Live Share
```

Figure 5: Code