

A COMPARISON BETWEEN K-MEANS AND TREE ALGORITHM IN R SUPPORT VECTOR OF DATA MINING CLASSIFICATION PROBLEMS

Abdelatti A Blg^{1*}, Muner Athaba¹, Sana Abouljam¹, Abdalla Mohamed Alasoud²

¹Computer Science Department, Faculty of Science - Ajaylat, University of Zawia

²Computer Department, Faculty of Science, Al-Asmarya Islamic University

* Corresponding author: a.blg@uni.za.ly

ABSTRACT

Correct analysis of the data to create a logical relationship that summarizes the data in a new way that is understandable and useful, the data will yield huge benefits. The process of analyzing and transforming data into knowledge is called knowledge discovery in the database. In the various steps of knowledge discovery in the database, the richness and precision of the algorithms make it difficult. When analyzing big data, effective user support is essential, and it is even more important now. Metadata is a necessary component to enhance user support [1]. In this paper, we will address problems of data mining classification, four machine learning validation methods used to build and test 4 distinct datasets with the same amount of training and testing data for each predictive model. Besides, calculated the accuracy average and standard deviation of 20 trials, visualized the accuracy.

Keywords: K-MEANS, Data Mining, Machine Learning, KNIME, Tree algorithm, R.

1- INTRODUCTION

The target of the information extracting method is to extract data from a dataset and makeover it into a clear construction for additional use. this is a diagnostic method planned to scrutinize the information in seek of reliable patterns or organized associations connecting variables, and then to confirm the findings by applying the detected patterns. the focal point of this documents to concern a variety of categorization methods such as J48, kNN, CART, and SVM.

2- MATERIALS AND METHODS

Data mining is a method of exploring, identifying, and modeling large amounts of data that reveal unspecified patterns or relationships that lead to a true result.

There are several data classification algorithms available in DM. We will discuss some of the algorithms used later in this article.

2.1 KNIME Analytics

KNIME (Konstanz Information Miner) It is a modular computing environment that enables easy visual aggregation, interactive data analysis, and data manipulation. It is an open-source predictive analytics platform (published under the GNU General Public License v3) that is suitable for processing a wide variety of data formats, from simple CSV or XLSX files to more complex data structures such as XML, URL, and relational databases such as DB2, Oracle, and MySQL. Surprisingly, widespread use has not been found in the earth sciences.

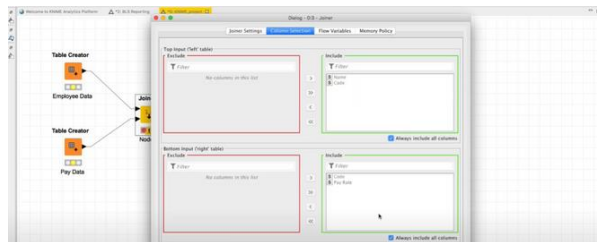


Figure 1: KNIME Analytics platform

KNIME is supported by a wide community of users and developers. Since KNIME is built on top of Eclipse, it shares the benefit of a plugin architecture that makes it easily extensible. Several custom-designed nodes

are available and are easily accessible through the community contributions area.

2.2 DecisionTree algorithm in R

A decision tree is a graph that shapes options and their outcomes in a tree format. A chart contract represents an event or a choice, margins, a graph, decision rules, or decisions. It is used mainly in machine learning and data mining with R.



Figure 2: Decision Tree algorithm in R

Moreover, R is a programming language and programming environment for statistical analysis, graphs, and reporting. R was created by Ross Ihaka and Robert Gentleman at the University of Auckland, New Zealand, and is now being developed by the Core team in development. The core of R is an interpretable computer language that allows for branching and redundancy as well as modular programming using functions. R allows integration with procedures written in C, C++, and .Net, Python, or FORTRAN to increase efficiency.

3- METHODOLOGY

With the increasing availability of information, the increase in data volume requires the use of data extraction techniques to collect useful information from the dataset. Data mining technology is described as paying special attention to classification technology as an important supervised learning

technique [3]. A Decision Tree algorithm has been used to train and test the datasets for this problem. Classification and prediction are techniques used to create categories of important data and predict potential trends. A decision tree is an important classification method for classifying data extraction. It is commonly used for marketing, surveillance, fraud detection, and scientific discovery [4]. This research will study two problems of data mining which hierarchical clustering and classification methods of data.

3.1 Problem 1: Hierarchical clustering:(single, average, complete, centroid and Ward linkages)

Hierarchical techniques often find clusters nested into groups or subdivisions. Clustering is where each data point begins with its group, and then a pair of similar groups are successively combined to form a hierarchical block. Instead, the group begins by dividing all the data points into one group, then repeatedly dividing each group into smaller groups. After splitting or merging, it will be irreversible, so the hierarchical group cannot be modified [5].

An agglomerative hierarchical clustering algorithm, the following steps are generally implemented [2]:

Step 1: Each observation is considered to be an initial cluster.

Step 2: Distances between clusters are computed.

Step 3: Two clusters that have minimum distance are combined and replaced by a single cluster. Step 4: Repeat Steps 2 and 3 until there is only a single cluster containing all observations.

The need to measure distance or proximity to determine the similarity between objects. The most common is the Euclidean distance. The output of

the hierarchical clustering algorithm is a tree diagram, which is a two-dimensional tree structure describing the arrangement of the nested groups[6].

- Single: Single linkage computes the smallest dissimilarity between two objects.
- Complete: Complete linkage, the opposite of Single linkage, computes the largest dissimilarity between two objects.
- Average: Average linkage is the intermediate between the maximum and minimum distance methods.
- Centroid: Centroid linkage is defined as the distance between centers of gravity (centroids) of two clusters.
- Ward: Ward linkage (or Ward minimum variance method) [2].

3.1.1 Task 1 KNIME Solution

The corresponding KNIME workflow and the R code have been used to calculate some statistics about the dataset "teeth.csv".

KNIME Solution:

The number of records and number of attributes has been calculated from data from this dataset as shown in figure 1.

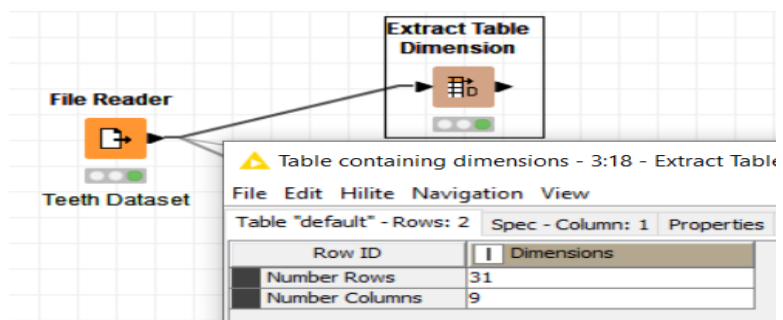


TABLE 1: THE RANGE AND MEAN VALUE OF EACH ATTRIBUTE

Animal	TopInc	BotInc	TopCan	BotCan	TopPre	BotPre	TopMol	BotMol
Length:31	Min. :0.000	Min. :1.000	Min. :0.0000	Min. :0.0000	Min. :0.000	Min. :0.000	Min. :1.000	Min. :1.000
Class :character	1st Qu.:1.000	1st Qu.:1.500	1st Qu.:0.5000	1st Qu.:0.0000	1st Qu.:2.000	1st Qu.:2.000	1st Qu.:1.000	1st Qu.:2.000
Mode :character	Median :2.000	Median :3.000	Median :1.0000	Median :1.0000	Median :3.000	Median :3.000	Median :3.000	Median :3.000
	Mean :2.097	Mean :2.419	Mean :0.7419	Mean :0.6452	Mean :2.806	Mean :2.677	Mean :2.194	Mean :2.419
	3rd Qu.:3.000	3rd Qu.:3.000	3rd Qu.:1.0000	3rd Qu.:1.0000	3rd Qu.:4.000	3rd Qu.:3.500	3rd Qu.:3.000	3rd Qu.:3.000
	Max. :3.000	Max. :4.000	Max. :1.0000	Max. :1.0000	Max. :4.000	Max. :4.000	Max. :3.000	Max. :3.000

The following figure shows a histogram of each feature of the dataset using the hist() function.

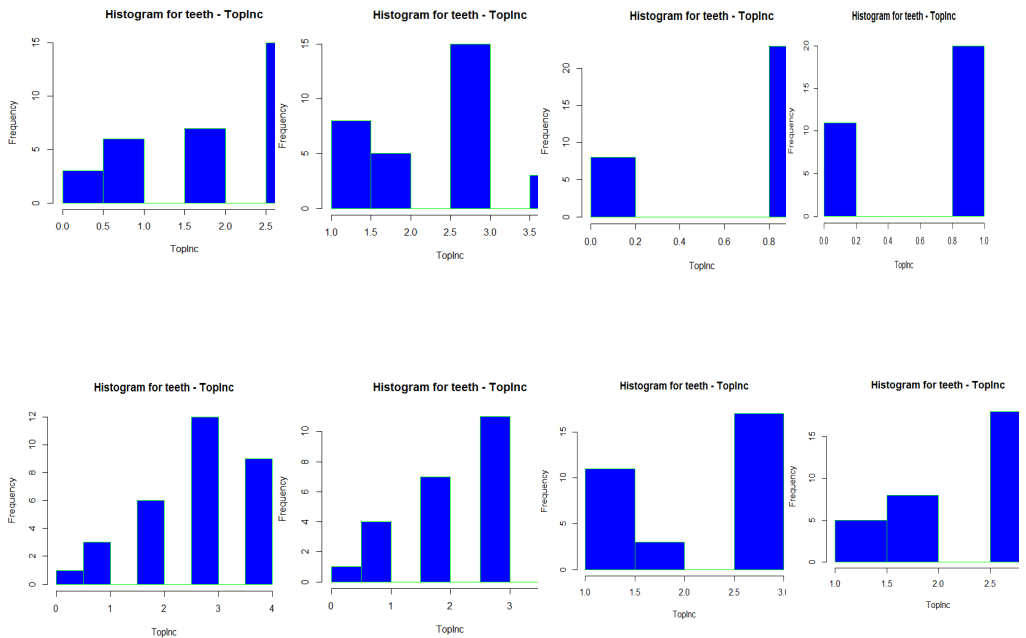


Figure 5: Histogram of each feature of the dataset

3.1.2 Task2 Clustering the teeth dataset

The generation of the dendrogram with KNIME.

Clustering the teeth dataset using hierarchical clustering.

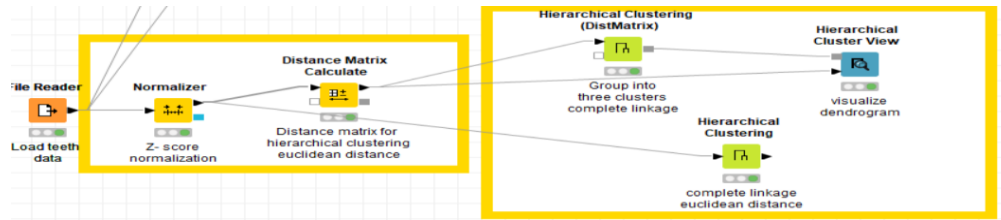


Figure 6: Clustering the teeth dataset using hierarchical clustering

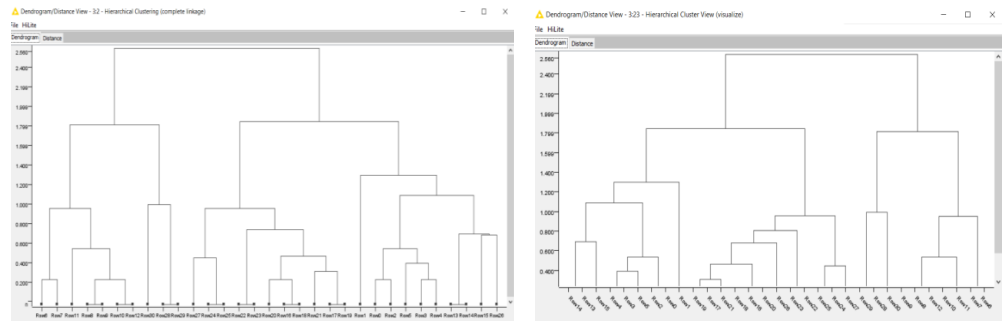


Figure7: Complete linkage Euclidean distance

3.1.3 Task dendrogram with R

The generation of the dendrogram with R.

We applied hierarchical clustering by normalizing all records and calculating the Euclidean distance between all records in the data set as shown in table2.

TABLE 2: ILLUSTRATE THE EUCLIDEAN DISTANCE MATRIX BETWEEN (DISSIMILARITY MATRIX) ALL THE ROWS (OBJECTS).

ROWS	COLS	1	2	3	4	5	6	7	8	9	10	11	12	13	14	15	16	17	18	19	20	21	22	23	24	25	26	27	28	29	30	31	32	33	34	35	36	37	38	39	40	41	42	43	44	45	46	47	48	49	50	51	52	53	54	55	56	57	58	59	60	61	62	63	64	65	66	67	68	69	70	71	72	73	74	75	76	77	78	79	80	81	82	83	84	85	86	87	88	89	90	91	92	93	94	95	96	97	98	99	100
ROWS	COLS	1	2	3	4	5	6	7	8	9	10	11	12	13	14	15	16	17	18	19	20	21	22	23	24	25	26	27	28	29	30	31	32	33	34	35	36	37	38	39	40	41	42	43	44	45	46	47	48	49	50	51	52	53	54	55	56	57	58	59	60	61	62	63	64	65	66	67	68	69	70	71	72	73	74	75	76	77	78	79	80	81	82	83	84	85	86	87	88	89	90	91	92	93	94	95	96	97	98	99	100

Hierarchical groups were applied to the dataset, and in Table 3, the methods used with generating a dendrogram by hierarchical clustering.

The data has been divided into three subgroups (k=3).

#Hierarchical clustering using complete

Linkage.

```
hc1=hclust(d, method = "complete")
```

#plot the obtained dendrogram

```
plot(hc1, cex=0.6, hang=-1,
main="Dendrogram of complete
```

Linkage")

#compute with agnes

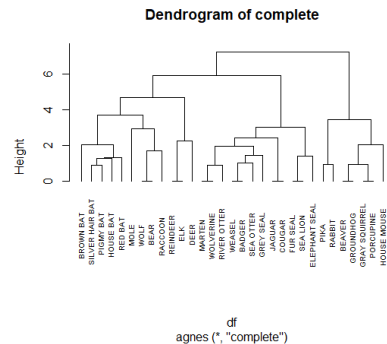
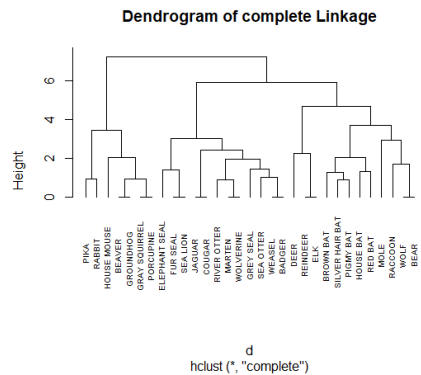
```
hc2=agnes(df,method="complete")
```

#Agglomerative coefficient- strength of clustering

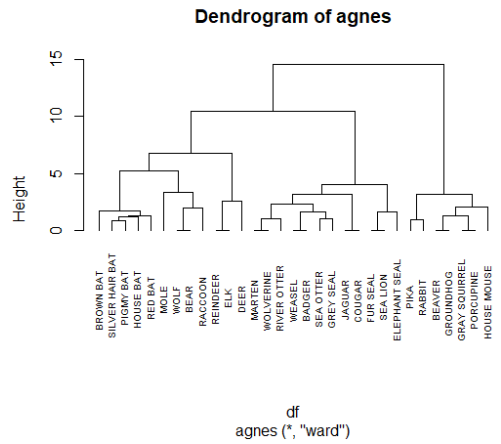
```
pltree(hc2,cex=0.6, hang=-1,
main="Dendrogram of complete")
```

hc2\$ac

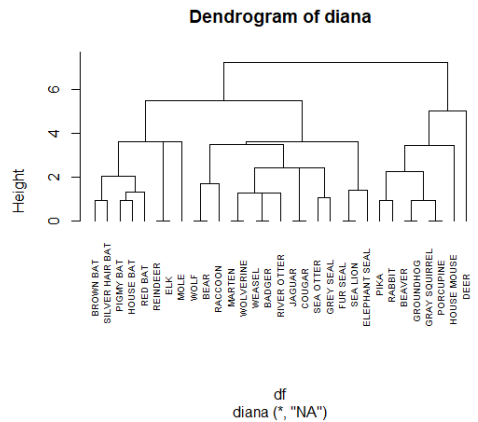
0.9014727



```
hc3=agnes(df, method="ward")
pltree(hc3,cex=0.6, hang=-1,
main="Dendrogram of agnes")
```



```
## Divisive HC
# compute divisive hierarchical clustering
hc4=diana(df)
# divide coefficient; amount of clustering
# structure found
pltree(hc4, cex=0.6, hang=-1,
main="Dendrogram of diana")
```



```
hc4$dc
0.8861321
## identify sub-groups
# ward's method
hc5<- hclust(d, method="ward.D2")
# cut tree into 3 groups
sub_grp<- cutree(hc5,k=3)# work with 3 clusters
```

Number of members in each cluster

Table(sub_grp) The number of members of each of the groups.

Group number	1	2	3
Number of members	12	7	12

library (tidyverse)

```
teeth%>%mutate(cluster=sub_grp)%>%h  
ead
```

```
plot(hc5, cex=0.6)
```

```
rect.hclust(hc5,k=5, border=2:5)
```

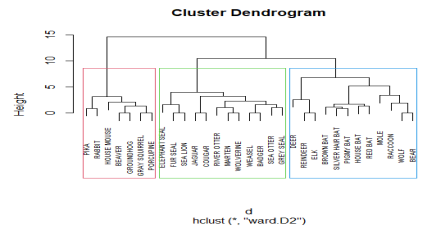
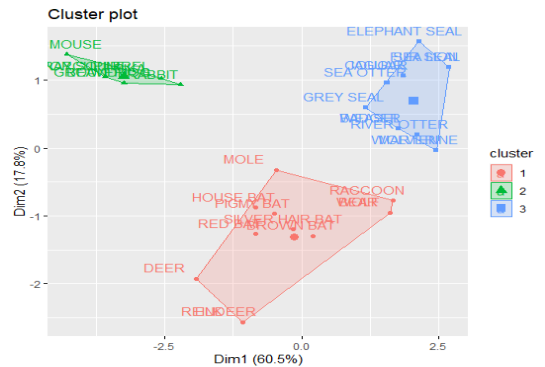


TABLE 3: DEMONSTRATE THE EVALUATION AND DETERMINE THE OPTIMAL CLUSTER NUMBER.

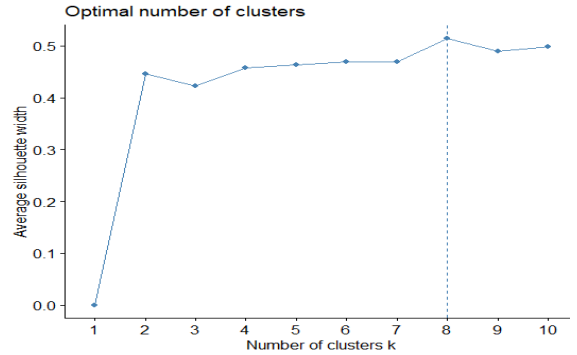
Plotted following plots for each cluster to evaluate.

```
fviz_cluster(list(data=df,cluster=sub_g  
rp))
```



Determine the optimal cluster number(silhouette)

```
fviz_nbclust(df, FUN=hcut,
method="silhouette")
```



Determine the optimal cluster number(wss)

```
fviz_nbclust(df, FUN=hcut,
method="wss")
```

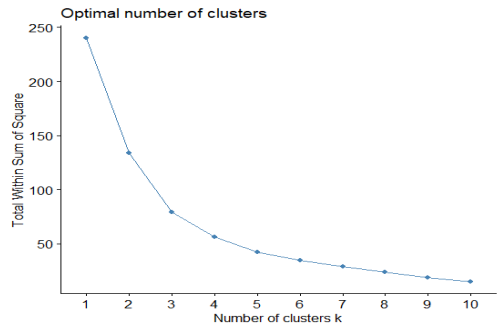


TABLE 4: EVALUATION AND DETERMINE THE OPTIMAL CLUSTER NUMBER

3.2 Problem #2 - Classification

Data mining is an interdisciplinary field, representing the integration of a range of disciplines, including database systems, statistics, machine learning, visualization, and information science.

Additionally, depending on the data extraction method used, techniques from other disciplines can be applied, such as neural networks, fuzzy and/or approximate set theory, knowledge representation, inductive logic programming, or high-performance computing.

3.2.1 Task 1 KNIME Solution

The iris dataset is used to illustrate solutions adopted by the four assessment methods. Use decision tree nodes and loops to train and test data sets for more than 20 paths for each method to calculate mean accuracy, standard deviation, mean, and standard deviation of calculation accuracy.

Resub: Resubstitution error method, Use all data for training and testing for this method.

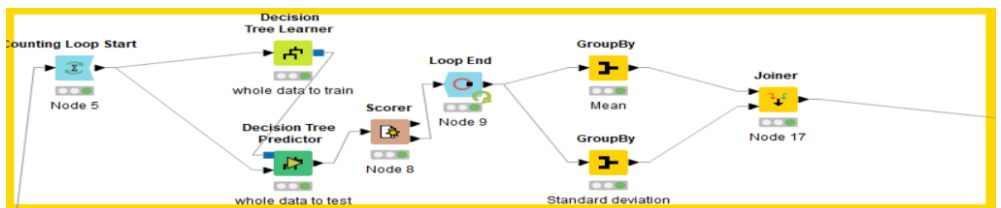


Figure 8: Resubstitution error method

hold-out-10%: hold-out method with 10% - 90% In this method, the partition split into 90% of the data for training and 10% for testing.

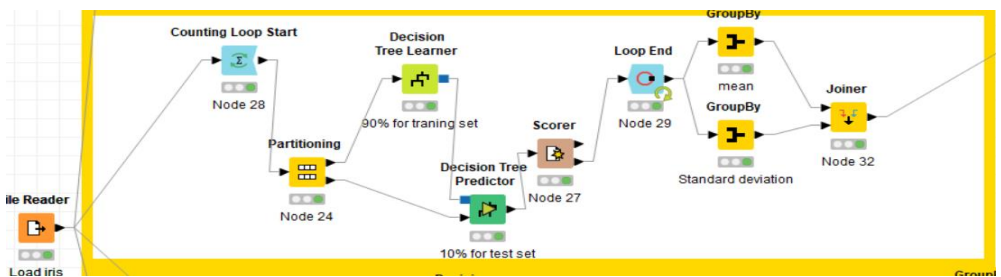


Figure 9: hold-out method

Val-10f: 10-fold cross-validation used the k-fold cross-validation method and needs to use the X-Partitioner and X-Aggregator nodes. Specifying the number of folds in the X-Partitioner node, k=10 used in

this case.

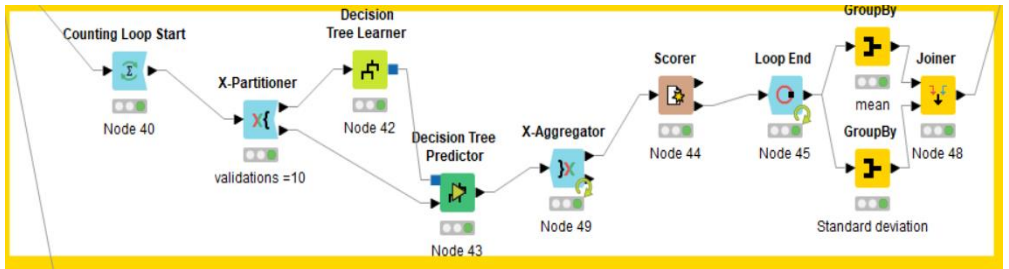


Figure10: 10-fold cross-validation

LOOCV: leave-one-out cross-validation method

LOOCV method was used in the same X-Partitioner and X-Aggregator nodes that used in the Val-10f method the distinction is selected the leave-one-out.

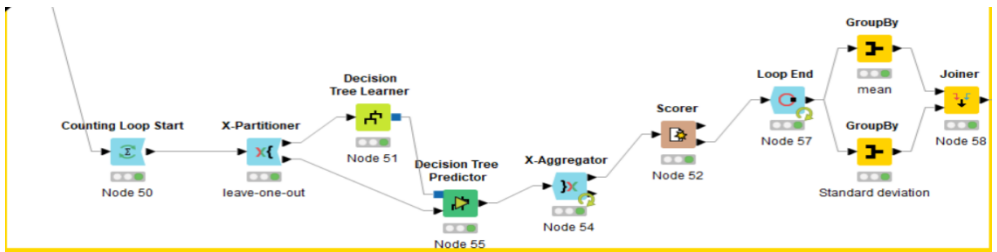
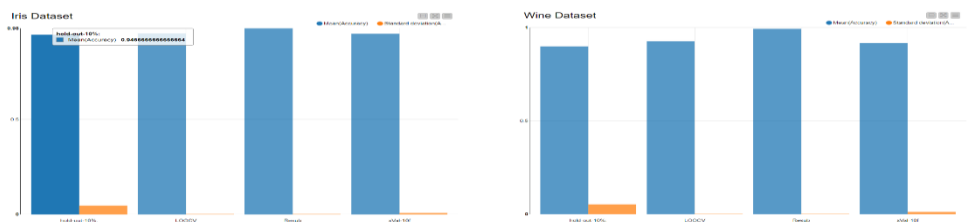


Figure 11: LOOCV method

The results are presented employing a bar chart for each dataset to compare the four methods.



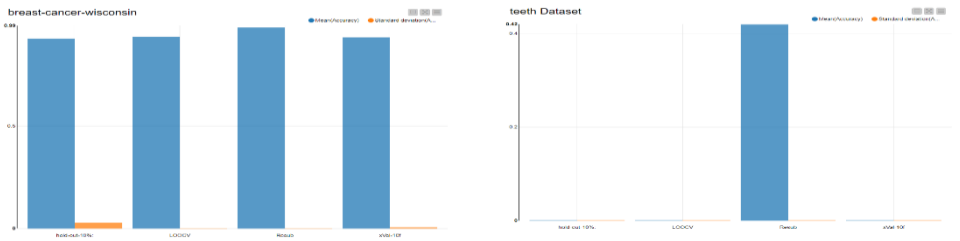


Figure 12 chart for each dataset to compare the four methods

Iris dataset:			Wine dataset:		
D	Mean(A...	D	Standa...	S	Method
0.98		0		Resub	
0.947		0.046		hold-out-10f	
0.951		0.007		xVal-10f	
0.953		0		LOOCV	
breast-cancer-Wisconsin dataset:			Teeth dataset:		
D	Mean(A...	D	Standa...	S	Method
0.981		0		Resub	
0.926		0.029		hold-out-10f	
0.932		0.006		xVal-10f	
0.936		0		LOOCV	

TABLE 5 RESULT OF DATASET TO COMPARE THE FOUR METHODS

3.2.2 Task 2 Decision of Tree algorithm in R

resub: resubstitution error method for this method.

```

holdout_acc<-c()

for (n in 1:20) {

Iris dataset    #H= holdout(1:5,ratio=0.9, mode="random",seed=NULL)

                #print (H)

                # random stratified holdout
    
```

```

H=holdout(iris$class,ratio=0.9,mode="random",seed=NULL)
train_data_iris<-(table(iris[H$str,]$class))
test_data_iris<-(table(iris[H$ts,]$class))
test_Y_iris<- (table(iris[H$ts,]$class))
train_data_iris
test_data_iris
test_Y_iris
iris_model<- train(class~ ., method="rpart",data=train_data_iris)
y_iris=predict(iris_model, newdata=test_data_iris)
y1_iris<-as.character(y_iris)
print (y1_iris)
holdout_acc$iris[n]<-mean(test_Y_iris==y1_iris)
}
for (n in 1:20){
train_data_teeth=teeth
train_test_teeth=teeth
test_Y_teeth=teeth$Animal
Teeth dataset teeth_model<- train(Animal~ ., method="rpart",data=train_data_teeth)
y_teeth=predict(teeth_model, newdata=train_test_teeth)
y1_teeth<-as.character(y_teeth)
resub_AccuracyData$teeth[n]<- mean(test_Y_teeth==y1_teeth)
}
Wine dataset
train_data_wine=wine

```

```

train_test_wine=wine

test_Y_wine=wine$Alcohol

wine_model<- train(Alcohol~ ., method="rpart",data=train_data_wine)

y_wine=predict(wine_model, newdata=train_test_wine)

y1_wine<-as.character(y_wine)

resub_AccuracyData$wine[n]<- mean(test_Y_wine==y1_wine)
}

for (n in 1:20){

train_data_breast-cancer-wisconsin=breast-cancer-wisconsin

train_test_breast-cancer-wisconsin=breast-cancer-wisconsin

test_Y_breast-cancer-wisconsin=breast-cancer-wisconsin$class

breast- breast-cancer-wisconsin_model<- train(class~ .,
cancer- method="rpart",data=train_data_breast-cancer-wisconsin)
wisconsin

y_breast-cancer-wisconsin=predict(breast-cancer-wisconsin_model,
dataset newdata=train_test_breast-cancer-wisconsin)

y1_breast-cancer-wisconsin<-as.character(y_breast-cancer-wisconsin)

resub_AccuracyData$breast-cancer-wisconsin[n]<- mean(test_Y_breast-
cancer-wisconsin==y1_breast-cancer-wisconsin)

}

```

hold-out-10%: hold-out method with 10% - 90%

```

holdout_acc<-c()

Teeth for (n in 1:20) {
dataset H=holdout(teeth$Animal, ratio=0.9, mode="random", seed=NULL)

train_data_teeth<-(table(teeth[H$str,]$Animal))

```

```

test_data_teeth<-(table(teeth[H$ts,]$Animal))
test_Y_teeth<- (table(teeth[H$ts,]$Animal))
train_data_teeth
test_data_teeth
test_Y_teeth
teeth_model<- train(Animal~ ., method="rpart",data=train_data_teeth)
y_teeth=predict(teeth_model, newdata=test_data_teeth)
y1_teeth<-as.character(y_teeth)
print (y1_teeth)
holdout_acc$teeth[n]<-mean(test_Y_teeth==y1_teeth)
}
holdout_acc<-c()
for (n in 1:20) {
  H=holdout(wine$Alcohol,ratio=0.9,mode="random",seed=NULL)
  train_data_wine<-(table(wine[H$tr,]$Alcohol))
  test_data_wine<-(table(wine[H$ts,]$Alcohol))
  test_Y_wine<- (table(wine[H$ts,]$Alcohol))
  train_data_wine
  test_data_wine
  test_Y_wine
  wine_model<- train(Alcohol~ ., method="rpart",data=train_data_wine)
  y_wine=predict(wine_model, newdata=test_data_wine)
  y1_wine<-as.character(y_wine)
  print (y1_wine)

```

Wine
dataset

```

        holdout_acc$wine[n]<-mean(test_Y_wine==y1_wine)
    }
    holdout_acc<-c()
    for (n in 1:20) {
        H=holdout(breast-cancer-
wisconsin$class,ratio=0.9,mode="random",seed=NULL)
        train_data_breast-cancer-wisconsin<-(table(breast-cancer-
wisconsin[H$str,]$class))
        test_data_breast-cancer-wisconsin<-(table(breast-cancer-
wisconsin[H$ts,]$class))
        test_Y_breast-cancer-wisconsin<-                (table(breast-cancer-
wisconsin[H$ts,]$class))
breast-
cancer-
wisconsin
dataset
        train_data_breast-cancer-wisconsin
        test_data_breast-cancer-wisconsin
        test_Y_breast-cancer-wisconsin
        breast-cancer-wisconsin_model<-                train(class~
        ..,
        method="rpart",data=train_data_breast-cancer-wisconsin)
        y_breast-cancer-wisconsin=predict(breast-cancer-wisconsin_model,
newdata=test_data_breast-cancer-wisconsin)
        y1_breast-cancer-wisconsin<-as.character(y_breast-cancer-wisconsin)
        print (y1_breast-cancer-wisconsin)
        holdout_acc$breast-cancer-wisconsin[n]<-mean(test_Y_breast-cancer-
wisconsin==y1_breast-cancer-wisconsin)
    }

```

xVal-10f: 10-fold cross-validation method.

```

xVal_acc<-c()
for(n in 1:20){
  iris_model<train(class ~., data=iris_data,
  Iris dataset method="report",parms=list(split="information", tr.control=trctrl))
  accuracy<iris_model[["resample"]][["Accuracy"]]
  xVal_acc$iris[n]<- accuracy
}
xVal_acc<-c()
for(n in 1:20){
  teeth_model<train(Animal ~., data=teeth_data,
  Teeth dataset method="report",parms=list(split="information", tr.control=trctrl))
  accuracy<teeth_model[["resample"]][["Accuracy"]]
  xVal_acc$teeth[n]<- accuracy
}
xVal_acc<-c()
for(n in 1:20){
  wine_model<train(class ~., data=wine_data,
  Wine dataset method="report",parms=list(split="information", tr.control=trctrl))
  accuracy<wine_model[["resample"]][["Accuracy"]]
  xVal_acc$wine[n]<- accuracy
}
xVal_acc<-c()
for(n in 1:20){
  breast-cancer- breast-cancer-wisconsin_model<train(Alcohol ~., data=breast-cancer-wisconsin_data,
  wisconsin method="report",parms=list(split="information", tr.control=trctrl))
  dataset accuracy<breast-cancer-wisconsin_model[["resample"]][["Accuracy"]]
  xVal_acc$breast-cancer-wisconsin[n]<- accuracy
}

```

LOOCV: leave-one-out cross-validation method

```

trctr1<- trainControl(method = "LOOCV")
for (n in 1:20){
  iris_model<- train(class~., data=iris_data, method="report",
Iris dataset parms=list(split="information"),tr.control=trctr1)
  accuracy<-iris_model[["results"]][[Accuracy]][1]
  loocv_acc$iris[n]<-accuracy
}
trctr1<- trainControl(method = "LOOCV")
for (n in 1:20){
  teeth_model<- train(Animal~., data=teeth_data, method="report",
Teeth parms=list(split="information"),tr.control=trctr1)
dataset accuracy<-teeth_model[["results"]][[Accuracy]][1]
  loocv_acc$teeth[n]<-accuracy
}
trctr1<- trainControl(method = "LOOCV")
for (n in 1:20){
  wine_model<- train(class~., data=wine_data, method="report",
Wine dataset parms=list(split="information"),tr.control=trctr1)
  accuracy<-wine_model[["results"]][[Accuracy]][1]
  loocv_acc$wine[n]<-accuracy
}
for (n in 1:20){
breast- breast-cancer-wisconsin_model<- train(Alcohol~., data=breast-cancer-
cancer- wisconsin_data, method="report",
wisconsin parms=list(split="information"),tr.control=trctr1)
dataset accuracy<-breast-cancer-wisconsin_model[["results"]][[Accuracy]][1]
  loocv_acc$breast-cancer-wisconsin[n]<-accuracy}

```

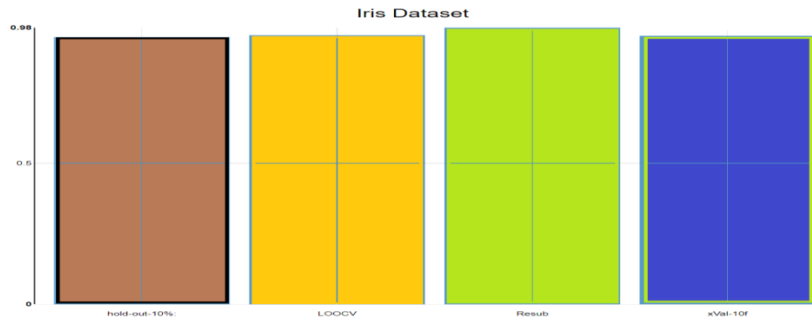


Figure13: chart of Iris Dataset

4- CONCLUSION

In this paper, we used the practical differences of the four validation methods using 4 distinct datasets: Iris dataset, the Wine dataset, the Breast-Cancer-Wisconsin dataset, and the Teeth dataset. These datasets are analyzed for classification and prediction data to find the optimal solution for data classifying. The performance indicators—accuracy, specificity, sensitivity, precision, and the error rate—are calculated for the given dataset. Accuracy, besides a proper data preprocessing technique, can get better the accuracy of the classifier. The function of data normalization had a noticeable impact on categorization performance and considerably enhanced the performance of the Iris dataset technique. The performance of the Teeth dataset method has minimum accuracy. Based on the analysis, the performances of the validation methods are analyzed. The results show that the performance of the Iris dataset technique is significantly superior to the other three techniques for the classification of data. To improve the overall accuracy, it is necessary to use more data set with a large number of attributes and use the best feature selection method in the future. Future works may also include hybrid classification models by combining some of the dataset techniques.

REFERENCES

- [1]. Bilalli, B., Abelló, A., Aluja-Banet, T. and Wrembel, R., 2016, April. Towards Intelligent Data Analysis: The Metadata Challenge. In *IoTBD* (pp. 331-338).
- [2]. Govender, P. and Sivakumar, V., 2020. Application of k-means and hierarchical clustering techniques for analysis of air pollution: A review (1980–2019). *Atmospheric Pollution Research*, 11(1), pp.40-56.
- [3]. David, S.K., Saeb, A.T., Rafiullah, M. and Rubeaan, K., 2019. Classification Techniques and Data Mining Tools Used in Medical Bioinformatics. In *Big Data Governance and Perspectives in Knowledge Management* (pp. 105-126). IGI Global.
- [4]. Brijain, M., Patel, R., Kushik, M. and Rana, K., 2014. A survey on decision tree algorithm for classification.
- [5]. Liao, T.W., 2005. Clustering of time series data—a survey. *Pattern recognition*, 38(11), pp.1857-1874.
- [6]. Dubes, R. and Jain, A.K., 1976. Clustering techniques: the user's dilemma. *Pattern Recognition*, 8(4), pp.247-260.
- [7]. Sanakal Ravi and Jayakumari T, 2014. Prognosis of Diabetes Using Data Mining Approach-Fuzzy C Means Clustering and Support Vector Machine. *International Journal of Computer Trends and Technology*.
- [8]. Sharma Arvind and Gupta PC, 2012. Predicting the Number of Blood Donors through their Age and Blood Group by using Data Mining Tool. *International Journal of Communication and Computer Technologies*. (01),pp.6-10
- [9]. Salama GI, Abdelhalim MB, Zeid MA 2012. Experimental comparison of classification for breast cancer diagnosis. *International Conference on Computer Engineering & Systems*. pp, (98-180-5)