

# Assignment-5: ZeroR, OneR, SVM, KNN, NB

Rezwan-Ul-Alam (ID: 2011659042)  
Md. Nur Alam Jowel (ID: 2012355042)  
Raian Ruku (ID: 2013409642)

## *Email addresses:*

rezwan.alam1@northsouth.edu  
alam.jowel@northsouth.edu  
raian.ruku@northsouth.edu

## I. INTRODUCTION

In this assignment, we explore the application of various supervised learning methods to different datasets. The datasets used for classification tasks are “car.csv” and “braincancer.csv”, while for regression tasks, we utilize “wage.csv” and “credit.csv”. The supervised learning methods employed includes ZeroR classifier, OneR classifier, K-nearest-neighbors (KNN) classifier, Naïve Bayesian Classifier, Support Vector Machine (SVM), and Support Vector Regression (SVR).

## II. METHOD

### A. Data Preprocessing

Data preprocessing involves handling missing values, encoding categorical variables and scaling numerical features to ensure compatibility with these learning Algorithm.

1) **Handling missing values:** First of all we check which column has how many missing values.

```
1 [for column in df.columns:
2     total = df[column].isnull().sum()
3     print (f"Column_{column}_has_{total}_null_values")
4     print ("\n") ]
```

Here, this code iterates over each column in the DataFrame df, calculates the number of null values in each column, and prints this information along with the column name.

2) **Encoding Categorical variables:** We use both label encoder and one hot encoder for encoding categorical variables into numerical representations.

**i. Label encoder:** A label encoder is a preprocessing techniques commonly used in machine learning to convert categorical data into numerical format. When dealing with categorical data many machine learning algorithm requires numerical input. Label encoders converts these categorical values into numerical labels. In our case we use label encoding these three “car.csv, credit.csv and wage.csv” datasets.

Below, the code for label encoding of the 3 datasets is shown separately.

### For “Car.csv” data set:

```
1 [data.classes.replace(('unacc', 'acc', 'good', 'vgood'), (0, 1, 2, 3),
2 inplace = True)
3 data['classes'].value_counts()]
```

Here, These lines of code use a Label Encoder to replace categorical values ('unacc', 'acc', 'good', 'vgood') in the 'classes' column with numerical labels (0, 1, 2, 3).

Similarly, the rest of the columns having categorical values are converted to numerical values. For example in "estimated safety" column we replace categorical values ('low', 'med', 'high') with numerical values (0,1,2).

Here is the code,

```
1 [data['estimated_safety'].replace(('low', 'med', 'high'), (0, 1, 2),
2 inplace = True)
3 data['estimated_safety'].value_counts()]
```

**For "credit.csv" dataset:** In this data set these four ('own', 'Student', 'Married', 'Region') columns have categorical values. So, we converted these categorical features into numerical features. Here is label encoding code.

```
1 df.Own.replace(('Yes', 'No'), (1, 0), inplace = True)
2 df.Student.replace(('Yes', 'No'), (1, 0), inplace = True)
3 df.Married.replace(('Yes', 'No'), (1, 0), inplace = True)
4 df.Region.replace(('South', 'West', 'East'), (0, 1, 2), inplace = True)
```

This code replaces categorical values in columns 'Own', 'Student', 'Married', and 'Region' with corresponding numerical labels (1 for 'Yes' and 0 for 'No' in the first three columns, and 0 for 'South', 1 for 'West', and 2 for 'East' in the 'Region' column).

**For "wage.csv" dataset:** In this data set we use imported labelencoder class from scikit learn libraries and initialize it. Here is the code,

```
1 from sklearn.preprocessing import LabelEncoder
2 label_enc=LabelEncoder()
3 df['maritl'] = label_enc.fit_transform(df['maritl'])
4 df['race'] = label_enc.fit_transform(df['race'])
5 df['education'] = label_enc.fit_transform(df['education'])
6 df['region'] = label_enc.fit_transform(df['region'])
7 df['jobclass'] = label_enc.fit_transform(df['jobclass'])
8 df['health'] = label_enc.fit_transform(df['health'])
9 df['health_ins'] = label_enc.fit_transform(df['health_ins'])
```

It applies the LabelEncoder to encode categorical variables ('maritl', 'race', 'education', 'region', 'jobclass', 'health', and 'health\_ins') into numerical labels. Each categorical column is transformed independently, assigning a unique numerical label to each category.

### B. Splitting Dataset

As per the instruction, the first 70% of the data was used for training, and 15% of the data was used for validation and last 15% of the data is used for testing. The code is given below:

```
1 X = data.iloc[:, :6]
2 y = data.iloc[:, 6]
3
4 print("Shape_of_X:", X.shape)
5 print("Shape_of_Y:", y.shape)
6 X_train, X_temp, y_train, y_temp = train_test_split(X, y, test_size=0.3,
7 random_state=42)
8 X_val, X_test, y_val, y_test = train_test_split(X_temp, y_temp, test_size=0.5,
9 random_state=42)
```

### C. Scaling features

In our case, we use both **MinMaxScaler()** and **StandardScaler()** in our all four datasets for scaling features.

1) **MinMaxScaler()** : MinMaxScaler is used for scaling features to specified range, typically between 0 and 1. The code below introduces a min max scaling function to perform min-max scaling on numeric columns in our data frame. It calculates the scaled values, transforming X features to a common scale between 0 and 1.

Here below the code is,

```
1 from sklearn.preprocessing import MinMaxScaler
2 scaler = MinMaxScaler()
3 X_train_scaled = scaler.fit_transform(X_train)
4 X_val_scaled = scaler.transform(X_val)
5 X_test_scaled = scaler.transform(X_test)
```

2) **StandardScaler()** : The StandardScaler in scikit-learn is used for standardizing features by removing the mean and scaling them to unit variance.

Here below the code is,

```
1 from sklearn.preprocessing import StandardScaler
2 scaler = StandardScaler()
3 X_train_scaled = scaler.fit_transform(X_train)
4 X_val_scaled = scaler.transform(X_val)
5 X_test_scaled = scaler.transform(X_test)
```

## III. MODEL TRAINING

### A. Classification

For the classification tasks involving datasets "**Car\_evaluation.csv**" and "**Braincancer.csv**," the following supervised learning methods are used to trained:

NOTE regarding to the report: It's important to understand that hyperparameter tuning isn't applicable to the **ZeroR** classifier. Here's why:

- **ZeroR Classifier:** This classifier simply predicts the majority class for all instances. It doesn't involve any learning from the data and hence doesn't have any hyperparameters to tune.

Since the ZeroR classifier always predicts the majority class, its performance is solely dependent on the class distribution in the training data. There are no parameters to adjust to improve its performance on unseen data. Similar to the ZeroR classifier, OneR (One Rule) doesn't involve hyperparameter tuning in the traditional sense. Here's why:

- **OneR Classifier:** This classifier assigns data points to the class of the attribute that has the minimum within-class variance. It focuses on finding a single rule to classify instances based on the attribute with the best separation between classes.

1) **ZeroR Classifier:** This classifier predicts the most frequent class label for all instances. It serves as a baseline for comparing the performance of other classifiers. Here is the code for ZeroR classifier learning algorithm,

```

1 from sklearn.dummy import DummyClassifier
2
3 zeror_classifier = DummyClassifier(strategy="most_frequent")
4 zeror_classifier.fit(X_train_scaled, y_train)
5
6 #on train
7 y_train_pred_zeror= zeror_classifier.predict(X_train_scaled)
8 train_accuracy_zeror = calculate_accuracy(y_train, y_train_pred_zeror)
9 print("ZeroR_Classifier_Accuracy_on_Training_Set:", train_accuracy_zeror)
10
11 #on validation
12 y_val_pred_zeror= zeror_classifier.predict(X_val_scaled)
13 val_accuracy_zeror = calculate_accuracy(y_val, y_val_pred_zeror)
14 print("ZeroR_Classifier_Accuracy_on_Validation_Set:", val_accuracy_zeror)
15
16 #on test
17 y_test_pred_zeror = zeror_classifier.predict(X_test_scaled)
18 test_accuracy_zeror = calculate_accuracy(y_test, y_test_pred_zeror)
19 print("ZeroR_Classifier_Accuracy_on_Test_Set:", test_accuracy_zeror)

```

2) **OneR Classifier:** OneR builds a simple rule-based model using a single feature. It selects the single feature with the best accuracy for predicting the target variable. Here is the code for OneR classifier learning algorithm,

```

1 from mlxtend.classifier import OneRClassifier
2 oner = OneRClassifier()
3 oner.fit(X_train_scaled, y_train);
4
5 y_train_pred_oner = oner.predict(X_train_scaled)
6 train_accuracy_oner = calculate_accuracy(y_train, y_train_pred_oner)
7 print("OneR_Classifier_Accuracy_on_Training_Set:", train_accuracy_oner)
8
9 #on val
10 y_val_pred_oner = oner.predict(X_val_scaled)
11 val_accuracy_oner = calculate_accuracy(y_val, y_val_pred_oner)
12 print("OneR_Classifier_Accuracy_on_Validation_Set:", val_accuracy_oner)
13
14 #on test
15 y_test_pred_oner = oner.predict(X_test_scaled)
16 test_accuracy_oner = calculate_accuracy(y_test, y_test_pred_oner)
17 print("OneR_Classifier_Accuracy_on_Test_Set:", test_accuracy_oner)

```

3) **K-Nearest-Neighbor (KNN) Classifiers:** KNN predicts the class of an instance based on the majority class among its k-nearest neighbors. It is a non-parametric method that does not make any assumptions about the underlying data distribution. Here is the code for KNN classifier learning algorithm.

First of we want try to find out optimal hyperparameter such as ( number of neighbours and weight scheme ) that maximize the accuracy of the KNN classifier on validation data set.

Here is the code for finding optimal hyperparameter,

```

1 from sklearn.neighbors import KNeighborsClassifier
2 best_accuracy_knn = 0
3 best_neighbors = 0
4 best_weights = ''
5
6 neighbor_values = [1, 3, 5, 7, 9]
7 weight_values = ['uniform', 'distance']
8
9 for neighbors in neighbor_values:
10     for weights in weight_values:
11         knn_classifier = KNeighborsClassifier(n_neighbors=neighbors, weights=weights)
12         knn_classifier.fit(X_train_scaled, y_train)
13         y_val_pred_knn = knn_classifier.predict(X_val_scaled)
14         accuracy_knn = calculate_accuracy(y_val, y_val_pred_knn)
15         print(f"Neighbors={neighbors},_Weights={weights},_accuracy={accuracy_knn}")
16
17         if accuracy_knn > best_accuracy_knn:
18             best_accuracy_knn = accuracy_knn
19             best_neighbors = neighbors
20             best_weights = weights
21
22 print("Best_number_of_neighbors_for_KNN:", best_neighbors)
23 print("Best_weights_for_KNN:", best_weights)
24 print("Best_validation_accuracy_for_KNN:", best_accuracy_knn)

```

After finding optimal hyperparameter, we use these parameters in final training set. Here below the code,

```

1 test_accuracy_knn = calculate_accuracy(y_test, y_test_pred_knn)
2 print("KNN_Test_accuracy:", test_accuracy_knn)

```

4) **Naive Bayesian Classifier:** This classifier applies Bayes' theorem with the "naive" assumption of independence between features. It is particularly effective for text classification tasks and is known for its simplicity and efficiency. Here below the code is Naïve Bayesian Classifier learning algorithm.

First of all, we want to try to find optimal hyperparameters such as (var\_smoothing) that maximize the accuracy of the Naive Bayesian classifier on validation data set.

Here is the code for finding optimal hyperparameter,

```

1 from sklearn.naive_bayes import GaussianNB
2 best_accuracy_nb = 0
3 best_var_smoothing = 0
4 var_smoothing_values = [1e-9, 1e-8, 1e-7, 1e-6]
5 for var_smoothing in var_smoothing_values:
6     nb_classifier = GaussianNB(var_smoothing=var_smoothing)
7     nb_classifier.fit(X_train_scaled, y_train)
8     y_val_pred_nb = nb_classifier.predict(X_val_scaled)
9     accuracy_nb = calculate_accuracy(y_val, y_val_pred_nb)
10    print(f"Var_smoothing={var_smoothing},_accuracy={accuracy_nb}")
11    if accuracy_nb > best_accuracy_nb:
12        best_accuracy_nb = accuracy_nb
13        best_var_smoothing = var_smoothing
14 print("Best_variance_smoothing_for_Gaussian_Naive_Bayes:", best_var_smoothing)
15 print("Best_validation_accuracy_for_Gaussian_Naive_Bayes:", best_accuracy_nb)

```

After finding optimal hyperparameters, we use these parameters in the final training set. Here below the code,

```

1 final_nb_classifier = GaussianNB(var_smoothing=best_var_smoothing)
2 final_nb_classifier.fit(X_train_scaled, y_train)
3 y_test_pred_nb = final_nb_classifier.predict(X_test_scaled)
4 test_accuracy_nb = calculate_accuracy(y_test, y_test_pred_nb)
5 print("naive_Bayes_Test_accuracy:", test_accuracy_nb)

```

5) **Support Vector Machine (SVM)** : SVM finds the hyperplane that best separates different classes in the feature space. It works well in high-dimensional spaces and is effective for both linear and non-linear classification tasks. Here below the code is SVM classifier learning algorithm.

First of all, we want to try to find optimal hyperparameters such as (C, Kernel, degree ) that maximize the accuracy of the Support Vector Machine (SVM) classifier on validation data set. Here is the code for finding optimal hyperparameter,

```

1 from sklearn.svm import SVC
2 best_accuracy = 0
3 best_c = 0
4 best_kernel = ''
5 best_degree = 0
6
7 C_values = [0.1, 1, 10, 100]
8 kernel_values = ['linear', 'rbf', 'poly']
9 degree_values = [2, 3, 4]
10
11 for kernel in kernel_values:
12     if kernel == 'poly':
13         for degree in degree_values:
14             for C in C_values:
15                 svm_classifier = SVC(C=C, kernel=kernel, degree=degree)
16                 svm_classifier.fit(X_train_scaled, y_train)
17                 y_val_pred_svm = svm_classifier.predict(X_val_scaled)
18                 accuracy = calculate_accuracy(y_val, y_val_pred_svm)
19                 print(f"Kernel={kernel}, Degree={degree}, C={C}, accuracy={accuracy}")
20                 if accuracy > best_accuracy:
21                     best_accuracy = accuracy
22                     best_c = C
23                     best_kernel = kernel
24                     best_degree = degree
25             else:
26                 for C in C_values:
27                     svm_classifier = SVC(C=C, kernel=kernel)
28                     svm_classifier.fit(X_train_scaled, y_train)
29                     y_val_pred_svm = svm_classifier.predict(X_val_scaled)
30                     accuracy = calculate_accuracy(y_val, y_val_pred_svm)
31                     print(f"Kernel={kernel}, C={C}, accuracy={accuracy}")
32                     if accuracy > best_accuracy:
33                         best_accuracy = accuracy
34                         best_c = C
35                         best_kernel = kernel
36
37 print("Best_Kernel_for_SVM:", best_kernel)
38 print("Best_Degree_for_SVM:", best_degree) if best_kernel == 'poly' else None
39 print("Best_C_for_SVM:", best_c)
40 print("Best_validation_accuracy_for_SVM:", best_accuracy)
41 print("naive_Bayes_Test_accuracy:", test_accuracy_nb)

```

After finding optimal hyperparameters, we use these parameters in the final training set. Here below the code,

```

1 final_svm_classifier = SVC(C=best_c, kernel=best_kernel, degree=best_degree)
2 final_svm_classifier.fit(X_train_scaled, y_train)
3 y_test_pred_svm = final_svm_classifier.predict(X_test_scaled)
4 test_accuracy_svm = calculate_accuracy(y_test, y_test_pred_svm)
5 print("SVM_Test_accuracy:", test_accuracy_svm)

```

## B. Regression

For regression tasks the involving datasets 'wage.csv' and 'credit.csv', the following supervised learning algorithm 'Support Vector Regression' (SVR) is used to train.

1) **Support Vector Machine (SVM)** : SVR predicts continuous values by finding the hyperplane that has the maximum margin to the nearest data points. It is effective for handling non-linear relationships between features and the target variable. Here below the code is SVR regression learning algorithm.

First, we want to try to find optimal hyperparameters such as (C, Kernel, degree ) that maximize the Support Vector Regression (SVR) accuracy on validation data set. Here is the code for finding the optimal hyperparameter,

```

1 from sklearn.svm import SVR
2 best_mse = float('inf')
3 best_c = 0
4 best_kernel = ''
5 best_degree = 0
6 C_values = [0.1, 1, 10, 100]
7 kernel_values = ['linear', 'rbf', 'poly']
8 degree_values = [2, 3, 4]
9 for kernel in kernel_values:
10     if kernel == 'poly':
11         for degree in degree_values:
12             for C in C_values:
13                 svr_regressor = SVR(C=C, kernel=kernel, degree=degree)
14                 svr_regressor.fit(X_train_scaled, y_train)
15                 y_val_pred_svr = svr_regressor.predict(X_val_scaled)
16                 mse = calculate_mse(y_val, y_val_pred_svr)
17                 print(f"Kernel={kernel}, Degree={degree}, C={C}, MSE={mse}")
18                 if mse < best_mse:
19                     best_mse = mse
20                     best_c = C
21                     best_kernel = kernel
22                     best_degree = degree
23     else:
24         for C in C_values:
25             svr_regressor = SVR(C=C, kernel=kernel)
26             svr_regressor.fit(X_train_scaled, y_train)
27             y_val_pred_svr = svr_regressor.predict(X_val_scaled)
28             mse = calculate_mse(y_val, y_val_pred_svr)
29             print(f"Kernel={kernel}, C={C}, MSE={mse}")
30             if mse < best_mse:
31                 best_mse = mse
32                 best_c = C
33                 best_kernel = kernel
34 print("Best_Kernel_for_SVR:", best_kernel)
35 print("Best_Degree_for_SVR:", best_degree) if best_kernel == 'poly' else None
36 print("Best_C_for_SVR:", best_c)
37 print("Best_validation_MSE_for_SVR:", best_mse)

```

After finding optimal hyperparameters, we use these parameters in the final training set. Here below the code,

```

1 best_svr_model = SVR(C=best_c, kernel=best_kernel)
2 if best_kernel == 'poly':
3     best_svr_model.degree = best_degree
4 best_svr_model.fit(X_train_scaled, y_train)
5
6 # Make predictions on the test set
7 y_test_pred_svr = best_svr_model.predict(X_test_scaled)
8
9 # Calculate MSE on the test set
10 test_mse = calculate_mse(y_test, y_test_pred_svr)
11 print("Test_MSE_for_SVR:", test_mse)

```

#### IV. EVALUATION

##### A. For Classification (*car\_evaluation.csv*, *Braincancer.csv*):

For evaluation in classification tasks we calculate accuracy, precision, recall, F1-score and confusion matrix.

- **Accuracy:** accuracy is the ratio of the number of correctly classified instances (true positives and true negatives) to the total number of instances in the dataset.
- **Precision:** Precision measures the proportion of true positive predictions among all positive predictions made by the classifier. It answers the question: "Out of all the instances predicted as positive, how many are actually positive?"
- **Recall:** Recall measures the proportion of true positive predictions among all actual positive instances in the dataset. It answers the question: "Out of all the actual positive instances, how many were correctly predicted as positive?"
- **F1-score:** The F1-score is the harmonic mean of precision and recall. It provides a balance between precision and recall. F1-score reaches its best value at 1 and worst at 0.
- **Confusion matrix** A confusion matrix is a table that is often used to describe the performance of a classification model on a set of test data for which the true values are known. It presents the count of true positive (TP), true negative (TN), false positive (FP), and false negative (FN) predictions made by the classifier.

*Calculation process:: here for calculating accuracy, precision, recall, F1-score and confusion matrix we use our own created function. In below we provide all the created function with necessary equation.*

NOTE here,

**TP (True positive):** Number of correctly predicted positive instances.

**TN (True negative):** Number of correctly predicted negative instances.

**FP (False positive):** Number of incorrectly predicted positive instances.

**FN (False negative):** Number of incorrectly predicted negative instances.

1) **Precision**  $[TP/(TP+FP)]$  : Here, this function calculate\_precision(y\_true, y\_pred, label): takes the true labels, predicted labels, and a specific label for which precision is calculated. It then computes and returns the precision value for that label. If there are no predicted positives, it returns 0 to avoid division by zero.

```

1 def calculate_precision(y_true, y_pred, label):
2     true_positives = sum((y_true == label) & (y_pred == label))
3     predicted_positives = sum(y_pred == label)
4     if predicted_positives > 0:
5         precision = true_positives / predicted_positives
6     else :0
7     return precision

```



2) **Recall**  $[TP/(TP + FN)]$ : Here, The function `calculate_recall(y_true, y_pred, label)`: essentially computes the recall metric for a given label based on the true and predicted labels provided. It ensures consistence by handling the scenario where there are no actual positive instances for the specified label.

```
1 def calculate_recall(y_true, y_pred, label):
2     true_positives = sum((y_true == label) & (y_pred == label))
3     actual_positives = sum(y_true == label)
4     if actual_positives > 0:
5         recall = true_positives / actual_positives
6     else :0
7     return recall
```

3) **F1-Score**  $[2 * (precision * recall)/(precision + recall)]$ : Here, the function `calculate_f1_score(precision, recall)`: essentially computes the F1-score metric based on the provided precision and recall values. It ensures consistence by handling the scenario where both precision and recall are zero to avoid division by zero.

```
1 def calculate_f1_score(precision, recall):
2     if (precision + recall) > 0 :
3         f1_score = 2 * (precision * recall) / (precision + recall)
4     else :0
5     return f1_score
```

4) **Accuracy**  $[(TP + TN)/(TP + TN + FP + FN)]$ : Here the function `calculate_accuracy(y_true, y_pred)` essentially computes the accuracy metric based on the provided true labels and predicted labels. It counts the number of correct predictions and divides it by the total number of predictions to obtain the accuracy value.

```
1 def calculate_accuracy(y_true, y_pred):
2     . correct_predictions = sum(y_true == y_pred)
3     . total_predictions = len(y_true)
4     . accuracy = correct_predictions / total_predictions
5     . return accuracy
```

5) **Confusion matrix**  $\begin{bmatrix} TN & FN \\ FP & TP \end{bmatrix}$ : Here the function prints the confusion matrix to the console. It first prints a header indicating that the confusion matrix is for a decision tree model. Then, it iterates over each row in the confusion matrix (`conf_matrix`) and prints each row. The confusion matrix is typically a 2x2 matrix where rows represent the actual classes and columns represent the predicted classes.

```
1 conf_matrix = calculate_confusion_matrix(y_test, y_test_pred)
2     print("Decision_tree_Confusion_Matrix:")
3     for row in conf_matrix:
4         print(row)
```

#### B. For Regression (*credit.csv*, *wage.csv*):

For evaluation in regression task we calculate mean squared error (mse).

- **Mean Squared Error (MSE)**: MSE is a measure used to evaluate the average squared difference between the actual and predicted values in a regression problem. It quantifies the average squared deviation of predictions from the actual values, providing a measure of the model's predictive accuracy. Here is the calculation process of calculating MSE, **MSE**

equation,

$$\text{MSE} = 1/n \sum_{i=0}^n (y - y^i)^2$$

Here is the code for implementing MSE:

```

1 import numpy as np
2 def calculate_mse(y_true, y_pred):
3     mse = np.mean((y_true - y_pred) ** 2)
4     return mse

```

## V. RESULTS

In this section we will show the results of models trained with different learning algorithms such as for classification oneR, zeroR, KNN, Naïve Bayesian and SVM and for regression SVR.

### A. Results For Classification Tasks

Data set	Model (classifier)	Accuracy (Test set)	Precision	Recall	F1-Score	Confusion matrix
Car.csv	ZeroR	0.69230	0.173	0.250	0.205	[180, 0, 0, 0]
	OneR	0.69230	0.173	0.250	0.205	[180, 0, 0, 0]
	KNN	0.96153	0.928	0.907	0.909	[176, 4, 0, 0]
	Naive Bayesian	0.78076	0.562	0.657	0.557	[168, 10, 1, 1]
	SVM	0.98076	0.945	0.961	0.950	[179, 1, 0, 0]
Braincancer.csv	ZeroR	0.64285	0.642	1.0	0.782	[ , ]
	OneR	0.65285	0.642	1.0	0.782	[ , ]
	KNN	0.78571	0.875	0.67	0.78	[7.2, ]
	Naive Bayesian	0.92857	1.0	0.83	0.89	[8, 1]
	SVM	0.92857	1.0	0.83	0.88	[8, 1]

Here, Looking at this table, we can observe that, for

#### 1) Car.csv Dataset::

- KNN and SVM classifiers achieved the highest accuracy of 96.15% and 98.08% respectively, outperforming ZeroR, OneR, and Naive Bayesian classifiers.
- KNN and SVM also show high precision, recall, and F1-score values, that's indicates good performance across all metrics.

#### 2) Braincancer.csv Dataset:

- Naive Bayesian, and SVM classifiers achieved the highest accuracy of 92.85% for the binary classification task, outperforming ZeroR and OneR classifiers.
- Naive Bayesian and SVM classifiers exhibit perfect precision for class 0, indicating they correctly classified all instances of class 0.

### B. Results for Regression Tasks

Data set	Model	Validation MSE	Test MSE
'credit.csv'	Support Vector Regression (SVR)	10945.229594997983	12033.412463087958
'wage.csv'	Support Vector Regression (SVR)	1.26394695472098	0.8096585861480677

Here, Looking at this table, we can observe that, for

#### 1) 'credit.csv' Dataset:

- The Support Vector Regression (SVR) model achieved a validation Mean Squared Error (MSE) of approximately 10945.23 and a test MSE of approximately 12033.41.
- The test MSE is slightly higher than the validation MSE, this indicates that the model may have slightly overfit the training data.

#### 2) 'wage.csv' Dataset::

- The SVR model achieved a validation MSE of approximately 1.26 and a test MSE of approximately 0.81.
- Here, the test MSE is slightly lower than the validation MSE, this indicates that the model's performance on unseen data is slightly better than on the validation set.

### C. Graph (precision Recall curve) for Binary classification Task (BrainCancer.csv):

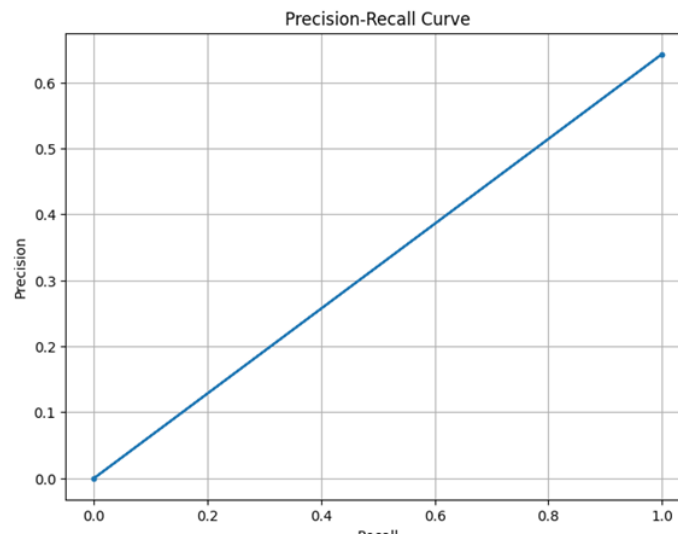


Fig. 1: ZeroR Classifier

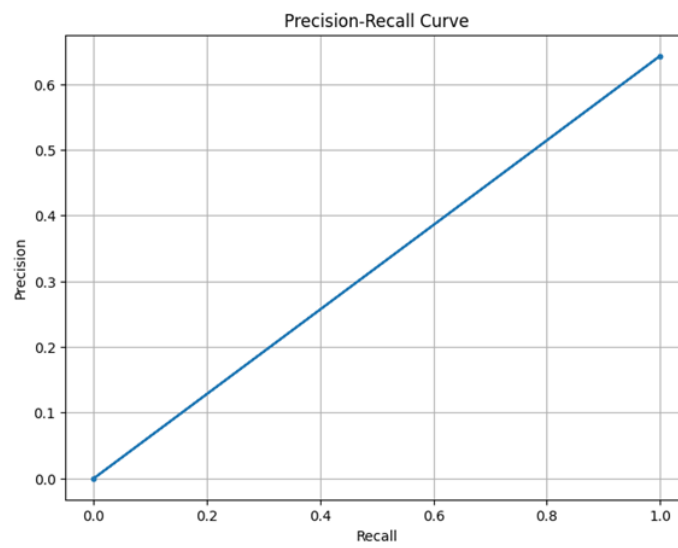


Fig. 2: OneR Classifier

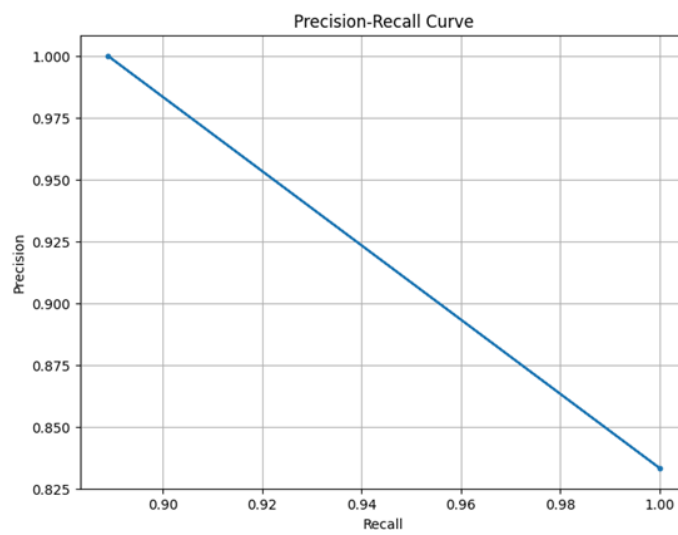


Fig. 3: SVM Classifier

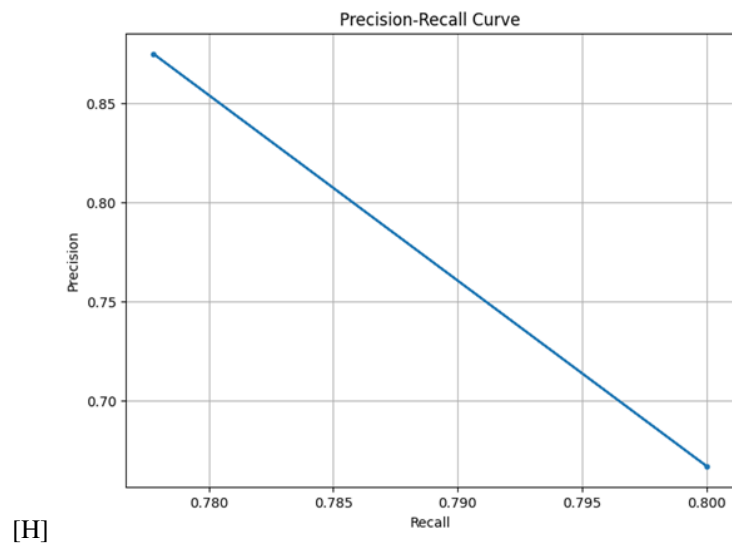


Fig. 4: KNN Classifier

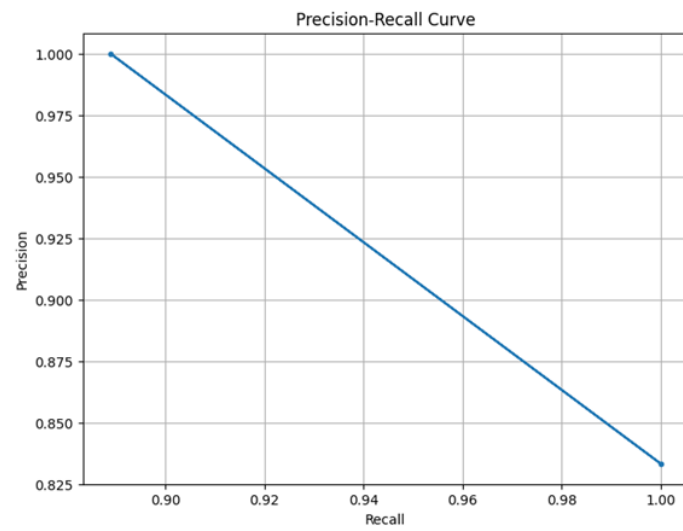


Fig. 5: Naïve Bayes Classifier

## VI. DISCUSSION

The evaluation results provide valuable insights into the performance of various supervised learning methods applied to both classification and regression tasks.

### A. *Classification Tasks:*

- For the "Car.csv" dataset, K-Nearest Neighbors (KNN) and Support Vector Machine (SVM) classifiers emerged as the top performers, achieving the highest accuracy of 96.15% and 98.08% respectively. These models also demonstrated high precision, recall, and F1-score values, indicating their robustness in correctly classifying instances across all metrics.
- Similarly, for the "Braincancer.csv" dataset, Naive Bayesian and SVM classifiers outperformed ZeroR and OneR classifiers, achieving the highest accuracy of 92.85% for the binary classification task. Notably, Naive Bayesian and SVM classifiers exhibited perfect precision for class 0, showcasing their ability to accurately classify all instances of class 0.

### B. *Regression Tasks:*

- In the "credit.csv" dataset, the Support Vector Regression (SVR) model achieved a validation Mean Squared Error (MSE) of approximately 10945.23 and a test MSE of approximately 12033.41. Although the test MSE is slightly higher than the validation MSE, indicating potential overfitting, the model still provides valuable insights into predicting the target variable.
- Conversely, in the "wage.csv" dataset, the SVR model achieved a validation MSE of approximately 1.26 and a test MSE of approximately 0.81. The lower test MSE compared to the validation MSE suggests that the model's performance on unseen data is slightly better.

Finally, these results highlight the effectiveness of KNN and SVM classifiers in handling complex classification tasks, while Naive Bayesian classifiers also demonstrate strong performance.