

Assignment-5: Other Classical Supervised Machine Learning

Sajan Kumer Sarker

Id: 2111131642

Email: sajan.sarker@northsouth.edu

Rafsan Jani Chowdhury

Id: 2011424642

Email: rafsan.chowdhury@northsouth.edu

Rosely Mohammad

Id: 2014219642

Email: rosely.mohammad@northsouth.edu

November 22, 2024

1. INTRODUCTION

In this assignment report, we will discuss about the training, evaluation, and testing process of six different model: ZeroR Classifier, OneR Classifier, K-Nearest-Neighbor Classifier, Naive Bayesian Classifier, Support Vector Machine Classifier, and Support Vector Machine Regression on four different datasets: 'Car Evaluation (for Multi-class Classification)', 'DT-BrainCancer (for Binary Classification)', 'DT-Wage (for Regression)', 'DT-Credits (for Regression)'. For these models we will report Accuracy, Precision, Average Precision, Recall, Average Recall, Confusion Matrix, F-Score, Precision-Recall Curve and Mean-Squared Error. Here, we split this assignment into 4 parts: 1. Binary Classification on 'DT-BrainCancer' dataset, 2. Multi-class Classification on 'Car Evaluation' dataset, 3. Regression on 'DT-Wage' dataset, 4. Regression on 'DT-Credits' dataset. Each part is handled in a separate notebook file, with each file focused on a single dataset. After implementing Binary Classification problem, we follow same step for Multi-class classification problem as well, and for Regression problem we follow same step in both two files. To implement these models, we have follow multiples resources. [1] [2] [3] [4] [5] [6]

2. METHODS

Necessary Imports:

```
1 import pandas as pd
2 import numpy as np
3 import matplotlib.pyplot as plt
4 %matplotlib inline
5 import math
6 import seaborn as sns
7 import warnings
8 warnings.filterwarnings('ignore')
9 from sklearn.model_selection import train_test_split, GridSearchCV
10 from sklearn.preprocessing import LabelEncoder, OneHotEncoder, StandardScaler
11 from sklearn.svm import SVC, SVR
12 from sklearn.neighbors import KNeighborsClassifier
13 from sklearn.naive_bayes import GaussianNB
14 from collections import defaultdict
```

We have imported all necessary libraries and modules to be used in each notebook.

User Define Evaluation Metrics Functions:

```

1  # functions to calculate: accuracy, average precision, average recall, average f-score.
2
3  # Accuracy
4  def accuracy_score(y_actual, y_predict):
5      correct = 0
6      total_samples = len(y_actual)
7
8      for predict, actual in zip(y_predict, y_actual):
9          if predict == actual:
10             correct += 1
11
12     return (correct/ total_samples)
13
14 # Precision
15 def precision_score(y_actual, y_predict):
16     precision_scores = []
17
18     for class_label in set(y_actual):
19         true_positive = 0
20         false_positive = 0
21
22         for actual, predict in zip(y_actual, y_predict):
23             if predict == class_label:
24                 if actual == class_label:
25                     true_positive += 1
26                 else:
27                     false_positive += 1
28
29         if true_positive+false_positive ==0:
30             class_precision = 0.0
31         else:
32             class_precision = true_positive / (true_positive + false_positive)
33         precision_scores.append(class_precision)
34     return precision_scores
35
36 # average precision
37 def average_precision_score_macro(precision_score, y_actual):
38     return sum(precision_score)/len(set(y_actual))
39
40 # recall
41 def recall_score(y_actual, y_predict):
42     recall_scores = []
43
44     for class_label in set(y_actual):
45         true_positive = 0
46         false_negative = 0
47
48         for actual, predict in zip(y_actual, y_predict):
49             if actual == class_label:
50                 if predict == class_label:
51                     true_positive += 1
52                 else:
53                     false_negative += 1
54

```

```

55     if true_positive+false_negative ==0:
56         class_recall = 0.0
57     else:
58         class_recall = true_positive / (true_positive + false_negative)
59
60     recall_scores.append(class_recall)
61     return recall_scores
62
63 # average recall
64 def average_recall_score_macro(recall_score, y_actual):
65     return sum(recall_score)/len(set(y_actual))
66
67 # f-score
68 def fscore(y_actual,y_predict):
69     f_score_total = []
70
71     for class_label in set(y_actual):
72         true_positive = 0
73         false_positive = 0
74         false_negative = 0
75
76         for actual, predict in zip(y_actual, y_predict):
77             if predict == class_label:
78                 if actual == class_label:
79                     true_positive += 1
80                 else:
81                     false_positive += 1
82             elif actual == class_label:
83                 false_negative += 1
84
85         if true_positive + false_positive == 0:
86             precision = 0.0
87         else:
88             precision = true_positive/(true_positive+false_positive)
89         if true_positive + false_negative == 0:
90             precision = 0.0
91         else:
92             recall = true_positive/(true_positive+false_negative)
93         if precision + recall ==0:
94             f_score = 0.0
95         else:
96             f_score = 2*true_positive/((2*true_positive)+false_positive+false_negative)
97         f_score_total.append(f_score)
98     return f_score_total
99
100 # fscore average
101 def fscore_average_macro(fscore, y_actual):
102     return sum(fscore)/len(set(y_actual))
103
104 # display confusion matrix
105 def display_confusion_matrix(y_test, y_predict):
106     num_class = len(set(y_test))
107     class_labels = list(set(y_test))
108
109     plt.figure(figsize=(8, 6))
110
111     conf_matrix = np.zeros((num_class, num_class), dtype=int)
112     for actual, predict in zip(y_test, y_predict):

```

```

113     conf_matrix[actual, predict] += 1
114
115     plt.imshow(conf_matrix, cmap='Blues', interpolation='nearest')
116     plt.title("Confusion_Matrix", fontsize=16, fontweight='bold')
117     plt.xlabel("Predicted_Labels", fontsize=12)
118     plt.ylabel("True_Labels", fontsize=12)
119
120     for i in range(num_class):
121         for j in range(num_class):
122             plt.text(j, i, str(conf_matrix[i, j]), ha='center', va='center', color='red', fontsi
123
124     plt.xticks(ticks=range(num_class), labels=class_labels, fontsize=10)
125     plt.yticks(ticks=range(num_class), labels=class_labels, fontsize=10)
126     plt.tight_layout()
127     plt.show()

```

Here we define some functions manually to report the evaluation matrices (f-score, precision, recall, average precision, average recall, confusion matrix) without using scikit learn's evaluation matrices library function.

2.1 Binary Classification on 'DT-BrainCancer' dataset

Loading Dataset:

```

1 df = pd.read_csv('/content/drive/MyDrive/CSE445-Assignment/dataset/DT-BrainCancer.csv')
2 print(df.shape)
3 print(df.head())

```

Load the dataset using pandas library and check first few data from the dataset and its shape.

Checking Missing Values:

```

1 print('Number_of_missing_values:_)
2 print(df.isnull().sum())

```

Checking if any missing values are found in any features.

Drop Column:

```

1 df = df.drop(columns=['Unnamed: 0', 'sex'])
2 df = df.dropna()
3 print(df.shape)
4 print(df.head())

```

Drop the unwanted column and remove the missing values row.

Encoded Categorical Data:

```

1 le = LabelEncoder()
2 df['diagnosis'] = le.fit_transform(df['diagnosis'])
3 df['loc'] = le.fit_transform(df['loc'])
4 print(df.head(), df.tail())

```

Since there are few complex relationships between the categories, we encoded them using LabelEncoder. Then we print some samples from the first and some from the last of this encoded dataset.

Split Data into Training, Validation and Test Sets:

```

1 X = df.drop(columns=['status'])
2 y = df['status']
3
4 X_train, X_temp, y_train, y_temp = train_test_split(X, y, test_size=0.3, random_state=42)
5 X_val, X_test, y_val, y_test = train_test_split(X_temp, y_temp, test_size=0.5, random_state=42)
6
7 # this will used after cross-validation
8 X_train_val = np.concatenate((X_train, X_val))
9 y_train_val = np.concatenate((y_train, y_val))
10
11 print(f'Train_Data_Shape_(X,_y):_{X_train.shape,_y_train.shape}')
12 print(f'Validation_Data_Shape_(X,_y):_{X_val.shape,_y_val.shape}')
13 print(f'Test_Data_Shape_(X,_y):_{X_test.shape,_y_test.shape}')

```

Splitting 70 percent data for training, 15 percent for validation and rest 15 percent for testing. And concatenate the train & validation set to retrain the model.

Implementing ZeroR Classifier Model:

```

1 class_counts = y.value_counts()
2 majority_class = class_counts.idxmax()
3 majority_count = class_counts.max()
4
5 baseline_accuracy = majority_count/len(y) * 100    # calculate baseline accuracy
6
7 print(f"Majority_Class:_{majority_class}")
8 print(f"ZeroR_Baseline_Accuracy:_{baseline_accuracy:.2f}%")
9
10 # plot data
11 plt.figure(figsize=(6, 4))
12 plt.title(f'ZeroR_Baseline_Class_Distribution')
13 plt.bar(class_counts.index, class_counts.values, color=['lightblue', 'lightcoral'])
14 plt.xticks([0,1], ['0', '1'])
15 plt.xlabel('class')
16 plt.ylabel('count')
17 plt.show()

```

Implemented the ZeroR classifier model that calculates the ZeroR baseline accuracy, which measures the accuracy of always predicting the majority class in a classification problem. Identifies the majority class in the target variable and computes the percentage of data points belonging to this class to determine the accuracy of the baseline. This serves as a simple benchmark to compare the performance of more complex and advance models.

Implementing OneR Classifier Model:

```

1 # OneR Classifier implementation
2 class OneRClassifier:
3     def __init__(self):
4         self.rule = None
5
6     def fit(self, X, y):
7         best_rule = None
8         best_error = float('inf')
9
10        for column in X.columns:
11            freq_table = defaultdict(lambda: defaultdict(int))
12            for value, label in zip(X[column], y):
13                freq_table[value][label] += 1

```

```

14         error_rate = sum(max(freq_table[value].values()) for value in freq_table)/len(y)
15
16     if error_rate < best_error:
17         best_rule = (column, freq_table)
18         best_error = error_rate
19     self.rule = best_rule
20
21
22     def predict(self, X):
23         if self.rule is None:
24             raise Exception("Classifier_has_not_been_trained_yet!")
25         column, freq_table = self.rule
26
27         predict = []
28         for value in X[column]:
29             predict.append(max(freq_table[value], key=freq_table[value].get))
30         return predict

```

Implemented the OneR Classifier model that is a simple rule-based algorithm that selects the single best predictor from the dataset based on its ability to minimize classification error. During training, it evaluates each feature by creating a frequency table of feature values and target labels, then selects the feature with the lowest error as the best rule. This rule is used to assign the majority class for new data based on the stored frequency table.

Implementing ZeroR Classifier Model:

```

1  #oneR classifier initialize
2  oner = OneRClassifier()
3  oner.fit(X_train, y_train)
4
5  y_val_predict = oner.predict(X_val)
6  print("OneR_validation_set_Accuracy:", accuracy_score(y_val, y_val_predict))
7
8  y_test_predict = oner.predict(X_test)
9  print("OneR_test_set_Accuracy:", accuracy_score(y_test, y_test_predict))

```

Initialize the OneR Classifier and train it using train set, then evaluate it on validation and test set, and find the overall accuracy.

Implementing KNN, Naive Bayesian and SVM Base Model:

```

1  # Base model train and accuracy on validation dataset
2  knn_base = KNeighborsClassifier()
3  svc_base = SVC()
4  nb_base = GaussianNB()
5
6  knn_base.fit(X_train, y_train)
7  svc_base.fit(X_train, y_train)
8  nb_base.fit(X_train, y_train)
9
10 knn_pred_val = knn_base.predict(X_val)
11 svc_pred_val = svc_base.predict(X_val)
12 nb_pred_val = nb_base.predict(X_val)
13
14 print("KNN_Accuracy:", accuracy_score(y_val, knn_pred_val))
15 print("SVM_Accuracy:", accuracy_score(y_val, svc_pred_val))
16 print("Naive_Bayes_Accuracy:", accuracy_score(y_val, nb_pred_val))

```

Implemented the KNN, Naive Bayesian and SVM Classification Base model and train these models. Then evaluate each of these models on validation set and find the overall accuracy of these models.

Optimizing KNN Model with GridSearchCV:

```

1  # initialize the classifier models
2  knn_classifier = KNeighborsClassifier(n_neighbors=5)
3
4  knn_param_grid = {
5      'n_neighbors': [3, 5, 7, 9],
6      'weights': ['uniform', 'distance'],
7      'metric': ['euclidean', 'manhattan', 'minkowski']
8  }
9
10 knn_grid = GridSearchCV(estimator=knn_classifier, param_grid=knn_param_grid, cv=5, scoring='accuracy')
11 knn_grid.fit(X_train, y_train)
12
13 knn_best_params = knn_grid.best_params_
14 print("Best_KNN_Parameters:", knn_best_params)

```

We initialize a KNN Classifier and perform hyper-parameter tuning using GridSearchCV to identify the best configuration model. This GridSearchCV evaluates all combinations of these hyper-parameters using 5-fold cross-validation and accuracy as the scoring metric. Then model was trained on the training set, and best set of hyper-parameters is selected and printed.

Initialize New KNN model with best Hyper-parameters:

```

1  knn_final_model = KNeighborsClassifier(**knn_best_params)
2  knn_final_model.fit(X_train, y_train)
3
4  knn_pred_val = knn_final_model.predict(X_val)
5  print("KNN_Accuracy:", accuracy_score(y_val, knn_pred_val))

```

Initialize a new KNN model with the best hyper-parameters that we found from the cross-validation. And find the models accuracy on validation set.

Retrain the KNN model and Evaluate it on Testset:

```

1  # retrain knn final model using train_val set and evaluate on test set
2  knn_final_model.fit(X_train_val, y_train_val)
3
4  y_test_pred_knn = knn_final_model.predict(X_test)
5
6  print("KNN:")
7  print(f'Accuracy: {accuracy_score(y_test, y_test_pred_knn)}')
8  print(f'Precision: {precision_score(y_test, y_test_pred_knn)}')
9  print(f'Recall: {recall_score(y_test, y_test_pred_knn)}')
10 print(f'F1-Score: {f1_score(y_test, y_test_pred_knn)}')
11 print()
12
13 display_confusion_matrix(y_test, y_test_pred_knn)
14
15 from sklearn.metrics import precision_recall_curve, average_precision_score
16 y_scores_knn = knn_final_model.predict_proba(X_test)[:, 1]
17
18 precision, recall, thresholds = precision_recall_curve(y_test, y_scores_knn)
19 average_precision = average_precision_score(y_test, y_scores_knn)
20
21 plt.figure(figsize=(8, 6))
22 plt.plot(recall, precision, label=f'Precision-Recall_curve (AP={average_precision:.2f})')
23 plt.xlabel('Recall', fontsize=14)
24 plt.ylabel('Precision', fontsize=14)
25 plt.title('Precision-Recall_Curve', fontsize=16)
26 plt.legend(loc='best')

```

```

27 plt.grid()
28 plt.show()

```

Retrain the new KNN model with the concatenated train, validation set. Then we Evaluate the model on test set and find and plot the models classification accuracy, precision score, recall score, f1-score, confusion matrix, and the precision-recall curve. Here we used the scikit-learn library function to plot the precision-recall curve, because implementing it manually was too much tough.

Optimizing Naive Bayesian Model with GridSearchCV:

```

1  # initialize the classifier models
2  nb_classifier = GaussianNB()
3
4  nb_param_grid = {
5      'var_smoothing': [1e-9, 1e-8, 1e-7, 1e-6, 1e-5]
6  }
7
8  nb_grid = GridSearchCV(estimator=nb_classifier, param_grid=nb_param_grid, cv=5, scoring='accuracy')
9
10 nb_grid.fit(X_train, y_train)
11
12 nb_best_params = nb_grid.best_params_
13 print("Best_KNN_Parameters:", nb_best_params)

```

We initialize a Naive Bayesian Classifier and perform hyper-parameter tuning using GridSearchCV to identify the best configuration model. This GridSearchCV evaluates all combinations of these hyper-parameters using 5-fold cross-validation and accuracy as the scoring metric. Then model was trained on the training set, and best set of hyper-parameters is selected and printed.

Initialize New Naive Bayesian model with best Hyper-parameters:

```

1  nb_final_model = GaussianNB(**nb_best_params)
2  nb_final_model.fit(X_train, y_train)
3
4  nb_pred_val = nb_final_model.predict(X_val)
5  print("Naive_Bayesian_Accuracy:", accuracy_score(y_val, nb_pred_val))

```

Initialize a new Naive Bayesian model with the best hyper-parameters that we found from the cross-validation. And find the models accuracy on validation set.

Retrain the Naive Bayesian model and Evaluate it on Testset:

```

1  # retrain naive bayesian final model using train_val set and evaluate on test set
2  nb_final_model.fit(X_train_val, y_train_val)
3
4  y_test_pred_nb = nb_final_model.predict(X_test)
5
6  print("Naive_Bayesian:")
7  print(f'Accuracy: {accuracy_score(y_test, y_test_pred_nb)}')
8  print(f'Precision: {precision_score(y_test, y_test_pred_nb)}')
9  print(f'Recall: {recall_score(y_test, y_test_pred_nb)}')
10 print(f'F1-Score: {f1_score(y_test, y_test_pred_nb)}')
11 print()
12
13 display_confusion_matrix(y_test, y_test_pred_nb)
14
15 y_scores_nb = nb_final_model.predict_proba(X_test)[:, 1]
16
17 precision, recall, thresholds = precision_recall_curve(y_test, y_scores_nb)

```



```

18 average_precision = average_precision_score(y_test, y_scores_nb)
19
20 plt.figure(figsize=(8, 6))
21 plt.plot(recall, precision, label=f'Precision-Recall_curve_(AP={average_precision:.2f})')
22 plt.xlabel('Recall', fontsize=14)
23 plt.ylabel('Precision', fontsize=14)
24 plt.title('Precision-Recall_Curve', fontsize=16)
25 plt.legend(loc='best')
26 plt.grid()
27 plt.show()

```

Retrain the new Naive Bayesian model with the concatenated train, validation set. Then we Evaluate the model on test set and find and plot the models classification accuracy, precision score, recall score, f1-score, confusion matrix, and the precision-recall curve.

Optimizing SVM Model with GridSearchCV:

```

1 # initialize the classifier models
2 svm_classifier = SVC()
3
4 svm_param_grid = {
5     'C': [0.1, 1, 10, 100],
6     'kernel': ['linear', 'poly', 'rbf', 'sigmoid'],
7     'gamma': ['scale', 'auto']
8 }
9
10 svm_grid = GridSearchCV(estimator=svm_classifier, param_grid=svm_param_grid, cv=5, scoring='accuracy')
11 svm_grid.fit(X_train, y_train)
12
13 svm_best_params = svm_grid.best_params_
14 print("Best_SVM_Parameters:", svm_best_params)

```

Here, we do the same thing again. We initialize a SVM Classifier and perform hyper-parameter tuning using GridSearchCV to identify the best configuration model. This GridSearchCV evaluates all combinations of these hyper-parameters using 5-fold cross-validation and accuracy as the scoring metric. Then model was trained on the training set, and best set of hyper-parameters is selected and printed.

Initialize New SVM model with best Hyper-parameters:

```

1 svm_final_model = SVC(probability=True, **svm_best_params)
2 svm_final_model.fit(X_train, y_train)
3
4 svm_pred_val = svm_final_model.predict(X_val)
5 print("SVM_validation_set_Accuracy:", accuracy_score(y_val, svm_pred_val))

```

Here also we do the same thing again. Initialize a new SVM model with the best hyper-parameters that we found from the cross-validation. And find the models accuracy on validation set.

Retrain the SVM model and Evaluate it on Testset:

```

1 # retrain svm final model using train_val set and evaluate on test set
2 svm_final_model.fit(X_train_val, y_train_val)
3
4 y_test_pred_svm = svm_final_model.predict(X_test)
5
6 print("SVM_Test:")
7 print(f'Accuracy:{accuracy_score(y_test, y_test_pred_svm)}')
8 print(f'Precision:{precision_score(y_test, y_test_pred_svm)}')
9 print(f'Recall:{recall_score(y_test, y_test_pred_svm)}')

```

```

10 print(f'F1-Score:_{fscore(y_test,_y_test_pred_svm)}')
11 print()
12
13 display_confusion_matrix(y_test, y_test_pred_svm)
14
15 y_scores_svm = svm_final_model.predict_proba(X_test)[: , 1]
16
17 precision, recall, thresholds = precision_recall_curve(y_test, y_scores_svm)
18 average_precision = average_precision_score(y_test, y_scores_svm)
19
20 plt.figure(figsize=(8, 6))
21 plt.plot(recall, precision, label=f'Precision-Recall_curve_(AP=_{average_precision:.2f})')
22 plt.xlabel('Recall', fontsize=14)
23 plt.ylabel('Precision', fontsize=14)
24 plt.title('Precision-Recall_Curve', fontsize=16)
25 plt.legend(loc='best')
26 plt.grid()
27 plt.show()

```

Here again Retrain the new SVM model with the concatenated train, validation set. Then we Evaluate the model on test set and find and plot the models classification accuracy, precision score, recall score, f1-score, confusion matrix, and the precision-recall curve.

2.2 Multi-class Classification on 'Car Evaluation' dataset

Loading Dataset:

```

1 car_evaluation = pd.read_csv('/content/drive/MyDrive/CSE445-Assignment/dataset/car.data')
2
3 car_evaluation.head()

```

Load the Dataset using pandas library and check first few data from the dataset.

Assign Column name:

```

1 column_names = ['buying', 'maint', 'doors', 'persons', 'lug_boot', 'safety', 'class']
2 car_evaluation.columns = column_names
3 car_evaluation.head()

```

Assign each columns with a name, since the original dataset didn't have any headings.

Checking Missing Values:

```

1 print('Number_of_missing_values:_)
2 print(car_evaluation.isnull().sum())

```

Checking if any missing values are founds in any features.

Mapping Categorical Data:

```

1 car_evaluation['buying'] = car_evaluation['buying'].map({'vhigh':3, 'high':2, 'med':1, 'low':0})
2 car_evaluation['maint'] = car_evaluation['maint'].map({'vhigh':3, 'high':2, 'med':1, 'low':0})
3 car_evaluation['doors'] = car_evaluation['doors'].map({'2':0, '3':1, '4':2, '5more':3})
4 car_evaluation['persons'] = car_evaluation['persons'].map({'2':0, '4':1, 'more':2})
5 car_evaluation['lug_boot'] = car_evaluation['lug_boot'].map({'small':0, 'med':1, 'big':2})

```

```

6 car_evaluation['safety'] = car_evaluation['safety'].map({'low':0, 'med':1, 'high': 2})
7
8 car_evaluation['class'] = car_evaluation['class'].map({'unacc':0, 'acc':1, 'good': 2, 'vgood':
9 car_evaluation.head()

```

Since there is no complex relationship between the categories, so we just mapped them to values ranging from 0 to 4. Then each categorical columns are encoded in 0 to 4.

Split Data into Training, Validation and Test Sets:

```

1 X = car_evaluation.drop(columns='class')
2 y = car_evaluation['class']
3 X_train, X_temp, y_train, y_temp = train_test_split(X, y, test_size=0.3, random_state=42)
4 X_val, X_test, y_val, y_test = train_test_split(X_temp, y_temp, test_size=0.5, random_state=42)
5
6 # this will used after cross-validation
7 X_train_val = np.concatenate((X_train, X_val))
8 y_train_val = np.concatenate((y_train, y_val))
9
10 print(f'Train_Data_Shape_(X,y):_{X_train.shape,_y_train.shape}')
11 print(f'Validation_Data_Shape_(X,y):_{X_val.shape,_y_val.shape}')
12 print(f'Test_Data_Shape_(X,y):_{X_test.shape,_y_test.shape}')

```

Splitting 70 percent data for training, 15 percent for validation and rest 15 percent for testing. And concatenate the train & validation set to retrain the model.

Implementing ZeroR Classifier Model:

```

1 class_counts = y.value_counts()
2 majority_class = class_counts.idxmax()
3 majority_count = class_counts.max()
4
5 baseline_accuracy = majority_count/len(y) * 100    # calculate baseline accuracy
6
7 print(f"Majority_Class:_{majority_class}")
8 print(f"ZeroR_Baseline_Accuracy:_{baseline_accuracy:.2f}%")
9
10 # plot data
11 plt.figure(figsize=(6, 4))
12 plt.title(f'ZeroR_Baseline_Class_Distribution')
13 plt.bar(class_counts.index, class_counts.values, color=['lightblue', 'lightcoral'])
14 plt.xticks([0,1,2,3], ['0', '1', '2', '3'])
15 plt.xlabel('class')
16 plt.ylabel('count')
17 plt.show()

```

Implemented the ZeroR classifier model that calculates the ZeroR baseline accuracy, which measures the accuracy of always predicting the majority class in a classification problem. Identifies the majority class in the target variable and computes the percentage of data points belonging to this class to determine the accuracy of the baseline. This serves as a simple benchmark to compare the performance of more complex and advance models.

Implementing OneR Classifier Model:

```

1 # OneR Classifier implementation
2 class OneRClassifier:
3     def __init__(self):
4         self.rule = None
5

```

```

6  def fit(self, X, y):
7      best_rule = None
8      best_error = float('inf')
9
10     for column in X.columns:
11         freq_table = defaultdict(lambda: defaultdict(int))
12         for value, label in zip(X[column], y):
13             freq_table[value][label] += 1
14
15         error_rate = sum(max(freq_table[value].values()) for value in freq_table)/len(y)
16
17         if error_rate < best_error:
18             best_rule = (column, freq_table)
19             best_error = error_rate
20     self.rule = best_rule
21
22     def predict(self, X):
23         if self.rule is None:
24             raise Exception("Classifier_has_not_been_trained_yet!")
25         column, freq_table = self.rule
26
27         predict = []
28         for value in X[column]:
29             predict.append(max(freq_table[value], key=freq_table[value].get))
30     return predict

```

Implemented the OneR Classifier model that is a simple rule-based algorithm that selects the single best predictor from the dataset based on its ability to minimize classification error. During training, it evaluates each feature by creating a frequency table of feature values and target labels, then selects the feature with the lowest error as the best rule. This rule is used to assign the majority class for new data based on the stored frequency table.

Implementing ZeroR Classifier Model:

```

1  #oneR classifier initialize
2  oner = OneRClassifier()
3  oner.fit(X_train, y_train)
4
5  y_val_predict = oner.predict(X_val)
6  print("OneR_validation_set_Accuracy:", accuracy_score(y_val, y_val_predict))
7
8  y_test_predict = oner.predict(X_test)
9  print("OneR_test_set_Accuracy:", accuracy_score(y_test, y_test_predict))

```

Initialize the OneR Classifier and train it using train set, then evaluate it on validation and test set, and find the overall accuracy.

Implementing KNN, Naive Bayesian and SVM Base Model:

```

1  # Base model train and accuracy on validation dataset
2  knn_base = KNeighborsClassifier()
3  svc_base = SVC()
4  nb_base = GaussianNB()
5
6  knn_base.fit(X_train, y_train)
7  svc_base.fit(X_train, y_train)
8  nb_base.fit(X_train, y_train)
9
10 knn_pred_val = knn_base.predict(X_val)
11 svc_pred_val = svc_base.predict(X_val)

```

```

12 nb_pred_val = nb_base.predict(X_val)
13
14 print("KNN_Accuracy:", accuracy_score(y_val, knn_pred_val))
15 print("SVM_Accuracy:", accuracy_score(y_val, svc_pred_val))
16 print("Naive_Bayes_Accuracy:", accuracy_score(y_val, nb_pred_val))

```

Implemented the KNN, Naive Bayesian and SVM Classification Base model and train these models. Then evaluate each of these models on validation set and find the overall accuracy of these models.

Optimizing KNN Model with GridSearchCV:

```

1 # initialize the classifier models
2 knn_classifier = KNeighborsClassifier(n_neighbors=5)
3
4 knn_param_grid = {
5     'n_neighbors': [3, 5, 7, 9],
6     'weights': ['uniform', 'distance'],
7     'metric': ['euclidean', 'manhattan', 'minkowski']
8 }
9
10 knn_grid = GridSearchCV(estimator=knn_classifier, param_grid=knn_param_grid, cv=5, scoring='accuracy')
11 knn_grid.fit(X_train, y_train)
12
13 knn_best_params = knn_grid.best_params_
14 print("Best_KNN_Parameters:", knn_best_params)

```

We initialize a KNN Classifier and perform hyper-parameter tuning using GridSearchCV to identify the best configuration model. This GridSearchCV evaluates all combinations of these hyper-parameters using 5-fold cross-validation and accuracy as the scoring metric. Then model was trained on the training set, and best set of hyper-parameters is selected and printed.

Initialize New KNN model with best Hyper-parameters:

```

1 knn_final_model = KNeighborsClassifier(**knn_best_params)
2 knn_final_model.fit(X_train, y_train)
3
4 knn_pred_val = knn_final_model.predict(X_val)
5 print("KNN_Accuracy:", accuracy_score(y_val, knn_pred_val))

```

Initialize a new KNN model with the best hyper-parameters that we found from the cross-validation. And find the model's accuracy on validation set.

Retrain the KNN model and Evaluate it on Testset:

```

1 # retrain knn final model using train_val set and evaluate on test set
2 knn_final_model.fit(X_train_val, y_train_val)
3
4 y_test_pred_knn = knn_final_model.predict(X_test)
5
6 print("KNN:")
7 print(f'Accuracy: {accuracy_score(y_test, y_test_pred_knn)}')
8 precision = precision_score(y_test, y_test_pred_knn)
9 print(f'Average Precision: {average_precision_score_macro(precision, y_test)}')
10 recall = recall_score(y_test, y_test_pred_knn)
11 print(f'Average Recall: {average_recall_score_macro(recall, y_test)}')
12 f_score = f_score(y_test, y_test_pred_knn)
13 print(f'Average F1-Score: {f_score_average_macro(f_score, y_test)}')
14
15 display_confusion_matrix(y_test, y_test_pred_knn)

```

Retrain the new KNN model with the concatenated train, validation set. Then we evaluate the model on test set and find and plot the model's classification accuracy, precision score, recall score, f1-score, confusion matrix.

Optimizing Naive Bayesian Model with GridSearchCV:

```

1 # initialize the classifier models
2 nb_classifier = GaussianNB()
3
4 nb_param_grid = {
5     'var_smoothing': [1e-9, 1e-8, 1e-7, 1e-6, 1e-5]
6 }
7
8 nb_grid = GridSearchCV(estimator=nb_classifier, param_grid=nb_param_grid, cv=5, scoring='accuracy')
9
10 nb_grid.fit(X_train, y_train)
11
12 nb_best_params = nb_grid.best_params_
13 print("Best_KNN_Parameters:", nb_best_params)

```

We initialize a Naive Bayesian Classifier and perform hyper-parameter tuning using GridSearchCV to identify the best configuration model. This GridSearchCV evaluates all combinations of these hyper-parameters using 5-fold cross-validation and accuracy as the scoring metric. Then model was trained on the training set, and best set of hyper-parameters is selected and printed.

Initialize New Naive Bayesian model with best Hyper-parameters:

```

1 nb_final_model = GaussianNB(**nb_best_params)
2 nb_final_model.fit(X_train, y_train)
3
4 nb_pred_val = nb_final_model.predict(X_val)
5 print("Naive_Bayesian_Accuracy:", accuracy_score(y_val, nb_pred_val))

```

Initialize a new Naive Bayesian model with the best hyper-parameters that we found from the cross-validation. And find the models accuracy on validation set.

Retrain the Naive Bayesian model and Evaluate it on Testset:

```

1 # retrain naive bayesian final model using train_val set and evaluate on test set
2 nb_final_model.fit(X_train_val, y_train_val)
3
4 y_test_pred_nb = nb_final_model.predict(X_test)
5
6 print("Naive_Bayesian:")
7 print(f'Accuracy: {accuracy_score(y_test, y_test_pred_nb)}')
8 precision = precision_score(y_test, y_test_pred_nb)
9 print(f'Average_Precision: {average_precision_score_macro(precision, y_test)}')
10 recall = recall_score(y_test, y_test_pred_nb)
11 print(f'Average_Recall: {average_recall_score_macro(recall, y_test)}')
12 f_score = fscore(y_test, y_test_pred_nb)
13 print(f'Average_F1-Score: {fscore_average_macro(f_score, y_test)}')
14
15 print()
16
17 display_confusion_matrix(y_test, y_test_pred_nb)

```

Retrain the new Naive Bayesian model with the concatenated train, validation set. Then we Evaluate the model on test set and find and plot the models classification accuracy, precision score, recall score, f1-score, confusion matrix.

Optimizing SVM Model with GridSearchCV:

```

1 # initialize the classifier models
2 svm_classifier = SVC()
3

```

```

4 svm_param_grid = {
5     'C': [0.1, 1, 10, 100],
6     'kernel': ['linear', 'poly', 'rbf', 'sigmoid'],
7     'gamma': ['scale', 'auto']
8 }
9
10 svm_grid = GridSearchCV(estimator=svm_classifier, param_grid=svm_param_grid, cv=5, scoring='acc')
11 svm_grid.fit(X_train, y_train)
12
13 svm_best_params = svm_grid.best_params_
14 print("Best_SVM_Parameters:", svm_best_params)

```

Here, we do the same thing again. We initialize a SVM Classifier and perform hyper-parameter tuning using GridSearchCV to identify the best configuration model. This GridSearchCV evaluates all combinations of these hyper-parameters using 5-fold cross-validation and accuracy as the scoring metric. Then model was trained on the training set, and best set of hyper-parameters is selected and printed.

Initialize New SVM model with best Hyper-parameters:

```

1 svm_final_model = SVC(probability=True, **svm_best_params)
2 svm_final_model.fit(X_train, y_train)
3
4 svm_pred_val = svm_final_model.predict(X_val)
5 print("SVM_validation_set_Accuracy:", accuracy_score(y_val, svm_pred_val))

```

Here also we do the same thing again. Initialize a new SVM model with the best hyper-parameters that we found from the cross-validation. And find the models accuracy on validation set.

Retrain the SVM model and Evaluate it on Testset:

```

1 # retrain svm final model using train_val set and evaluate on test set
2 svm_final_model.fit(X_train_val, y_train_val)
3
4 y_test_pred_svm = svm_final_model.predict(X_test)
5
6 print("SVM:")
7 print(f'Accuracy: {accuracy_score(y_test, y_test_pred_svm)}')
8 precision = precision_score(y_test, y_test_pred_svm)
9 print(f'Average Precision: {average_precision_score_macro(precision, y_test)}')
10 recall = recall_score(y_test, y_test_pred_svm)
11 print(f'Average Recall: {average_recall_score_macro(recall, y_test)}')
12 f_score = f_score(y_test, y_test_pred_svm)
13 print(f'Aveerage_F1-Score: {f_score_average_macro(f_score, y_test)}')
14
15 print()
16
17 display_confusion_matrix(y_test, y_test_pred_svm)

```

Here again Retrain the new SVM model with the concatenated train, validation set. Then we Evaluate the model on test set and find and plot the models classification accuracy, precision score, recall score, f1-score, confusion matrix.

2.3 Regression on 'DT-Wage' dataset

Loading Dataset:

```

1 df = pd.read_csv('/content/drive/MyDrive/CSE445-Assignment/dataset/DT-Wage.csv')
2 print(df.shape)
3 print(df.head())

```

Load the Dataset using pandas library and check first few data from the dataset and its Shape.

Checking Missing Values:

```

1 print('Number_of_missing_values:_)')
2 print(df.isnull().sum())

```

Checking if any missing values are found in any features.

Encoded Categorical Data:

```

1 df = pd.get_dummies(df, columns=['maritl', 'race', 'education', 'region', 'jobclass', 'health'])
2 df.columns = df.columns.str.replace(' ', '_').str.replace(']', '_').str.replace('<', '_')
3 print(df.shape)
4 print(df.head())

```

Since there has few complex relationship between the categories, so we used `pd.get_dummies` to convert the categorical columns into one-hot encoded columns Then, we replace special characters ([,], <) in the column names with underscores. Then we print few samples from this encoded dataset.

Split Data into Training, Validation and Test Sets:

```

1 X = df.drop(columns=['wage'])
2 y = df['wage']
3
4 X_train, X_temp, y_train, y_temp = train_test_split(X, y, test_size=0.3, random_state=42)
5 X_val, X_test, y_val, y_test = train_test_split(X_temp, y_temp, test_size=0.5, random_state=42)
6
7 print(f'Train_Data_Shape_(X,y):_{X_train.shape,_y_train.shape}')
8 print(f'Validation_Data_Shape_(X,y):_{X_val.shape,_y_val.shape}')
9 print(f'Test_Data_Shape_(X,y):_{X_test.shape,_y_test.shape}')

```

Splitting 70 percent data for training, 15 percent for validation and rest 15 percent for testing.

User Define Function to Calculate MSE:

```

1 # mse calculation
2 def mse(y_actual, y_pred):
3     error = 0.0
4
5     for actual, predicted in zip(y_actual, y_pred):
6         error += (actual - predicted) ** 2
7     mse = error / len(y_actual)
8     return mse

```

This function will calculate the MSE of the model.

Initialize SVM Regressor Base Model:

```

1 # SVM base model
2 svm_base = SVR()
3 svm_base.fit(X_train, y_train)
4
5 y_val_pred = svm_base.predict(X_val)
6 print('Validation_Set_MSE:', mse(y_val, y_val_pred))

```

Initialize a SVM Regressor base model and train it using train dataset and then evaluate it on validation dataset and find the MSE.

Optimizing SVM Regressor Model with GridSearchCV:

```

1 param_grid = {
2     'C': [0.1, 1, 10, 100],
3     'epsilon': [0.01, 0.1, 0.2, 0.5],
4     'kernel': ['linear', 'rbf', 'poly']
5 }
6
7 svr = SVR()
8 svm_grid_search = GridSearchCV(estimator=svr, param_grid=param_grid, cv=5, scoring='neg_mean_s
9 svm_grid_search.fit(X_train, y_train)
10
11 best_params = svm_grid_search.best_params_
12 print("Best_Parameters:", best_params)

```

We initialize a SVM Regressor model and perform hyper-parameter tuning using GridSearchCV to identify the best configuration model. This GridSearchCV evaluates all combinations of these hyper-parameters using 5-fold cross-validation and `neg_mean_squared_error` as the scoring metric. Then the model was trained on the training set, and the best set of hyper-parameters is selected and printed.

Initialize the SVM Regressor Model with Best Hyper-parameters:

```

1 svm_final = SVR(**best_params)
2 svm_final.fit(X_train, y_train)
3
4 y_val_pred = svm_final.predict(X_val)
5 print('Validation_Set_MSE:', mse(y_val, y_val_pred))

```

Initialize a new SVM Regressor model with the best hyper-parameters that we found from the cross-validation. And calculate the MSE for the validation set.

Retrain the SVM Regressor model and Evaluate it on Testset:

```

1 # now retrain the model on the combined training and validation set with the best hyperparameter
2 X_train_val = np.concatenate((X_train, X_val))
3 y_train_val = np.concatenate((y_train, y_val))
4
5 svm_final.fit(X_train_val, y_train_val)
6
7 y_test_pred = svm_final.predict(X_test)
8 print('Test_Set_MSE:', mse(y_test, y_test_pred))

```

Here, we retrain the new SVM Regressor model with the concatenated train, validation set. Then we again calculate the MSE for the testset.

2.4 Regression on 'DT-Credit' dataset

Loading Dataset:

```

1 df = pd.read_csv('/content/drive/MyDrive/CSE445-Assignment/dataset/DT-Credit.csv')
2 print(df.shape)
3 print(df.head())

```

Load the Dataset using pandas library and check first few data from the dataset and its Shape.

Checking Missing Values:

```

1 print('Number_of_missing_values:_)
2 print(df.isnull().sum())

```

Checking if any missing values are found in any features.

Encoded Categorical Data:

```

1 df = pd.get_dummies(df, columns=['Own', 'Student', 'Married', 'Region'])
2 print(df.shape)
3 print(df.head())

```

Since there has few complex relationship between the categories, so we used `pd.get_dummies` to convert the categorical columns into one-hot encoded columns. Then we print few samples from this encoded dataset.

Split Data into Training, Validation and Test Sets:

```

1 X = df.drop(columns=['Balance'])
2 y = df['Balance']
3
4 X_train, X_temp, y_train, y_temp = train_test_split(X, y, test_size=0.3, random_state=42)
5 X_val, X_test, y_val, y_test = train_test_split(X_temp, y_temp, test_size=0.5, random_state=42)
6
7 print(f'Train_Data_Shape_(X,y):_{X_train.shape,_y_train.shape}')
8 print(f'Validation_Data_Shape_(X,y):_{X_val.shape,_y_val.shape}')
9 print(f'Test_Data_Shape_(X,y):_{X_test.shape,_y_test.shape}')

```

Splitting 70 percent data for training, 15 percent for validation and rest 15 percent for testing.

User Define Function to Calculate MSE:

```

1 # mse calculation
2 def mse(y_actual, y_pred):
3     error = 0.0
4
5     for actual, predicted in zip(y_actual, y_pred):
6         error += (actual - predicted) ** 2
7     mse = error / len(y_actual)
8     return mse

```

This function will calculate the MSE of the model.

Initialize SVM Regressor Base Model:

```

1 # SVM base model
2 svm_base = SVR()
3 svm_base.fit(X_train, y_train)
4
5 y_val_pred = svm_base.predict(X_val)
6 print('Validation_Set_MSE:', mse(y_val, y_val_pred))

```

Initialize a SVM Regressor base model and train it using train dataset and then evaluate it on validation dataset and find the MSE.

Optimizing SVM Regressor Model with GridSearchCV:

```

1 param_grid = {
2     'C': [0.1, 1, 10],
3     'epsilon': [0.01, 0.1, 0.5],
4     'kernel': ['linear', 'rbf', 'poly']
5 }
6
7 svr = SVR()

```

```

8 svm_grid_search = GridSearchCV(estimator=svr, param_grid=param_grid, cv=5, scoring='neg_mean_s
9 svm_grid_search.fit(X_train, y_train)
10
11 best_params = svm_grid_search.best_params_
12 print("Best_Parameters:", best_params)

```

We initialize a SVM Regressor model and perform hyper-parameter tuning using GridSearchCV to identify the best configuration model. This GridSearchCV evaluates all combinations of these hyper-parameters using 5-fold cross-validation and neg_mean_squared_error as the scoring metric. Then model was trained on the training set, and best set of hyper-parameters is selected and printed.

Initialize the SVM Regressor Model with Best Hyper-parameters:

```

1 svm_final = SVR(**best_params)
2 svm_final.fit(X_train, y_train)
3
4 y_val_pred = svm_final.predict(X_val)
5 print('Validation_Set_MSE:', mse(y_val, y_val_pred))

```

Initialize a new SVM Regressor model with the best hyper-parameters that we found from the cross-validation. And calculate the MSE for validation set.

Retrain the SVM Regressor model and Evaluate it on Testset:

```

1 # now retrain the model on the combined training and validation set with the best hyperparameter
2 X_train_val = np.concatenate((X_train, X_val))
3 y_train_val = np.concatenate((y_train, y_val))
4
5 svm_final.fit(X_train_val, y_train_val)
6
7 y_test_pred = svm_final.predict(X_test)
8 print('Test_Set_MSE:', mse(y_test, y_test_pred))

```

Here, we retrain the new SVM Regressor model with the concatenated train, validation set. Then we again calculate the MSE for the testset.

3. EXPERIMENT RESULTS (TEST SET)

3.1 Binary Classification on 'DT-BrainCancer' Dataset Results

Accuracy:

ZeroR	OneR	KNN	SVM	NB
59.77%	64.28%	78.57%	64.28%	92.85%

For Binary Classification dataset Naive Bayes Classifier model perform better than rest four models.

Precision:

ZeroR		OneR		KNN		SVM		NB	
Class 0	Class 1	Class 0	Class 1	Class 0	Class 1	Class 0	Class 1	Class 0	Class 1
-	-	-	-	0.80	0.75	0.64	0.00	1.00	0.83

Overall Naive Bayes has higher precision score compare to other models.

Recall:

ZeroR		OneR		KNN		SVM		NB	
Class 0	Class 1	Class 0	Class 1	Class 0	Class 1	Class 0	Class 1	Class 0	Class 1
-	-	-	-	0.88	0.60	1.00	0.00	0.88	1.00

Overall Naive Bayes has higher recall score compare to other models.

F-Score:

ZeroR		OneR		KNN		SVM		NB	
Class 0	Class 1	Class 0	Class 1	Class 0	Class 1	Class 0	Class 1	Class 0	Class 1
-	-	-	-	0.84	0.66	0.78	0.00	0.94	0.90

Overall Naive Bayes has higher f-score compare to other models.

Confusion Matrix and Precision-Recall Curve:

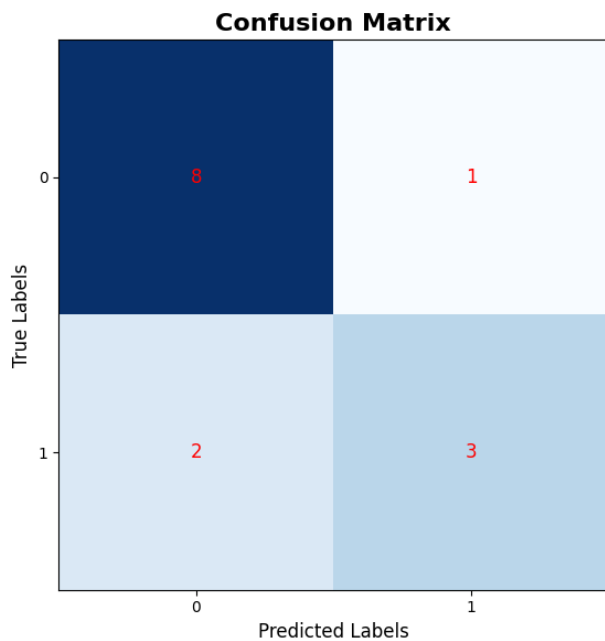


Fig. 1: KNN Models Confusion Matrix

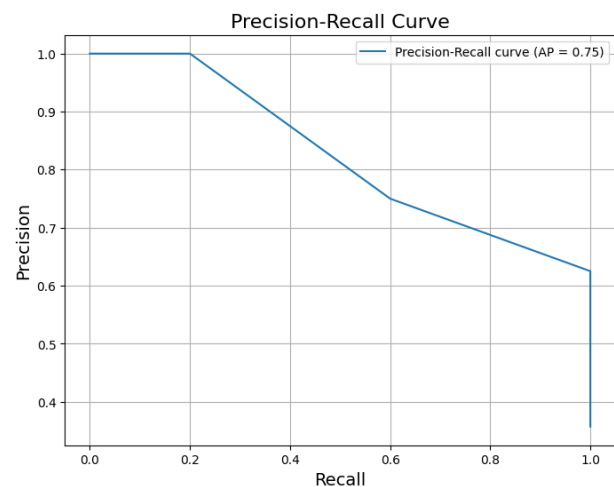


Fig. 2: KNN Models Precision-Recall Curve

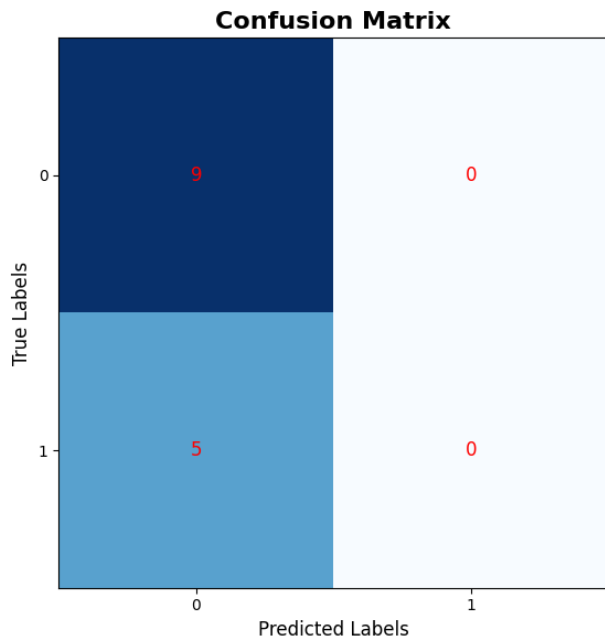


Fig. 3: SVM Models Confusion Matrix

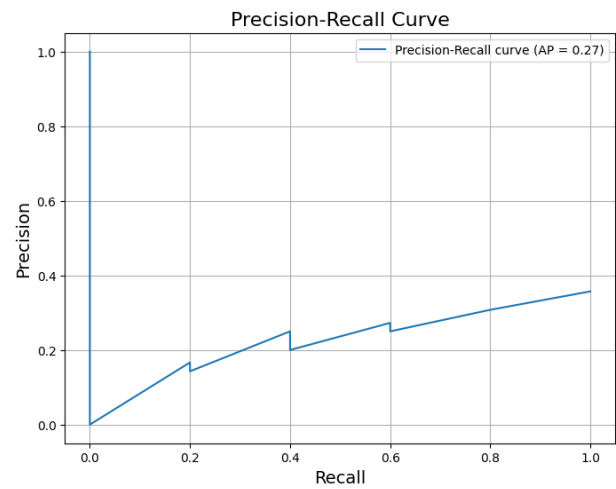


Fig. 4: SVM Models Precision-Recall Curve

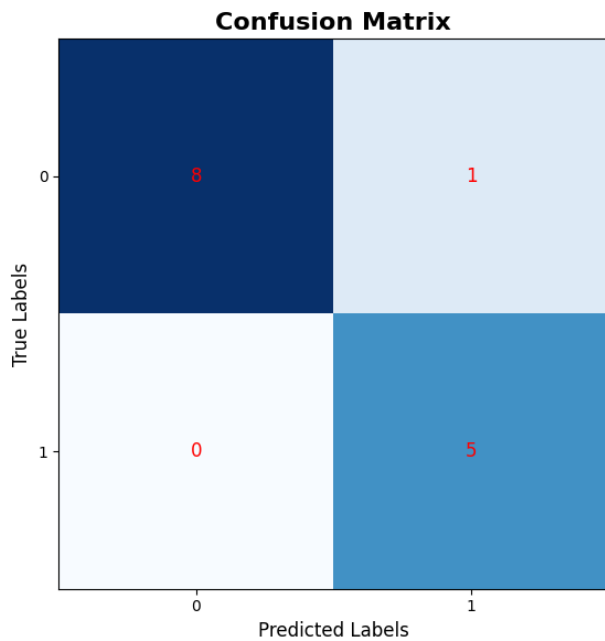


Fig. 5: Naive Bayes Models Confusion Matrix

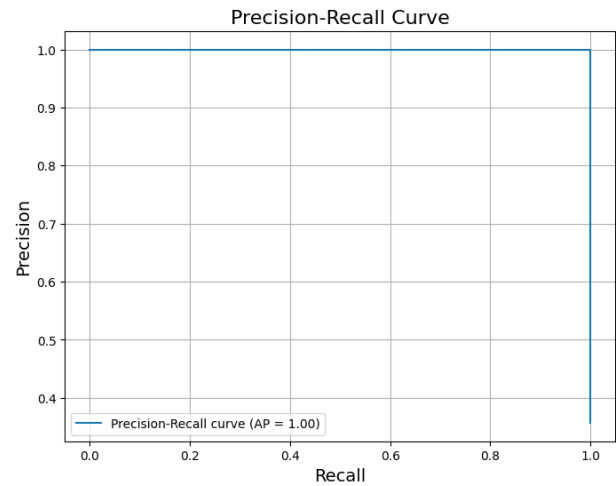


Fig. 6: Naive Bayes Precision-Recall Curve

3.2 Multi-class Classification on 'Car Evaluation' dataset

Accuracy:

ZeroR	OneR	KNN	SVM	NB
70.01%	70.38%	93.07%	99.23%	80.38%

For Multi-class Classification dataset SVM Classifier model perform better then rest four model.

Average Precision:

ZeroR	OneR	KNN	SVM	NB
-	-	0.87	0.96	0.65

For Multi-class Classification dataset SVM model has higher Precision.

Average Recall:

ZeroR	OneR	KNN	SVM	NB
-	-	0.73	0.96	0.64

For Multi-class Classification dataset SVM model has higher Average Recall.

Average F-score:

ZeroR	OneR	KNN	SVM	NB
-	-	0.77	0.96	0.58

Here, SVM Model has higher F-score compare to other models.

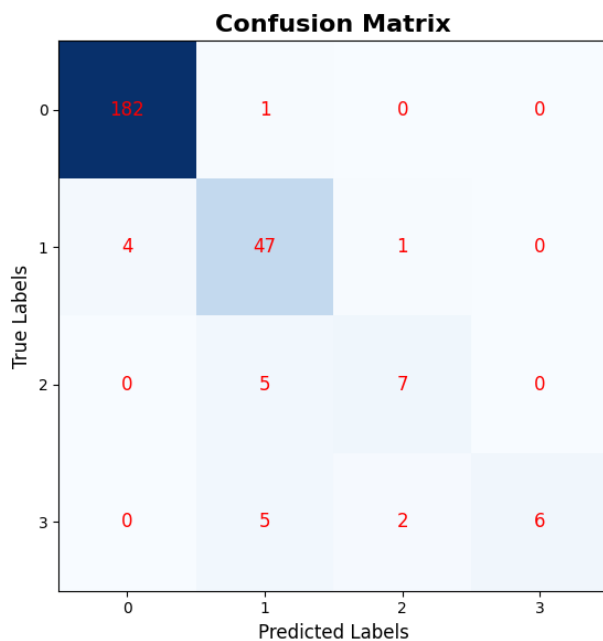
Confusion Matrix:

Fig. 7: KNN Confusion Matrix

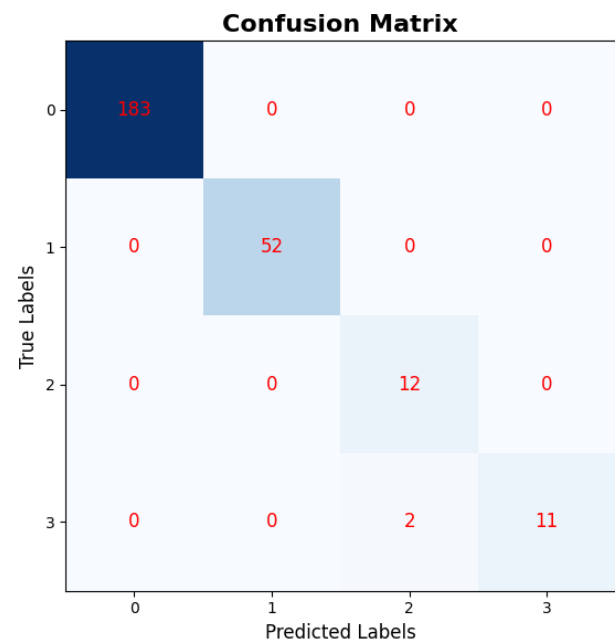


Fig. 8: SVM Confusion Matrix

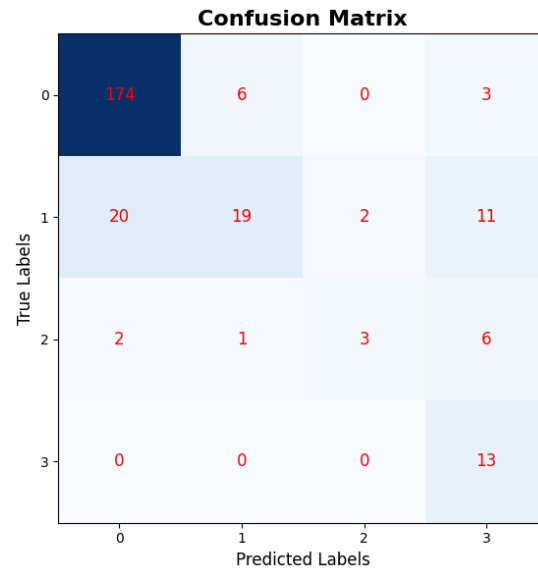


Fig. 9: Naive Bayes Confusion Matrix

3.3 Regression on 'DT-Wage' Dataset and 'DT-Credit' Dataset:

Mean Squared Error (Test Set):

DT-Wage	DT-Credit
358.61	13779.73

For the Regression Problem of two dataset, Dt-Wage datasets model has Lower MSE compare to Dt-Credit datasets model. For both model we apply same method.

4. DISCUSSION

During this assignment, we learned more about some Classical Machine Learning algorithm which are OneR classifier, ZeroR classifier, KNN, Naive Bayes, SVM, etc. During implementation we first implement the base model using test set and check the models performance on validation set. Then we perform hyper-parameter tuning on test set and then check its performance using validation set. after hyper-parameter tuning the model was performing well better then the base model. Then we merge the training set and validation set and return our models with this dataset and lastly we check the models classification performance on test set. Now the model was performing more better. Here, for Regression task we check the MSE to check our models performance. During this whole assignment we follow multiple resources, whole links are given below on the References section.

REFERENCES

- [1] S. Joshi. Gridsearchcv implementation, 2024. <https://www.kaggle.com/code/shrushtijoshi/gridsearchcv-implementation>.
- [2] mfreyeso. oner-scratch, 2024. <https://github.com/mfreyeso/oner-scratch>.
- [3] prabhat kumar singh. Oner zeror, 2020. <https://www.kaggle.com/code/prabhat12/oner-zeror>.
- [4] scikit-learn developers. Naive bayes — scikit-learn 1.3.0 documentation, 2024. https://scikit-learn.org/stable/modules/naive_bayes.html.
- [5] scikit-learn developers. Nearest neighbors — scikit-learn 1.3.0 documentation, 2024. <https://scikit-learn.org/stable/modules/neighbors.html>.
- [6] scikit-learn developers. Support vector machines — scikit-learn 1.3.0 documentation, 2024. <https://scikit-learn.org/stable/modules/svm.html>.