

CHRONIC KIDNEY DISEASE PREDICTION USING MACHINE LEARNING

*Report submitted to the SASTRA Deemed to be University
as the requirement for the course*

EEE300 - MINI PROJECT

Submitted by

Muttumula Sai Vardhan Reddy
(Reg. No.: 123005092, EEE)
Edamadaka Jahnavi
(Reg. No.: 123005032, EEE)

December 2022



**SCHOOL OF ELECTRICAL &
ELECTRONICS ENGINEERING**

THANJAVUR, TAMIL NADU, INDIA – 613 401



SASTRA

ENGINEERING · MANAGEMENT · LAW · SCIENCES · HUMANITIES · EDUCATION

DEEMED TO BE UNIVERSITY
(U/S 3 OF THE UGC ACT, 1956)

THINK MERIT | THINK TRANSPARENCY | THINK SASTRA

SCHOOL OF ELECTRICAL & ELECTRONICS ENGINEERING

THANJAVUR - 613401

Bonafide Certificate

This is to certify that the report titled "**Chronic Kidney Disease Prediction Using Machine Learning**" submitted as a requirement for the course, EEE300: **MINI PROJECT** for B.Tech. is a bonafide record of the work done by

Mr. Muttumula Sai Vardhan Reddy (Reg. No. 123005092, EEE)

Ms. Edamadaka Jahnavi (Reg. No. 123005032, EEE)

during the academic year 2022-23, in the School of ELECTRICAL AND ELECTRONICS
ENGINEERING, under my supervision.

Signature of Project Supervisor :

Name with Affiliation : Dr. Venkatesh T, Asst Professor – III, SEEE

Date : 06/12/22

Mini Project *Viva voce* held on _____

Examiner 1

Examiner 2

Acknowledgements

First of all, I would like to thank *God Almighty* for his endless blessings.

I would like to express my sincere gratitude to **Dr S. Vaidya Subramaniam, Vice-Chancellor** for his encouragement during the span of my academic life at SASTRA Deemed University.

I would forever remain grateful and I would like to thank **Dr K. Thenmozhi, Dean, School of Electrical and Electronics Engineering** and **R. Chandramouli, Registrar** for their overwhelming support provided during my course span in *SASTRA Deemed University*.

I am extremely grateful to Lt. **Dr K. vijayrekha, Associate Dean, School of Electrical and Electronics Engineering** for her constant support, motivation and academic help extended for the past three years of my life in *School of Electrical and Electronics Engineering*.

I would specially thank and express my gratitude to **Dr T. Venkatesh, Assistant Professor-III, School of Electrical and Electronics Engineering** for providing me an opportunity to do this project and for his guidance and support to successfully complete the project.

I also thank all the *Teaching and Non-teaching faculty*, and all other people who have directly or indirectly help me through their support, encouragement and all other assistance extended for completion of my project and for successful completion of all courses during my academic life at *SASTRA Deemed University*.

List of Figures

Figure No.	Title	Page No.
1.1	Work Flow model for CKD Prediction before applying SMOTE	5
1.2	Work Flow model for CKD Prediction after applying SMOTE	5
4.1	Sample Image of Execution of Source code in Jupyter Notebook	50
4.2	Heat Maps of Correlation between Selected Features in case of Full feature selection	51
4.3	Heat Maps of Correlation between Selected Features in case of CFS	51
4.4	Heat Maps of Correlation between Selected Features in case of Wrapper FS	51
4.5	Heat Maps of Correlation between Selected Features in case of LASSO FS	52
4.6	Heat Maps of Correlation between Selected Features in case of Full features with SMOTE	52
4.7	Heat Maps of Correlation between Selected Features in case of LASSO FS with SMOTE	52
4.8	Accuracy, Precision, Recall for ML classifiers in case of Full feature selection	53
4.9	Accuracy, Precision, Recall for ML classifiers in case of CFS	53
4.10	Accuracy, Precision, Recall for ML classifiers in case of Wrapper FS	54
4.11	Accuracy, Precision, Recall for ML classifiers in case of LASSO FS	54
4.12	Accuracy, Precision, Recall for ML classifiers in case of Full features with SMOTE	55
4.13	Accuracy, Precision, Recall for ML classifiers in case of LASSO FS with SMOTE	55
4.14	Classification Error for ML classifiers in case of Full feature selection	56
4.15	Classification Error for ML classifiers in case of CFS	56
4.16	Classification Error for ML classifiers in case of Wrapper FS	56
4.17	Classification Error for ML classifiers in case of LASSO FS	57
4.18	Classification Error for ML classifiers in case of Full features with SMOTE	57
4.19	Classification Error for ML classifiers in case of LASSO FS with SMOTE	57
4.20	F- Measure for ML classifiers in case of Full feature selection	58

4.21	F- Measure for ML classifiers in case of CFS	58
Figure No.	Title	Page No.
4.22	F- Measure for ML classifiers in case of Wrapper FS	58
4.23	F- Measure for ML classifiers in case of LASSO FS	59
4.24	F- Measure for ML classifiers in case of Full features with SMOTE	59
4.25	F- Measure for ML classifiers in case of LASSO FS with SMOTE	59
4.26	AUC for ML classifiers in case of Full feature selection	60
4.27	AUC for ML classifiers in case of CFS	60
4.28	AUC for ML classifiers in case of Wrapper FS	60
4.29	AUC for ML classifiers in case of LASSO FS	61
4.30	AUC for ML classifiers in case of Full features with SMOTE	61
4.31	AUC for ML classifiers in case of LASSO FS with SMOTE	61
4.32	GINI Index for ML classifiers in case of Full feature selection	62
4.33	GINI Index for ML classifiers in case of CFS	62
4.34	GINI Index for ML classifiers in case of Wrapper FS	62
4.35	GINI Index for ML classifiers in case of LASSO FS	63
4.36	GINI Index for ML classifiers in case of Full features with SMOTE	63
4.37	GINI Index for ML classifiers in case of LASSO FS with SMOTE	63
6.1	Screenshot of Indexing	66

List of Tables

Table No.	Table name	Page No.
1.1	Confusion Matrix for prediction of CKD	3
1.2	Description of Attributes in the Dataset	6
2.1	Accuracy Table of Classifiers in case of different FS	9

Abbreviations

AI	Artificial Intelligence
ANN	Artificial Neural Network
AUC	Area Under Curve
CFS	Correlation based Feature Selection
CHAID	Chi-square Automatic Interaction Detection
CKD	Chronic Kidney disease
FN	False Negative
FP	False Positive
FS	Feature Selection
KNN	K-Nearest Neighbor
LASSO	Least Absolute Shrinkage and Selection Operator
LASSO FS	Least Absolute Shrinkage and Selection Operator Feature Selection
LR	Logistic Regression
LSVM	Linear Support Vector Machine
LSVM L1	Linear Support Vector Machine with penalty L1
LSVM L2	Linear Support Vector Machine with penalty L2
ML	Machine Learning
RF	Random Forest
ROC	Receiver Operating a Characteristic
SMOTE	Synthetic Minority class Oversampling Technique
TN	True Negative
TP	True Positive
UCI	UC Irvine repository

Abstract

Illness detection has become a crucial component in the medical industry because treating a disease can only be accomplished once it has been identified. One of the primary diseases that needs to be caught early on is chronic kidney disease (CKD). In today's world, machine learning (ML) is essential for disease identification. Therefore, the major goal of this work is to identify the top machine learning model for CKD prediction. A dataset from the UCI repository was gathered for this viewpoint. We used the K-Nearest Neighbor (KNN), Logistic Regression (LR), Linear Support Vector Machine (LSVM) with penalties L1 and L2, Random Forest (RF), CHAID, C 4.5, and Artificial Neural Network (ANN) classifier methods, as well as ML for prediction. To pick the features, we use feature selection (FS) techniques as full feature selection, correlation-based feature selection (CFS), wrapper FS, and LASSO FS. The target class is balanced using a data balancing technique on a dataset for the SMOTE. Accuracy, precision, recall, error, F-measure, AUC, and GINI index are used to evaluate the performance of classifiers.

The findings for each feature selection approach and all classifier methods are computed in this study. The highest-performing algorithms with the best feature selection techniques are chosen based on the results, and their performances are computed in each instance using SMOTE. We deduced from the data that LR, LSVM L1, LSVM L2, CHAID, C 4.5, and RF outperform Full Features and LASSO Feature Selection in terms of accuracy measure and other evaluation metrics. Random Forest performs best with an accuracy of 98.49% after applying SMOTE with LASSO Feature Selection, whereas LSVM with penalty L2 has an accuracy of 97.73%.

KEY WORDS: CKD, prediction, ML Classifier algorithms, FS methods, SMOTE.

Table of Contents

Title	Page No.
Bonafide Certificate	ii
Acknowledgements	iii
List of Figures	iv
List of Tables	vi
Abbreviations	vii
Abstract	viii
1. Summary of the base paper	1
2 Merits and Demerits of the base paper	7
3 Source Code	10
4 Snapshots	50
5 Conclusion and Future Plans	64
6 References	. 65
7 Appendix -Base Paper	66

CHAPTER 1

SUMMARY OF THE BASE PAPER

1.1. Problem Statement

CKD cannot be easily recognised unless the kidney is severely damaged. We are unable to routinely assess our disease status in this high-risk population.

We are unable to anticipate the sickness using current techniques. We are able to accurately predict the disease using ML classifier techniques. Data set collection & Pre-processing

1.1.1. Data Set

The UC Irvine Machine Learning Repository was offering the CKD dataset, which was obtained for this study from the UCI website. 400 instances and 25 attributes make up this dataset. One of those 25 qualities is a goal attribute, while the other 24 are feature attributes. Both CKD and Not CKD are valid values for the goal attribute. Table 1.2 lists the 24 feature qualities and their descriptions.

1.1.2. KNN Imputation

The first and most important thing we need to do with the dataset is find any missing values and replace them. Data cleansing is what this is. Predicting the missing values with an ML model is one of the finest methods for cleaning data. In this study, we predict missing values using the KNN method.

1.1.3. Label Encoder

Label encoding is the process of turning every label into a numeric form that will then be transformed into a machine-readable form. It is among the greatest methods for preparing data. In this paper, all the labels are converted into float or int types.

1.1.4. Train- Test split

Train-Test Split is the most frequently used approach for preprocessing. In this process, the data is divided into two parts, one for training the model and the other for testing. It is mainly used for estimating the performance of ML classifiers. In this paper, we consider the train-test split ratio to be 50:50.

1.2. Proposed Methodology

1.2.1. Feature Selection Methods

1.2.1.1. Full Feature Selection

In this case, we will consider all features as input for ML classifiers for prediction.

1.2.1.2. Correlation based FS (CFS)

CFS is one of the filter methods to select the features based on the correlation

between them. In this process, it will eliminate the highly correlated feature with the other feature. It selects only the remaining features after eliminating the highly correlated ones.

1.2.1.3. Wrapper FS

The feature selection in Wrapper FS is based on an accurate ML classifier. For creating subgroups, it employs a greedy strategy. Both the forward selection approach and the backward elimination method can be used to create subsets. In this paper, "forward selection" is used. Forward selection starts with a null subset and incrementally adds features.

1.2.1.4. LASSO FS

LASSO (least absolute shrinkage and selection operator) (least absolute shrinkage and selection operator) FS is a procedure that is integrated in FS. In the embedded approach, decision trees are used to pick features. This method involves the regularisation and feature selection phases. The coefficients of insignificant features are reduced to zero during regularisation, and the zero-coefficient features are disregarded during feature selection.

1.2.2. Machine Learning Classifiers

1.2.2.1. LR

Logistic regression is one of supervised ML classifiers. It's a mathematical model. It predicts the probability of achieving the target value. It categorizes the target attribute as either successful or unsuccessful. It returns 1 when successful, and 0 when unsuccessful.

1.2.2.2. KNN

KNN is a simplest supervised algorithm. For both classification and regression problems KNN is applicable. It is, however, primarily used for solving classification problems. The distance between two already labelled data sets is calculated. The distance aids in locating the new data nearest neighbor. Distance is calculated using the Euclidian method.

1.2.2.3. LSVM

LSVM is a latest, uniquely fast machine learning algorithm that uses a simple iterative approach to solve multiclass classification problems for large datasets. The SVM model is built in the dataset's linear CPU time. There are two penalty cases, L1 and L2.

1.2.2.4. CHAID

This decision tree technique is chi-square automatic interaction detection (CHAID). It's used to figure out how variables are related to one another. CHAID can find the outcome using nominal, ordinal, and continuous data.

1.2.2.5. C 4.5

C4.5 is a decision tree in the sense that from input it generates the decision tree. The tree will have the specified count of branches. The structure of tree is used to find out the relationship between potential outcomes and features. This process is iterative and continues until no further splits are found.

1.2.2.6.ANN

ANN is a type of AI. It falls under the category of supervised ML. It has the same structure as the human brain. It can solve problems that human or statistical standards have failed to solve.

1.2.2.7.RF

Random forests, also called as random decision forests, is an ensemble learning classifier for regression and classification that works by constructing a big number of decision trees during training. For classification tasks, the random forest output is the class chosen by the majority of trees. The mean prediction value of the different trees is returned for regression tasks. RF outperform decision trees in general.

1.2.3. SMOTE

Data balancing is the major problem one is facing while using dataset for prediction. There are 2 methods for data balancing, oversampling and under sampling. In this case we use oversampling. For oversampling the minority class, the SMOTE is used. It is also referred to as a data balancer. It accepts the entire dataset as input but works only on the minority class. It raises the proportion of minorities in the population. This methodology expands the features available for target class.

1.2.4. Performance Evaluation Metrics

1.2.4.1. Confusion Matrix for prediction of CKD

	CKD (1)	Not CKD (0)
CKD (1)	TP	FN
Not CKD (0)	FP	TN

Table. 1.1. Confusion Matrix for prediction of CKD

TP: It indicates that output was CKD and the prediction was classified correctly.

TN: It indicates that output was not CKD and the prediction was classified correctly.

FP: It indicates that output was CKD and the prediction was classified incorrectly.

FN: It indicates that output was not CKD and the prediction was classified incorrectly.

1.2.4.2. Accuracy:

It shows the right rate of prediction findings.

$$\text{Accuracy} = ((\text{TP} + \text{TN}) / (\text{TP} + \text{TN} + \text{FP} + \text{FN})) * 100$$

1.2.4.3.Precision:

It is the proportion of related instances in the total number of retrieved instances. It is a predicted positive value.

$$\text{Precision} = (\text{TP} / (\text{TP} + \text{FP})) * 100$$

1.2.4.4.Recall:

It is the proportion of related instances among all retrieved instances.

$$\text{Recall} = (\text{TP} / (\text{TP} + \text{FN})) * 100$$

1.2.4.5.F- Measure:

It is referred to as the F Score. It is used to determine the test accuracy.

$$\text{F- Measure} = 2 * ((\text{precision} * \text{recall}) / (\text{precision} + \text{recall}))$$

1.2.4.6.AUC:

ROC curve, a graph that compares the true positive rate to the false positive rate at various classification thresholds. The whole area under the ROC curve is referred to as the AUC.

1.2.4.7.GINI Index:

It is used to assess the disparity in attribute value distributions. It calculates the impurity of a specific attribute as probability.

$$\text{GINI} = 2 * \text{AUC} - 1$$

1.2.4.8.Error:

The rate of incorrect predicted results is shown by the classification error.

$$\text{Error} = 1 - \text{accuracy}$$

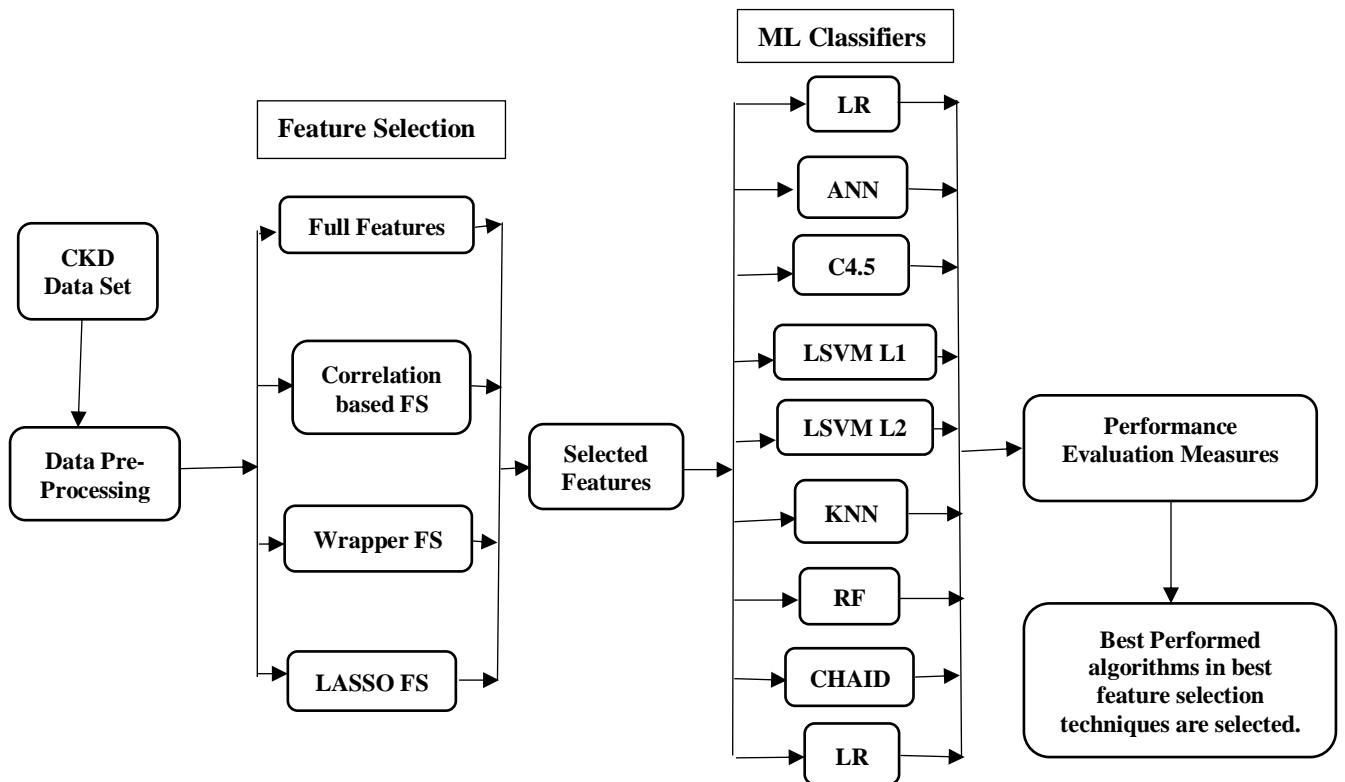


Fig. 1.1. Work Flow model for CKD Prediction before applying SMOTE

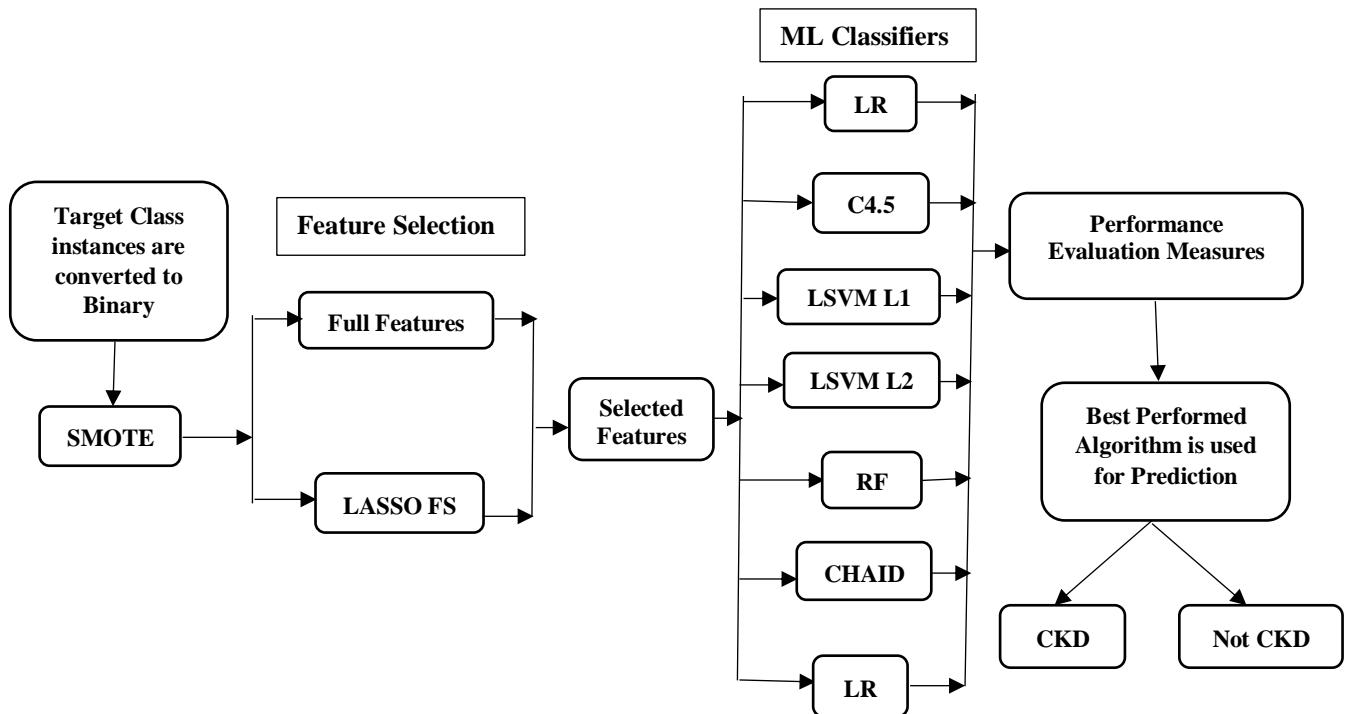


Fig. 1.2. Work Flow model for CKD Prediction after applying SMOTE

S. No.	Attribute Name	Description
1	Age	Age of patient
2	Bp	Blood pressure of patient
3	Sg	Urine specific gravity of patient
4	Al	Albumin level of patient
5	Su	Sugar level of patient
6	Rbc	Red blood cells normality in patient
7	Pc	Pus cell normality in patient
8	Pcc	Pus cell clumps existence in patient
9	Ba	Bacteria existence in patient
10	Bgr	Blood glucose random value of patient
11	Bu	Blood urea value of patient
12	Sc	Serum creatinine value of patient
13	Sod	Sodium level of patient
14	Pot	Potassium level of patient
15	Hemo	Hemoglobin level of patient
16	Pcv	Packed cell volume % of RBCs in circulating blood of patient
17	Wc	White blood cells count in patient
18	Rc	Red blood cells count in patient
19	Htn	Hypertension Yes or No for patient
20	Dm	Diabetes mellitus Yes or No for patient
21	Cad	Coronary artery disease Yes or No for patient
22	Appet	Appetite Yes or No for patient
23	Pe	Pedal edema Yes or No for patient
24	Ane	Anemia Yes or No for patient
25	Class	Target Variable (CKD or Not CKD)

Table. 1.2. Description of Attributes in the Dataset.

CHAPTER 2

MERITS AND DEMERITS OF THE BASE PAPER

2.1. Literature Review

G. Parthiban and K. R. A. Padmanaban et al. [9] identified Kidney disorder in its early stages by using Nave Bayes and Decision tree techniques on the selected features. However, the neural network prediction system's accuracy could be improved.

K. D. M. Perera, W. Gunarathne and K. A. D. C. P Kahandawaarachchi et al. [8] conducted a performance evaluation on Machine Learning Classification Techniques for classification of disease and Forecasting via Data Analytics for CKD and discovered that the Multiclass Decision Forest algorithm has the best accuracy of 98.1 percent.

The effects of using clinical features to classify patients with chronic kidney disease using the SVM algorithm were investigated by S. Shamiluulu, Y. Amirkaliyev and A. Serek et al. [7]. The implemented classifier is capable of achieving an overall performance of 94.602 percent.

M. Almasoud and T. E. Ward et al. [6] used 10-fold cross-validation to train and test logistic regression, support vector machines, random forest, and gradient boosting algorithms in this study. The gradient boosting algorithm outperformed other algorithms in terms of F1-measure (99.1 percent), sensitivity (98.8 percent), and specificity (98.8 percent) (99.3 percent).

L. Lessa, A. Peixoto, L. Kilvia De Almeida, R. Gomes, and J. Celestino, among others, [3] They devised a method for detecting kidney failure early on, increasing the chances of successful treatment for these patients. Several Machine Learning techniques (DT, SVM, and RF) were employed for this purpose. However, the accuracy of the prediction was low and could be improved further.

S. Elavarthy, T. Kiran, S. Shankar, S. Verma, P. Ghuli, and colleagues [4] Logistic Regression, Random Forest Tree, K-Nearest Neighbor, and Neural Network are four ML techniques that they used. To determine the correct classifier for CKD diagnosis, they conducted a systematic study of the outcomes of all three classifiers. However, certain noisy values are missing from the data collection.

2.2. Merits

ML classifiers are commonly used for prediction. We used several FS methods to improve accuracy. Not every algorithm performs well with a well-balanced data set. As a result, we used SMOTE to balance the dataset.

2.3. Demerits

In general, ANN and KNN are better performing classifiers, but their performance on our data set was poor.

2.4. Results

In case of Full feature selection, for LR the accuracy, precision, recall, F-measure, AUC, GINI Index and Error are 97, 100, 94.96, 0.96, 97.47, 0.94 and 3.0 respectively. for KNN the accuracy, precision, recall, F-measure, AUC, GINI Index and Error are 66.5, 66.25, 89.07, 0.44, 61.20, 0.22 and 33.5 respectively. for LSVM penalty L1 the accuracy, precision, recall, F-measure, AUC, GINI Index and Error are 97.0, 99.1, 95.7, 0.96, 97.2, 0.94 and 3.0 respectively. for LSVM penalty L2 the accuracy, precision, recall, F-measure, AUC, GINI Index and Error are 97.5, 100, 95.8, 0.97, 97.9, 0.95 and 2.5 respectively. for RF the accuracy, precision, recall, F-measure, AUC, GINI Index and Error are 99.5, 99.1, 100, 0.99, 99.3, 0.98 and 0.5

respectively. for CHAID the accuracy, precision, recall, F-measure, AUC, GINI Index and Error are 96.2, 97.5, 96.3, 0.96, 96.1, 0.92 and 3.78 respectively. for C 4.5 the accuracy, precision, recall, F-measure, AUC, GINI Index and Error are 96.9, 97.5, 97.5, 0.97, 96.7, 0.93 and 3.03 respectively. for ANN the accuracy, precision, recall, F-measure, AUC, GINI Index and Error are 83, 94.3, 61.7, 0.74, 79.6, 0.59 and 17.0 respectively.

In case of CFS, for LR the accuracy, precision, recall, F-measure, AUC, GINI Index and Error are 94.5, 97.3, 93.27, 0.93, 94.7, 0.89 and 5.5 respectively. for KNN the accuracy, precision, recall, F-measure, AUC, GINI Index and Error are 64.5, 68.4, 74.7, 0.52, 62.08, 0.24 and 35.5 respectively. for LSVM penalty L1 the accuracy, precision, recall, F-measure, AUC, GINI Index and Error are 95.5, 97.4, 94.9, 0.94, 95.6, 0.91 and 4.5 respectively. for LSVM penalty L2 the accuracy, precision, recall, F-measure, AUC, GINI Index and Error are 91.5, 97.2, 88.2, 0.90, 92.2, 0.84 and 8.5 respectively. for RF the accuracy, precision, recall, F-measure, AUC, GINI Index and Error are 97, 95.1, 100, 0.96, 96.29, 0.92 and 3.0 respectively. for CHAID the accuracy, precision, recall, F-measure, AUC, GINI Index and Error are 95.45, 95.23, 97.56, 0.96, 94.7, 0.89 and 4.54 respectively. for C 4.5 the accuracy, precision, recall, F-measure, AUC, GINI Index and Error are 96.9, 98.75, 96.34, 0.97, 97.1, 0.94 and 3.0 respectively. for ANN the accuracy, precision, recall, F-measure, AUC, GINI Index and Error are 67.5, 60.0, 59.25, 0.59, 66.1, 0.032 and 32.5 respectively.

In case of Wrapper FS, for LR the accuracy, precision, recall, F-measure, AUC, GINI Index and Error are 90.5, 92.3, 91.5, 0.88, 90.24, 0.80 and 9.5 respectively. for KNN the accuracy, precision, recall, F-measure, AUC, GINI Index and Error are 64.0, 68.2, 73.9, 0.52, 61.6, 0.23 and 36 respectively. for LSVM penalty L1 the accuracy, precision, recall, F-measure, AUC, GINI Index and Error are 90.0, 93.04, 89.9, 0.87, 90.0, 0.8 and 10 respectively. for LSVM penalty L2 the accuracy, precision, recall, F-measure, AUC, GINI Index and Error are 90.5, 92.3, 91.5, 0.8, 90.2, 0.8 and 9.5 respectively. for RF the accuracy, precision, recall, F-measure, AUC, GINI Index and Error are 91.5, 91.8, 94.1, 0.89, 90.8, 0.81 and 8.5 respectively. for CHAID the accuracy, precision, recall, F-measure, AUC, GINI Index and Error are 94.6, 94.1, 97.5, 0.95, 93.7, 0.87 and 5.3 respectively. for C 4.5 the accuracy, precision, recall, F-measure, AUC, GINI Index and Error are 96.2, 96.3, 97.5, 0.96, 95.7, 0.91 and 3.7 respectively. for ANN the accuracy, precision, recall, F-measure, AUC, GINI Index and Error are 79.5, 66.94, 97.5, 0.88, 90.88, 0.81 and 20.5 respectively.

In case of LASSO FS, for LR the accuracy, precision, recall, F-measure, AUC, GINI Index and Error are 97.5, 100, 95.7, 0.97, 97.8, 0.95 and 2.5 respectively. for KNN the accuracy, precision, recall, F-measure, AUC, GINI Index and Error are 67.0, 70.8, 75.6, 0.57, 64.97, 0.29 and 33 respectively. for LSVM penalty L1 the accuracy, precision, recall, F-measure, AUC, GINI Index and Error are 97.5, 99.1, 96.6, 0.96, 97.7, 0.95 and 2.5 respectively. for LSVM penalty L2 the accuracy, precision, recall, F-measure, AUC, GINI Index and Error are 97.5, 99.13, 96.63, 0.96, 97.7, 0.95 and 2.5 respectively. for RF the accuracy, precision, recall, F-measure, AUC, GINI Index and Error are 98.0, 98.3, 98.3, 0.97, 97.9, 0.95 and 2 respectively. for CHAID the accuracy, precision, recall, F-measure, AUC, GINI Index and Error are 98.4, 98.7, 98.7, 0.98, 98.3, 0.96 and 1.5 respectively. for C 4.5 the accuracy, precision, recall, F-measure, AUC, GINI Index and Error are 93.1, 95.0, 93.9, 0.94, 92.9, 0.85 and 6.8 respectively. for ANN the accuracy, precision, recall, F-measure, AUC, GINI Index and Error are 87.5, 87.8, 80.2, 0.83, 86.3, 0.72 and 12.5 respectively.

In case of Full features with SMOTE, for LR the accuracy, precision, recall, F-measure, AUC, GINI Index and Error are 96.5, 100, 94.1, 0.95, 97.0, 0.94 and 3.5 respectively. for LSVM penalty L1 the accuracy, precision, recall, F-measure, AUC, GINI Index and Error are 95.5, 98.3, 98.3, 0.97, 85.8, 0.91 and 4.5 respectively. for LSVM penalty L2 the accuracy, precision, recall, F-measure, AUC, GINI Index and Error are 96.5, 99.1, 94.9, 0.95, 96.8, 0.93 and 3.5 respectively. for RF the accuracy, precision, recall, F-measure, AUC, GINI Index and Error are 97.5, 98.3, 97.4, 0.96, 97.5, 0.95 and 2.5 respectively. for CHAID the accuracy, precision, recall, F-measure, AUC, GINI Index and Error are 95.4, 98.7, 93.9, 0.96, 95.9, 0.9 and 4.5 respectively. for C 4.5 the accuracy, precision, recall, F-measure, AUC, GINI Index and Error are 62.1, 62.1, 100, 0.76, 50, 0 and 37.8 respectively.

In case of LASSO FS with SMOTE, for LR the accuracy, precision, recall, F-measure, AUC, GINI Index and Error are 96.9, 98.75, 96.3, 0.97, 97.1, 0.91 and 3.03 respectively. for LSVM penalty L1 the accuracy, precision, recall, F-measure, AUC, GINI Index and Error are 96.9, 98.7, 96.3, 0.97, 97.1, 0.91 and 3.03 respectively. for LSVM penalty L2 the accuracy, precision, recall, F-measure, AUC, GINI Index and Error are 97.7, 98.7, 97.5, 0.98, 97.7, 0.93 and 2.27 respectively. for RF the accuracy, precision, recall, F-measure, AUC, GINI Index and Error are 98.4, 97.6, 100, 0.98, 98, 0.95 and 1.51 respectively. for CHAID the accuracy, precision, recall, F-measure, AUC, GINI Index and Error are 95.45, 96.34, 96.34, 0.96, 95.17, 0.91 and 4.5 respectively. for C 4.5 the accuracy, precision, recall, F-measure, AUC, GINI Index and Error are 62.12, 62.12, 100, 0.76, 50, 0 and 37.8 respectively.

In the Table. 2.1. the accuracies of all classifiers in all cases are compared and for each case the best accurate algorithm is marked.

Classifiers / FS	Full Features	CFS	Wrapper	LASSO	Full Features with SMOTE	LASSO with SMOTE
LR	97.0%	94.5%	90.5%	97.5%	96.5%	96.97%
LSVM L1	97.0%	95.5%	90.0%	97.5%	95.5%	96.97%
LSVM L2	97.5%	91.5%	90.5%	97.5%	96.5%	97.73%
CHAID	96.22%	95.45%	94.67%	98.49%	95.46%	95.46%
C 4.5	96.97%	96.97%	96.22%	93.19%	62.13%	62.13%
RF	99.5%	97.0%	91.5%	98.0%	97.5%	98.49%
KNN	66.5%	64.5%	64.0%	67.0%	----	----
ANN	83.0%	67.5%	79.5%	87.5%	----	----
Best Accuracy	RF	RF	C 4.5	CHAID	RF	RF

Table. 2.1. Accuracy Table of Classifiers in case of different FS

CHAPTER 3

SOURCE CODE

```

# Data Pre-Processing

import pandas as pd
data = pd.read_csv(&quot;DataSet.csv&quot;)
data.head()
data.info()
dec = data[[&quot;SK_ID_CURR&quot;]]
list = []
for i in data:
    if data[i].dtype == &lt;class 'object'>:
        list.append(i)
list

# Label Encoder
from sklearn.preprocessing import LabelEncoder
labelencoder_RF = LabelEncoder()
for i in list:
    data[i] = labelencoder_RF.fit_transform(data[i])
data
type(data)
list = [x for x in data]
data.info()
# KNN Imputer
from sklearn.impute import KNNImputer
imputer = KNNImputer(n_neighbors=20, metric = &quot;euclidean&quot;)
DataNoNulls = DataFrame(imputer.fit_transform(data))
Data = pd.DataFrame()
for i in range(len(list)):
    Data[list[i]] = DataNoNulls[i]
    Data[<class 'pandas.core.series.Series'>] = Data[<class 'pandas.core.series.Series'>].dtypes
    Data.head()
    Data.info()

print(confusion_matrix(y_test,y_predict))
print(&quot;Correct Predictions :&quot;+str(y_predict==y_test))
print(&quot;Incorrect Predictions :&quot;+str(y_predict!=y_test))

# SVM 13 in case of Full Features
from sklearn.svm import LinearSVC
lsvm1 = LinearSVC(C=0.01, False, penalty = &quot;l1&quot;,&gt;max_iter = 10000)
lsvm1.fit(x_train,y_train)

y_predict = lsvm1.predict(x_test)
accuracy_lvm1 = accuracy_score(y_test,y_predict)
precision_lvm1 = precision_score(y_test,y_predict,pos_label=0)
recall_lvm1 = recall_score(y_test,y_predict,pos_label=0)
fmea_lvm1 = f1_score(y_test,y_predict,zero_division = 1)
accuracy_lvm1+=100
precision_lvm1+=100
recall_lvm1+=100
error_lvm1 = 100 - accuracy_lvm1
print(&quot;Accuracy :&quot;+str(accuracy_lvm1))
print(&quot;Precision :&quot;+str(precision_lvm1))
print(&quot;Recall :&quot;+str(recall_lvm1))
print(&quot;F1 Score :&quot;+str(fmea_lvm1))
print(&quot;Error :&quot;+str(error_lvm1))
auc_lvm1 = roc_auc_score(y_test,y_predict)*100
print(&quot;AUC :&quot;+str(auc_lvm1))
gini_lvm1 = 2*(auc_lvm1)-1
print(&quot;GINI :&quot;+str(gini_lvm1))
print(&quot;Confusion matrix :&quot;+str(cnf_matrix))
print(confusion_matrix(y_test,y_predict))

# SVM 12 in case of Full Features
from sklearn.svm import LinearSVC

```

```

y_predict.append(1)
y_testing.append(0)
fp+=1
else:
    y_predict.append(0)
    y_testing.append(1)
fn+=1

accuracy_chaid = (tp+tn)/(tp+tn+fp+fn)
precision_chaid = tp/(tp+fp)
recall_chaid = tp/(tp+fn)
accuracy_chaid * 100
precision_chaid * 100
recall_chaid * 100
print(accuracy_chaid)
print(precision_chaid)
print(recall_chaid)
error_chaid = 100 - accuracy_chaid
print('Error:',error_chaid)

fmes_chaid = 2*((precision_chaid*recall_chaid)/[(precision_chaid+recall_chaid)])/(100)

```

```

y_predict = []
y_testing = []
for i in range(1, test_index):
    pred = c.predict([CATree,X1_test.loc[i,:]])
    if pred == 1:
        y_testing.append(1)
    else:
        y_testing.append(0)
    y_predict.append(1)
    if pred == 0:
        tp += 1
    else:
        fp += 1
    y_testing.append(0)
    tn += 1
    else:
        fn += 1
    y_predict.append(0)
    y_testing.append(1)
accuracy_c45 = (tp+tn)/(tp+tn+fp+fn)
precision_c45 = tp/(tp+fp)
recall_c45 = tp/(tp+fn)
print(precision_c45)
print(recall_c45)
print(accuracy_c45)

```

```

error_c45 = 100 - accuracy_c45
print(f'Error: {error_c45} | Error: {error_c45}')
fme_w_c45 = 2 * ((precision_c45 * recall_c45) / (precision_c45 + recall_c45)) / 100
print(f'me_w_c45')

auc_c45 = roc_auc_score(y_true=y_testing, y_pred=prediction)*100
print(f'AUC: {auc_c45}')
print(f'F1 Score: {fme_w_c45*100} | 1')

print(f'ANN in case of Full Features')

import keras
from keras.models import Sequential
from keras.layers import Dense
from sklearn.metrics import confusion_matrix, accuracy_score, precision_score, recall_score
import tensorflow as tf

# Create Sequential Model
classifier = Sequential()
# Add first layer
classifier.add(Dense(units = 6, kernel_initializer = tf.keras.initializers.he_uniform(), activation = tf.keras.activations.relu, input_dim = 2))
# Add second layer
classifier.add(Dense(units = 5, kernel_initializer = tf.keras.initializers.he_uniform(), activation = tf.keras.activations.relu))
# Add third layer
classifier.add(Dense(units = 1, kernel_initializer = tf.keras.initializers.he_uniform(), activation = tf.keras.activations.sigmoid))

# Compile the ANN
classifier.compile(optimizer = 'adam', loss = 'binary_crossentropy', metrics = [accuracy])

# Fit the ANN
y_train = y_train
y_test = y_test
y_train = y_train.astype('category')
y_train = y_train.cat.codes
y_test = y_test.astype('category')
y_test = y_test.cat.codes

# Train the ANN
classifier.fit(x_train, y_train, batch_size = 10, epochs = 100)

# Predict the test set results
y_pred = classifier.predict(x_test)
y_pred = (y_pred > 0.5)

```

```

accuracy_arn = accuracy_scorely_test
precision_arn = 100
precision_knn = precision_scorely_test
precision_knn = 100
precision_rf = precision_scorely_test
precision_rf = 100

print(f'Accuracy: {accuracy_arn}, Accuracy: {accuracy_scorely_test}')
print(f'Precision: {precision_arn}, Precision: {precision_scorely_test}')
print(f'Recall: {recall_arn}, Recall: {recall_scorely_test}')
print(f'F1 Score: {f1_score_arn}, F1 Score: {f1_score_scorely_test}')
print(f'Confusion matrix: {confusion_matrix(y_true,y_pred)}')
print(f'Confusion matrix: {confusion_matrix(y_true,y_pred_scorely_test)}')

auc_arn = roc_auc_scorely_test
print(f'AUC: {auc_arn}, AUC: {roc_auc_scorely_test}')
print(f'GINI: {gini_arn}, GINI: {gini_scorely_test}')
gini_arn = 2*(auc_arn/100) - 1
print(f'GINI: {gini_arn}, GINI: {gini_scorely_test}')

# Evaluation Measures Plotting in case
import matplotlib.pyplot as plt
import numpy as np

barWidth = 0.25

rg = plt.subplots(figsize=[12, 8])
accuracy = [accuracy_lrc,accuracy_knn,accuracy_rf]
precision = [precision_lrc,precision_knn,precision_rf]
precision_threshold = [precision_threshold_lrc,precision_threshold_knn,precision_threshold_rf]
recall = [recall_lrc,recall_knn,recall_rf]
f1_score = [f1_score_lrc,f1_score_knn,f1_score_rf]

tbc = np.arange(len(accuracy))

plt.bar(tbc, accuracy, barWidth)
plt.bar(tbc, precision, barWidth)
plt.bar(tbc, precision_threshold, barWidth)
plt.bar(tbc, recall, barWidth)
plt.bar(tbc, f1_score, barWidth)

plt.xlabel('Model')
plt.ylabel('Value')
plt.title('Performance Metrics Comparison')
plt.legend(['Accuracy', 'Precision', 'Precision Threshold', 'Recall', 'F1 Score'])

plt.show()

```

```

br2 = [x + barWidth for x in br1]

br3 = [x + barWidth for x in br2]

plt.bar(br1, accuracy, color="#839933", width=barWidth,
       edgecolor="#839933", label="#39;accuracy#39;),
plt.bar(br2, precision, color="#839933", width=barWidth,
       edgecolor="#839933", label="#39;precision#39;),
plt.bar(br3, recall, color="#839933", width=barWidth,
       edgecolor="#839933", label="#39;recall#39;),
plt.xlabel('#39;Algorithm#39;, fontweight="#39;bold#39;", fontsize= 15)
plt.ylabel('#39;Performance#39;, fontweight="#39;bold#39;", fontsize= 15)
plt.xticks([x + barWidth for x in range(len(accuracy))],
          [accuracy, f1Score, AUC, Kappa, FMeasure, Jaccard, CHAID, DBE,
           KS, ANN, RFE, SVM, L1, SVM_L1, SVM_2, SVM_3, SVM_L2, SVM_2_L1,
           SVM_3_L1, SVM_2_L2, SVM_3_L2], rotation=45)
plt.title('Accuracy Precision Recall for all features')
plt.legend()
plt.show()

f_mes = plt.figure()
x = f_mes.add_axe([0,0,0,0,3,0])
algo = [#39;LR#39;, #39;KNN#39;, #39;SVM_L1#39;, #39;SVM_L2#39;, #39;SVM_2#39;,
        #39;SVM_3#39;, #39;ANN#39;, #39;RFE#39;]
f_measures = [fmes_lr, fmes_knn, fmes_svm_l1, fmes_svm_l2, fmes_chaid,
              fmes_c45, fmes_an, fmes_rfe]
x.bar(algo,f_measures)

plt.xlabel('#39;Algorithm#39;, fontweight="#39;bold#39;", fontsize= 15)
plt.ylabel('#39;F Measures#39;, fontweight="#39;bold#39;", fontsize= 15)
plt.title('Variation of F Measure for Various Classifiers using Full Features')
plt.show()

f_auc = plt.figure()
x = f_auc.add_axe([0,0,0,0,3,0])
algo = [#39;LR#39;, #39;KNN#39;, #39;SVM_L1#39;, #39;SVM_L2#39;, #39;SVM_2#39;,
        #39;SVM_3#39;, #39;CHAID#39;, #39;DBE#39;, #39;KS#39;, #39;ANN#39;, #39;RFE#39;]
auc = [auc_lr, auc_knn, auc_svm_l1, auc_svm_l2, auc_chaid, auc_c45, auc_an, auc_rfe]

```

```

plt.xlabel("AUC", fontweight = "bold", fontsize = 15)
plt.ylabel("AUC", fontweight = "bold", fontsize = 15)
plt.title("Variation of AUC for Various Classifiers using Full Features")
plt.show()

algo = [knn, knn_knn, knn_isvm1, knn_tsvm2, knn_chaid, knn_c45, knn_ann, knn_rf]
plt.rcParams["format.figure.figsize"] = (60,10)

plt.plot(algo, gini)
plt.xlabel("Algorithm", fontweight = "bold", fontsize = 25)
plt.ylabel("GINI", fontweight = "bold", fontsize = 25)
plt.title("Variation of GINI for Various Classifiers using Full Features")
plt.show()

x = np.add(xes,[0,0,3,0,3])
algo = [knn, knn_knn, knn_isvm1, knn_tsvm2, knn_chaid, knn_c45, knn_ann, knn_rf]
gini = [gini, gini_knn, gini_isvm1, gini_tsvm2, gini_chaid, gini_c45, gini_ann, gini_rf]

plt.bar(algo,gini)
plt.xlabel("Algorithm", fontweight = "bold", fontsize = 15)
plt.ylabel("GINI", fontweight = "bold", fontsize = 15)
plt.title("Variation of GINI for Various Classifiers using Full Features")
plt.show()

err = [100-accuracy_lr, 100-accuracy_knn, 100-accuracy_isvm1, 100-accuracy_tsvm2, 100-accuracy_chaid, 100-accuracy_c45, 100-accuracy_ann, 100-accuracy_rf]

x_bar=algo[0]
x_bar[algo,errors]

plt.xlabel("Algorithm", fontweight = "bold", fontsize = 15)
plt.ylabel("Error", fontweight = "bold", fontsize = 15)
plt.title("Variation of Error for Various Classifiers using Full Features")
plt.show()

# CFS

def correlation(dataSetT, threshold):
    col_corr = set()
    corr_matrix = dataSetT.corr()

```

```

for i in range(len(ler.get_corr_matrix().columns)):
    for j in range(i):
        if abs(ler.get_corr_matrix().iloc[i][j])>=threshold :
            col_corr.add(ler.get_corr_matrix().columns[[j]])
return col_corr

corr_features = correlation(x_train,0.5)

corr_features

x_train_corr = x_train.drop(corr_features, axis = 1]
x1_train_corr = x1_train.drop(corr_features, axis = 1]
x_test_corr = x_test.drop(corr_features, axis = 1]
x1_test_corr = x1_test.drop(corr_features, axis = 1]

x1_train_corr
y1_train
x_train_corr
y_train

sns.heatmap(x_train_corr.corr(), annot = True)

x_train_corr
y1_train
# LR

from sklearn.linear_model import LogisticRegression
from sklearn.metrics import precision_score, accuracy_score, recall_score
import sys

lr = LogisticRegression(solver='liblinear', max_iter = sys.maxsize)

lr.fit(x_train_corr,y_train)
y_predict = lr.predict(x_test_corr)

accuracy_lr_corr = accuracy_score(y_test,y_predict)
precision_lr_corr = precision_score(y_test, y_predict, pos_label=0)
recall_lr_corr = recall_score(y_test, y_predict, pos_label=0)
fmeas_lr_corr = f1_score(y_test,y_predict,zero_division=1)

accuracy_lr_corr *= 100
precision_lr_corr *= 100
recall_lr_corr *= 100
fmeas_lr_corr *= 100

```

```

error_lr_corr = 100 - accuracy_lr_corr

print(@quot;Error : &quot;error_lr_corr)
print(@quot;Accuracy : &quot;accuracy_lr_corr)
print(@quot;Precision : &quot;precision_lr_corr)
print(@quot;Recall : &quot;recall_lr_corr)
print(@quot;F1 Score : &quot;fmeas_lr_corr)

auc_lr_corr = roc_auc_score(y_test, y_predict)*100
print(@quot;AUC : &quot;auc_lr_corr)

print(@quot;Confusion matrix:&quot;
print(ConfusionMatrix(y_true,y_predict))

print(@quot;Correct Predictions : &quot;sum(y_predict == y_test))
print(@quot;Incorrect Predictions : &quot;sum(y_predict != y_test))
auc_lr_corr = roc_auc_score(y_test, y_predict)*100
print(@quot;AUC : &quot;auc_lr_corr)

gini_lr_corr = 2 * (auc_lr_corr / (100 - 1))

print(@quot;GINI : &quot;gini_lr_corr)

# KNN

from sklearn.neighbors import KNeighborsClassifier

knn = KNeighborsClassifier(n_neighbors=5)
knn.fit(x_train, y_train)

y_predict_knn = knn.predict(x_test_corr)

accuracy_knn_corr = accuracy_score(y_test,y_predict_knn_corr)

precision_knn_corr = precision_score(y_test,y_predict_corr, pos_label=0)

recall_knn_corr = recall_score(y_test,y_predict_corr, pos_label=0)

fmeas_knn_corr = f1_score(y_test,y_predict_corr, zero_division=1)

accuracy_knn_corr *= 100

precision_knn_corr *= 100

recall_knn_corr *= 100

print(@quot;Accuracy : &quot;accuracy_knn_corr)
print(@quot;Precision:&quot;precision_knn_corr)

```

```

print(#Call:Recall,Error,specificity,recall_lnn_corr)
error_lnn_corr = 100 - accuracy_lnn_corr
print(error_lnn_corr)

print(#Error:Error,accuracy,recall_lnn_corr)
print(error_lnn_corr)

print(#F1 Score : f1_score_lnn_lnn_corr)
print(f1_score_lnn_lnn_corr)

print(#Confusion matrix:conf_lnn)
print(confusion_matrix(y_true,y_pred))

print(#Accuracy:accuracy_lnn_corr,y_pred_lnn_corr)
print(accuracy_lnn_corr)

print(#Precision:precision_lnn_corr,y_pred_lnn_corr)
print(precision_lnn_corr)

print(#Recall:recall_lnn_corr,y_pred_lnn_corr)
print(recall_lnn_corr)

print(#Specificity:specificity_lnn_corr,y_pred_lnn_corr)
print(specificity_lnn_corr)

print(#AUC:auc_lnn_corr,y_pred_lnn_corr)
print(auc_lnn_corr)

print(#ROC:roc_auc_lnn_corr,y_pred_lnn_corr)
print(roc_auc_lnn_corr)

print(#AUROC:auROC_lnn_corr,y_pred_lnn_corr)
print(auROC_lnn_corr)

print(#GINI:gin_lnn_corr,y_pred_lnn_corr)
print(gin_lnn_corr)

#L1 SVM L1

from sklearn.svm import LinearSVC

lsvm1 = LinearSVC(C=1, penalty = "l1")
lsvm1.fit(x_train, y_train)

y_predict_lnn1 = lsvm1.predict(x_test_lnn)

accuracy_lnn1_corr = accuracy_score(y_test_lnn, y_predict_lnn1)

precision_lnn1_corr = precision_score(y_test_lnn, y_predict_lnn1, pos_label=1)

recall_lnn1_corr = recall_score(y_test_lnn, y_predict_lnn1, pos_label=0)

f1_score_lnn1_corr = f1_score(y_test_lnn, y_predict_lnn1, zero_division=1)

accuracy_lnn1_corr = 100

precision_lnn1_corr = 100

recall_lnn1_corr = 100

error_lnn1_corr = 100 - accuracy_lnn1_corr

print(#Call:Recall,Error,specificity,recall_lnn1_corr)
error_lnn1_corr = 100 - accuracy_lnn1_corr
print(error_lnn1_corr)

print(#Error:Error,accuracy,recall_lnn1_corr)
print(error_lnn1_corr)

print(#F1 Score : f1_score_lnn1_lnn1_corr)
print(f1_score_lnn1_lnn1_corr)

print(#Confusion matrix:conf_lnn1)
print(confusion_matrix(y_true,y_pred))

print(#Accuracy:accuracy_lnn1_corr,y_pred_lnn1_corr)
print(accuracy_lnn1_corr)

print(#Precision:precision_lnn1_corr,y_pred_lnn1_corr)
print(precision_lnn1_corr)

print(#Recall:recall_lnn1_corr,y_pred_lnn1_corr)
print(recall_lnn1_corr)

print(#Specificity:specificity_lnn1_corr,y_pred_lnn1_corr)
print(specificity_lnn1_corr)

print(#AUC:auc_lnn1_corr,y_pred_lnn1_corr)
print(auc_lnn1_corr)

print(#ROC:roc_auc_lnn1_corr,y_pred_lnn1_corr)
print(roc_auc_lnn1_corr)

print(#AUROC:auROC_lnn1_corr,y_pred_lnn1_corr)
print(auROC_lnn1_corr)

print(#GINI:gin_lnn1_corr,y_pred_lnn1_corr)
print(gin_lnn1_corr)

```

```

print(&quot;confusion_matrix_y_test,y_predict_corr)&quot;
print(&quot;Correct Predictions :&quot;,sum(y_predict_corr == y_test))
print(&quot;incorrect Predictions :&quot;,sum(y_predict_corr != y_test))
auc_lvm1_1_corr = roc_auc_score(y_test,y_predict_corr)*100
print(&quot;AUC :&quot;,auc_lvm1_1_corr)
gini_lvm1_1_corr = 2*(auc_lvm1_1_corr)/100 - 1
print(&quot;GINI :&quot;,gini_lvm1_1_corr)
# L SVM L2

from sklearn.svm import LinearSVC
import sys

lsvm2 = LinearSVC(dual=False,penalty=&quot;#&39;#&39;#&39;#&39;,max_iter = 10000)
lsvm2.fit(x,y_train)

y_predict_corr = lsvm2.predict(text_corr)

accuracy_lvm2_2_corr = accuracy_score(y_test,y_predict_corr)
precision_lvm2_2_corr = precision_score(y_test,y_predict_corr,pos_label=0)
recall_lvm2_2_corr = recall_score(y_test,y_predict_corr,pos_label=0)
fmes_lvm2_2_corr = f1_score(y_test,y_predict_corr,zero_division =1)

accuracy_lvm2_2_corr *= 100
precision_lvm2_2_corr *= 100
recall_lvm2_2_corr *= 100
print(&quot;Accuracy :&quot;,accuracy_lvm2_2_corr)
print(&quot;Precision:&quot;,precision_lvm2_2_corr)
print(&quot;Recall :&quot;,recall_lvm2_2_corr )
print(&quot;Error_lvm2_2_corr = 100 - accuracy_lvm2_2_corr
print(&quot;Error :&quot;,100-accuracy_lvm2_2_corr
print(&quot;F1 Score :&quot;,fmes_lvm2_2_corr)
print(&quot;Confusion matrix:&quot;)
print(&quot;confusion_matrix_y_test,y_predict_corr)&quot;
print(&quot;Correct Predictions :&quot;,sum(y_predict_corr == y_test))
print(&quot;incorrect Predictions :&quot;,sum(y_predict_corr != y_test))
auc_lvm2_2_corr = roc_auc_score(y_test,y_predict_corr)*100

```

```

print(&quot;AUC :&quot;,auc_,auc_.item_Ml2_corr))

gini_lvm2_corr = 2*(gini_lvm1.item_Ml2_corr)/100 - 1

print(&quot;GINI :&quot;,gini_lvm2_corr)

## RF

from sklearn.ensemble import RandomForestClassifier

classifier=RandomForestClassifier(n_estimators=10, criterion=&quot;entropy&quot;
classifier.fit_(train_x,y_train)

y_predict_rf=classifier.predict_(test_x)

accuracy_rf_corr = accuracy_score(y_test,y_predict_rf)

precision_rf_corr = precision_score(y_test,y_predict_rf,pref_label=0)

recall_rf_corr = recall_score(y_test,y_predict_rf,pref_label=0)

fmes_rf_corr = f1_score(y_test,y_predict_rf,corr_ratio_division=1)

accuracy_rf_corr = 100

precision_rf_corr *= 100

recall_rf_corr *= 100

print(&quot;Accuracy :&quot;,accuracy_rf_corr)

print(&quot;Precision:&quot;,precision_rf_corr)

print(&quot;Recall:&quot;,recall_rf_corr)

error_rf_corr = 100 - accuracy_rf_corr

print(&quot;Error :&quot;,error_rf_corr)

print(&quot;F1 Score :&quot;,fmes_rf_corr)

print(&quot;Confusion matrix:&quot;)

print(confusion_matrix(test_y,y_predict_rf))

print(&quot;Correct Predictions :&quot;,sum(y_predict_rf == y_test))

print(&quot;Incorrect Predictions :&quot;,sum(y_predict_rf != y_test))

auc_rf_corr = roc_auc_score(y_test,y_predict_rf)*100

print(&quot;AUC :&quot;,auc_rf_corr)

gini_rf_corr = 2*(auc_rf_corr/100) - 1

print(&quot;GINI :&quot;,gini_rf_corr)

## ANN

import keras

```

```

# C4.5
config = {"algorithm": "C4.5",
          "C4Tree": c4f.mf("sony1_train", config)
          #print(type(ChadTree))
tp = 0
tn = 0
fp = 0
fn = 0

y_predict = []
y_testing = []

for i in x1_test_corr.index:
    pred = c4f.predict(C4Tree,x1_test_corr.loc[i,:])
    if pred == y1_test[i]:
        tp += 1
    if pred == -1 && y1_test[i] == 1:
        fn += 1
    if pred == 1 && y1_test[i] == -1:
        fp += 1
    if pred == -1 && y1_test[i] == -1:
        tn += 1

y_predict.append(tp)
y_testing.append(tp)

tn += 1
fp += 1
fn += 1

y_predict.append(0)
y_testing.append(0)

fp+=1
else:
    y_predict.append(0)
    y_testing.append(1)

fn+=1

accuracy_C45_corr = (tp+tn)/(tp+tn+fp+fn)

```

```

accuracy_c45_corr*100
print(accuracy_c45_corr)

precision_c45_corr + tp/(tp+fp)
precision_c45_corr*100

recall_c45_corr + tp/(tp+fn)
recall_c45_corr*100
print(recall_c45_corr)

print(recall_c45_corr)

error_c45_corr = 100 - accuracy_c45_corr
print(f'@quot;Error : &quot;{error_c45_corr}'')

fmeas_c45_corr =
    2*(precision_c45_corr*recall_c45_corr)/(precision_c45_corr+recall_c45_corr)*100
print(fmeas_c45_corr)

auc_c45_corr = roc_auc_score(y_testing, y_predict)*100
print(f'@quot;AUC : &quot;{auc_c45_corr}'')

gini_c45_corr = 2*[auc_c45_corr/100] - 1
print(f'@quot;Gini : &quot;{gini_c45_corr}'')

# Evaluation Measures Plotting in case of CFS
import matplotlib.pyplot as plt
import numpy as np

barWidth = 0.25
fig2 = plt.subplots(figsize=(12, 8))

accuracy = [accuracy_U_corr,accuracy_knn_corr,accuracy_lsvm1_corr,
            accuracy_lsvm2_corr,accuracy_ann_corr,accuracy_chaid_corr,accuracy_c45_corr,accuracy_rf_corr]
accuracy_U_corr = [accuracy_U_corr,accuracy_knn_corr,accuracy_lsvm1_corr,accuracy_lsvm2_corr,accuracy_ann_corr,accuracy_chaid_corr,accuracy_c45_corr,accuracy_rf_corr]

precision = [precision_U_corr,precision_knn_corr,precision_lsvm1_corr,precision_lsvm2_corr,
            precision_ann_corr,precision_chaid_corr,precision_c45_corr,precision_rf_corr]
precision_U_corr = [precision_U_corr,precision_knn_corr,precision_lsvm1_corr,precision_lsvm2_corr,precision_ann_corr,precision_chaid_corr,precision_c45_corr,precision_rf_corr]

recall = [recall_U_corr,recall_knn_corr,recall_lsvm1_corr,recall_lsvm2_corr,recall_ann_corr,recall_chaid_corr,recall_c45_corr,recall_rf_corr]
recall_U_corr = [recall_U_corr,recall_knn_corr,recall_lsvm1_corr,recall_lsvm2_corr,recall_ann_corr,recall_chaid_corr,recall_c45_corr,recall_rf_corr]

fmeas = [fmeas_U_corr,fmeas_knn_corr,fmeas_lsvm1_corr,fmeas_lsvm2_corr,fmeas_ann_corr,fmeas_chaid_corr,fmeas_c45_corr,fmeas_rf_corr]
fmeas_U_corr = [fmeas_U_corr,fmeas_knn_corr,fmeas_lsvm1_corr,fmeas_lsvm2_corr,fmeas_ann_corr,fmeas_chaid_corr,fmeas_c45_corr,fmeas_rf_corr]

br1 = np.arange(len(accuracy))
br2 = [x + barWidth for x in br1]
br3 = [x + 2*barWidth for x in br2]

plt.bar(br1, accuracy, color='red', width=barWidth, label='Accuracy')
plt.bar(br2, precision, color='blue', width=barWidth, label='Precision')
plt.bar(br3, recall, color='green', width=barWidth, label='Recall')

plt.xlabel('Classification Techniques')
plt.ylabel('Evaluation Measures')
plt.title('Comparison of Evaluation Measures for CFS')
plt.legend()
plt.show()

```



```

pred = c.predict([C45Tree,x1_test_wrp.loc,:])
if pred == > y1_test[0]:
    if pred == > 839.ckpt&#39;
        y_predict.append(1)
    y_testing.append(1)
    tp += 1
else:
    y_predict.append(0)
    y_testing.append(0)
tn += 1
else:
    if pred == > 839.ckpt&#39;
        y_predict.append(1)
    y_testing.append(1)
    tn += 1
else:
    if pred == > 839.ckpt&#39;
        y_predict.append(0)
    y_testing.append(0)
fpr+=1
else:
    y_predict.append(0)
    y_testing.append(1)
tn+=1
accuracy_c45_wrp = (tp+tn)/(tp+tn+fpr)
print(accuracy_c45_wrp)
precision_c45_wrp = tp/(tp+fp)
recall_c45_wrp = tp/(tp+fn)
print(precision_c45_wrp)
print(recall_c45_wrp)
error_c45_wrp = 100 - accuracy_c45_wrp
print(error_c45_wrp)
fmes_c45_wrp = roc_auc_score(y_testing,y_predict)*100
print(fmes_c45_wrp)
auc_c45_wrp = roc_auc_score(y_testing,y_predict)*100
print(auc_c45_wrp)

plt.show()
error.add_exes([0,0,0,0,0,0])
x = error.add_exes[[0,0,0,0,0,0]]
alg = [839.1R&#39;, 839.XNN&#39;, 839.LSVM1&#39;, 839.LSVM2&#39;, 839.CHAD&#39;, 839.C4.5&#39;, 839.ANN&#39;, 839.RF&#39;]
auc = [100-accuracy_lr_wrp,100-accuracy_knn_wrp,100-accuracy_lsvm1_wrp,100-accuracy_lsvm2_wrp,100-accuracy_chad_wrp,100-accuracy_c45_wrp,100-accuracy_rf_wrp]
x_bar(algo,auc)
plt.xlabel([839.Algorithm&#39;, fontweight = 839.bold&#39;, fontsize = 15])
plt.ylabel([839.Measure&#39;, fontweight = 839.bold&#39;, fontsize = 15])
plt.title([839.Variation of Error for Various Classifiers using Wrapper Feature Selection&quot;])
plt.show()

# Lasso F5
from sklearn.linear_model import Lasso
features = [i for i in data]
features.remove([839.SCKDB&#39;])
print(features)
lasso = Lasso(alpha = 0.01)
lasso.fit(x_train,y_train)
coefflasso.coef_
for i in range(len(coeff)):
    coeff[i] = abs(coeff[i])
df_lasso = pd.DataFrame([{"Features":features,"Coefficient":coeff}])
df_lasso.sort_values(by='Coefficient',ascending=False)
list_lasso = []
for i in range(len(df_lasso)):
    if df_lasso[[839.Coefficient&#39;][i]] == 0:
        list_lasso.append(df_lasso[[839.Features&#39;][i]])
11
list_lasso.append(df_lasso[[839.Features&#39;][0]])

fmes_lsvm1_lso = f1_score(y_testing,y_predict_lsvo.zero_division=1)
print(accuracy_lsvo_lsvo)
print(precision_lsvo_lsvo)
print(recall_lsvo_lsvo)
error_lsvo_lsvo = 100 - accuracy_lsvo_lsvo
print(F1 Score : &quot;fmes_lsvo_lsvo")
print(confusion_matrix_lsvo_lsvo)
print(confusion_matrix_lsvo_lsvo)
print(accuracy_lsvo_lsvo)
print(precision_lsvo_lsvo)
print(recall_lsvo_lsvo)
print(error_lsvo_lsvo)
print(F1 Score : &quot;fmes_lsvo_lsvo")
print(gini_lsvo_lsvo)
print(gini_lsvo_lsvo)
# RF
from sklearn.ensemble import RandomForestClassifier
classifier = RandomForestClassifier(n_estimators=10, criterion=&quot;entropy&quot;)
classifier.fit(x_train,y_train)
y_predict_rf = classifier.predict(x_test)
accuracy_rf_lso = accuracy_score(y_testing,y_predict_rf_lso)
precision_rf_lso = precision_score(y_testing,y_predict_rf_lso,labels=0)*100
recall_rf_lso = recall_score(y_testing,y_predict_rf_lso,labels=0)*100
fmes_rf_lso = f1_score(y_testing,y_predict_rf_lso,zero_division=1)
print(accuracy_rf_lso)
print(precision_rf_lso)
print(recall_rf_lso)
print(error_rf_lso)
print(F1 Score : &quot;fmes_rf_lso")
print(gini_rf_lso)
print(gini_rf_lso)
# SVM L2
from sklearn.svm import LinearSVC
import sys
lsvm12 = LinearSVC(dual=False,penalty='l2',C=1,max_iter = 12000)
lsvm12.fit(x_train,y_train)
y_predict_lsvo = lsvo12.predict(x_test)
accuracy_lsvo_lsvo = accuracy_score(y_testing,y_predict_lsvo)*100
precision_lsvo_lsvo = precision_score(y_testing,y_predict_lsvo,labels=0)*100
recall_lsvo_lsvo = recall_score(y_testing,y_predict_lsvo,labels=0)*100
fmes_lsvo_lsvo = f1_score(y_testing,y_predict_lsvo,zero_division=1)
print(accuracy_lsvo_lsvo)
print(precision_lsvo_lsvo)
print(recall_lsvo_lsvo)
error_lsvo_lsvo = 100 - accuracy_lsvo_lsvo
print(F1 Score : &quot;fmes_lsvo_lsvo")
print(gini_lsvo_lsvo)
print(gini_lsvo_lsvo)

# KNN
from sklearn.neighbors import KNeighborsClassifier
knn = KNeighborsClassifier(n_neighbors=5)
knn.fit(x_train,y_train)
y_predict_knn_lso = knn.predict(x_test)
accuracy_knn_lso = accuracy_score(y_testing,y_predict_knn_lso)*100
precision_knn_lso = precision_score(y_testing,y_predict_knn_lso,labels=0)*100
recall_knn_lso = recall_score(y_testing,y_predict_knn_lso,zero_division=1)
fmes_knn_lso = f1_score(y_testing,y_predict_knn_lso,zero_division=1)
print(accuracy_knn_lso)
print(precision_knn_lso)
print(recall_knn_lso)
error_knn_lso = 100 - accuracy_knn_lso
print(error_knn_lso)
print(F1 Score : &quot;fmes_knn_lso")
print(confusion_matrix_knn_lso)
print(incorrect_predictions_knn_lso)
print(correct_predictions_knn_lso)
print(incorrect_predictions_knn_lso)
print(correct_predictions_knn_lso)
print(error_knn_lso)
print(F1 Score : &quot;fmes_knn_lso")
print(gini_knn_lso)
print(gini_knn_lso)
# LSVM L1
from sklearn.svm import LinearSVC
lsvo11 = LinearSVC(dual=False,penalty='l1',max_iter = 10000)
lsvo11.fit(x_train,y_train)
y_predict_lsvo11 = lsvo11.predict(x_test)
accuracy_lsvo11_lso = accuracy_score(y_testing,y_predict_lsvo11)
precision_lsvo11_lso = precision_score(y_testing,y_predict_lsvo11,labels=0)*100
recall_lsvo11_lso = recall_score(y_testing,y_predict_lsvo11,labels=0)*100
fmes_lsvo11_lso = f1_score(y_testing,y_predict_lsvo11,zero_division=1)
print(accuracy_lsvo11_lso)
print(precision_lsvo11_lso)
print(recall_lsvo11_lso)
error_lsvo11_lso = 100 - accuracy_lsvo11_lso
print(error_lsvo11_lso)
print(F1 Score : &quot;fmes_lsvo11_lso")
print(gini_lsvo11_lso)
print(gini_lsvo11_lso)
# ANN
import keras
from keras.models import Sequential
from keras.layers import Dense
from keras.metrics import confusion_matrix,accuracy_score,precision_score,recall_score
classifier = Sequential()
classifier.add(Dense(units = 6, kernel_initializer = 839.uniform&#39;, activation = 839.relu&#39;, input_dim = 16))
classifier.add(Dense(units = 6, kernel_initializer = 839.uniform&#39;, activation = 839.relu&#39;))
classifier.add(Dense(units = 1, kernel_initializer = 839.uniform&#39;, activation = 839.sigmoid&#39;))
classifier.compile(optimizer = 839.adam&#39;, loss = 839.binary_crossentropy&#39;, metrics = [839.accuracy&#39;])
Y_train = y_train
Y_test = y_test
y_train = y_train.astype(&#39;category&#39;)
y_train = y_train.cat.codes
y_test = y_test.cat.codes
classifier.fit(x_train,y_train,batch_size = 10, epochs = 100)
y_pred = classifier.predict(x_test)
y_pred = (y_pred > 0.5)
accuracy_ann_lso = accuracy_score(y_testing,y_pred)*100
precision_ann_lso = precision_score(y_testing,y_pred)*100
recall_ann_lso = recall_score(y_testing,y_pred)*100
fmes_ann_lso = f1_score(y_testing,y_pred,zero_division=1)
print(accuracy_ann_lso)
print(precision_ann_lso)
print(recall_ann_lso)
print(F1 Score : &quot;fmes_ann_lso")
print(gini_ann_lso)
print(gini_ann_lso)
# GINI
fmes_gini_lso = f1_score(y_testing,y_pred,zero_division=1)
print(F1 Score : &quot;fmes_gini_lso")

```

```

print("\nConfusion matrix:")
print(confusion_matrix(y_true,y_pred))

auc_anno_lso = roc_auc_score(y_true, y_pred)*100
print('AUC :',auc_anno_lso)
gini_anno_lso = 2*(auc_anno_lso)/100
print('GINI :',gini_anno_lso)

# CHAID
from chaid import ChaidBoost as cf
from pandas.api.types import CategoricalDtype
config = {"algorithm": "CHAID"}
df = pd.DataFrame()
df = x1_train_lso
ChaidTree = cf.fit(df.join(y1_train), config)

tp = 0
tn = 0
fp = 0
fn = 0

y_predict = []
y_testing = []
for i in x1_test.index:
    pred = predict(ChaidTree,x1_test_lso.loc[i,:])
    if pred == y1_test[i]:
        tp += 1
    if pred == 0:
        tn += 1
    else:
        fp += 1
    fn += 1
    y_predict.append(pred)
    y_testing.append(y1_test[i])

print('Predicted values:',y_predict)
print('Actual values:',y_testing)

accuracy_lso = (tp+tn)/(tp+tn+fp+fn)
precision_lso = tp/(tp+fp)
recall_lso = tp/(tp+fn)
print('Accuracy_lso:',accuracy_lso)
print('Precision_lso:',precision_lso)
print('Recall_lso:',recall_lso)
print('F-measure_lso:',f1_score(y_testing,y_predict))
print('AUC_lso:',roc_auc_score(y_testing,y_predict)*100)
gini_lso = 2*(auc_lso/100)-1
print('GINI_lso:',gini_lso)

# Evaluation Measures Plotting in case of LASSO FS
fig2 = plt.subplots(figsize=(12,8))
accuracy = [accuracy_lso,accuracy_knn_lso,accuracy_lr_lso,accuracy_rf_lso]
accuracy_lso2_lso = [accuracy_lso,accuracy_knn_lso,accuracy_chaid_lso,accuracy_c45_lso,accuracy_rf_lso]
precision = [precision_lr_lso,precision_knn_lso,precision_rf_lso,precision_chaid_lso,precision_c45_lso,precision_rf_lso]
recall = [recall_lr_lso,recall_knn_lso,recall_rf_lso,recall_chaid_lso,recall_c45_lso,recall_rf_lso]
recall_anno_lso = recall_chaid_lso
recall_anno_lso2_lso = recall_chaid_lso,recall_c45_lso,recall_rf_lso
x1 = np.arange(len(accuracy))
x2 = [x + barWidth for x in x1]
x3 = [x + 2 * barWidth for x in x2]
plt.bar(x1, accuracy, color='darkblue', width=barWidth,
        edgecolor='black', precision_color='darkblue', width=barWidth,
        edgecolor='black')
plt.bar(x2, precision, color='darkblue', width=barWidth,
        edgecolor='black', precision_color='darkblue', width=barWidth,
        edgecolor='black')
plt.bar(x3, recall, color='darkblue', width=barWidth,
        edgecolor='black', recall_color='darkblue', width=barWidth,
        edgecolor='black')
plt.title('Variation of F measure for various classifiers using Lasso feature selection')
plt.xlabel('Classifiers')
plt.ylabel('F-measure')
plt.show()

# Variation of AUC for various classifiers using Lasso feature selection
f_mse = plt.figure()
x = f_mse.add_axes([0.0,0.3,0.3])
algo = [8495,8495,8495,8495,8495,8495]
gini = [8495,8495,8495,8495,8495,8495]
f_mse_measures = [fmeas_lr_lso,fmeas_knn_lso,fmeas_chaid_lso,fmeas_c45_lso,fmeas_rf_lso,fmeas_anno_lso]
fmeas_lr_lso,fmeas_knn_lso,fmeas_chaid_lso,fmeas_c45_lso,fmeas_rf_lso,fmeas_anno_lso
xbar(algo,f_measures)
ph.xlabel('Algorithm',fontweight='bold',fontsize = 15)
ph.ylabel('AUC',fontweight='bold',fontsize = 15)
ph.title('Variation of AUC for various classifiers using Lasso feature selection')
ph.show()

# Variation of GINI for various classifiers using Lasso feature selection
f_gini = plt.figure()
x = f_gini.add_axes([0.0,0.3,0.3])
algo = [8495,8495,8495,8495,8495,8495]
gini = [8495,8495,8495,8495,8495,8495]
f_gini_measures = [fmeas_lr_lso,fmeas_knn_lso,fmeas_chaid_lso,fmeas_c45_lso,fmeas_rf_lso,fmeas_anno_lso]
fmeas_lr_lso,fmeas_knn_lso,fmeas_chaid_lso,fmeas_c45_lso,fmeas_rf_lso,fmeas_anno_lso
xbar(algo,gini_measures)
ph.xlabel('Algorithm',fontweight='bold',fontsize = 15)
ph.ylabel('GINI',fontweight='bold',fontsize = 15)
ph.title('Variation of GINI for various classifiers using Lasso feature selection')
ph.show()

# Variation of Error for various classifiers using Lasso feature selection
f_error = plt.figure()
x = f_error.add_axes([0.0,0.3,0.3])
algo = [8495,8495,8495,8495,8495,8495]
error = [8495,8495,8495,8495,8495,8495]
f_error_measures = [fmeas_lr_lso,fmeas_knn_lso,fmeas_chaid_lso,fmeas_c45_lso,fmeas_rf_lso,fmeas_anno_lso]
fmeas_lr_lso,fmeas_knn_lso,fmeas_chaid_lso,fmeas_c45_lso,fmeas_rf_lso,fmeas_anno_lso
xbar(algo,error_measures)
ph.xlabel('Algorithm',fontweight='bold',fontsize = 15)
ph.ylabel('Error',fontweight='bold',fontsize = 15)
ph.title('Variation of error for various classifiers using Lasso feature selection')
ph.show()

# SMOTE
Data.head()
Data
y_train
y_train
np.bencount(y_train)
np.bencount(y_test)
from imblearn.over_sampling import SMOTE
sm = SMOTE(random_state=42)
x1_smote_train,y1_smote_train = sm.fit_resample(x1_train,y1_train)
x_smote_train,y_smote_train = sm.fit_resample(x_train,y_train)
np.bencount(y_smote_train)
x_smote_train
# Full features with SMOTE
sns.heatmap(y_smote_train.corr(), annot=True)
# LR
from sklearn.linear_model import LogisticRegression

# Confusion matrix
print('Confusion matrix')
print(confusion_matrix(y_true,y_pred))
print('Correct Predictions :',sum(y_pred==y_true))
print('Incorrect Predictions :',sum(y_pred!=y_true))
print('True Positive :',auc_smote_lsm1)
print('True Negative :',accuracy_smote_lsm1)
print('False Positive :',error_smote_lsm1)
print('False Negative :',gini_smote_lsm1)
print('F1 Score :',f1_score(y_true,y_smote_lsm1))
print('Confusion matrix')
print(confusion_matrix(y_true,y_smote_lsm1))
print('Correct Predictions :',sum(y_smote_lsm1==y_true))
print('Incorrect Predictions :',sum(y_smote_lsm1!=y_true))
print('True Positive :',auc_smote_lsm1)
print('True Negative :',accuracy_smote_lsm1)
print('False Positive :',error_smote_lsm1)
print('False Negative :',gini_smote_lsm1)
print('F1 Score :',f1_score(y_true,y_smote_lsm1))
print('Confusion matrix')
print(confusion_matrix(y_true,y_smote_lsm2))
print('Correct Predictions :',sum(y_smote_lsm2==y_true))
print('Incorrect Predictions :',sum(y_smote_lsm2!=y_true))
print('True Positive :',auc_smote_lsm2)
print('True Negative :',accuracy_smote_lsm2)
print('False Positive :',error_smote_lsm2)
print('False Negative :',gini_smote_lsm2)
print('F1 Score :',f1_score(y_true,y_smote_lsm2))
print('Confusion matrix')
print(confusion_matrix(y_true,y_smote_lsm3))
print('Correct Predictions :',sum(y_smote_lsm3==y_true))
print('Incorrect Predictions :',sum(y_smote_lsm3!=y_true))
print('True Positive :',auc_smote_lsm3)
print('True Negative :',accuracy_smote_lsm3)
print('False Positive :',error_smote_lsm3)
print('False Negative :',gini_smote_lsm3)
print('F1 Score :',f1_score(y_true,y_smote_lsm3))
print('Confusion matrix')
print(confusion_matrix(y_true,y_smote_lsm4))
print('Correct Predictions :',sum(y_smote_lsm4==y_true))
print('Incorrect Predictions :',sum(y_smote_lsm4!=y_true))
print('True Positive :',auc_smote_lsm4)
print('True Negative :',accuracy_smote_lsm4)
print('False Positive :',error_smote_lsm4)
print('False Negative :',gini_smote_lsm4)
print('F1 Score :',f1_score(y_true,y_smote_lsm4))
print('Confusion matrix')
print(confusion_matrix(y_true,y_smote_lsm5))
print('Correct Predictions :',sum(y_smote_lsm5==y_true))
print('Incorrect Predictions :',sum(y_smote_lsm5!=y_true))
print('True Positive :',auc_smote_lsm5)
print('True Negative :',accuracy_smote_lsm5)
print('False Positive :',error_smote_lsm5)
print('False Negative :',gini_smote_lsm5)
print('F1 Score :',f1_score(y_true,y_smote_lsm5))
print('Confusion matrix')
print(confusion_matrix(y_true,y_smote_lsm6))
print('Correct Predictions :',sum(y_smote_lsm6==y_true))
print('Incorrect Predictions :',sum(y_smote_lsm6!=y_true))
print('True Positive :',auc_smote_lsm6)
print('True Negative :',accuracy_smote_lsm6)
print('False Positive :',error_smote_lsm6)
print('False Negative :',gini_smote_lsm6)
print('F1 Score :',f1_score(y_true,y_smote_lsm6))
print('Confusion matrix')
print(confusion_matrix(y_true,y_smote_lsm7))
print('Correct Predictions :',sum(y_smote_lsm7==y_true))
print('Incorrect Predictions :',sum(y_smote_lsm7!=y_true))
print('True Positive :',auc_smote_lsm7)
print('True Negative :',accuracy_smote_lsm7)
print('False Positive :',error_smote_lsm7)
print('False Negative :',gini_smote_lsm7)
print('F1 Score :',f1_score(y_true,y_smote_lsm7))
print('Confusion matrix')
print(confusion_matrix(y_true,y_smote_lsm8))
print('Correct Predictions :',sum(y_smote_lsm8==y_true))
print('Incorrect Predictions :',sum(y_smote_lsm8!=y_true))
print('True Positive :',auc_smote_lsm8)
print('True Negative :',accuracy_smote_lsm8)
print('False Positive :',error_smote_lsm8)
print('False Negative :',gini_smote_lsm8)
print('F1 Score :',f1_score(y_true,y_smote_lsm8))
print('Confusion matrix')
print(confusion_matrix(y_true,y_smote_lsm9))
print('Correct Predictions :',sum(y_smote_lsm9==y_true))
print('Incorrect Predictions :',sum(y_smote_lsm9!=y_true))
print('True Positive :',auc_smote_lsm9)
print('True Negative :',accuracy_smote_lsm9)
print('False Positive :',error_smote_lsm9)
print('False Negative :',gini_smote_lsm9)
print('F1 Score :',f1_score(y_true,y_smote_lsm9))
print('Confusion matrix')
print(confusion_matrix(y_true,y_smote_lsm10))
print('Correct Predictions :',sum(y_smote_lsm10==y_true))
print('Incorrect Predictions :',sum(y_smote_lsm10!=y_true))
print('True Positive :',auc_smote_lsm10)
print('True Negative :',accuracy_smote_lsm10)
print('False Positive :',error_smote_lsm10)
print('False Negative :',gini_smote_lsm10)
print('F1 Score :',f1_score(y_true,y_smote_lsm10))

```



```
accuracy_smote_lsvmL2_lss0,100-accuracy_smote_chaid_lss0,100-
accuracy_smote_c45_lss0,100-accuracy_smote_rf_lss0 ]
x.bar(algo,errorures)
plt.xlabel(&#39;Algorithm&#39;, fontweight ='bold', fontsize = 15)
plt.ylabel(&#39;Error&#39;, fontweight ='bold', fontsize = 15)
plt.title(&quot;Variation of Error for Various Classifiers using Lasso Feature selection with
Smote&quot;)
plt.show()
```

CHAPTER 4

SNAP SHOTS



Fig. 4.1. Sample Image of Execution of Source code in Jupyter Notebook

Heat Maps:

Here are the heat maps of correlation between selected features in case of Full feature selection, CFS, Wrapper FS, LASSO FS, Full features with SMOTE, LASSO with SMOTE respectively.

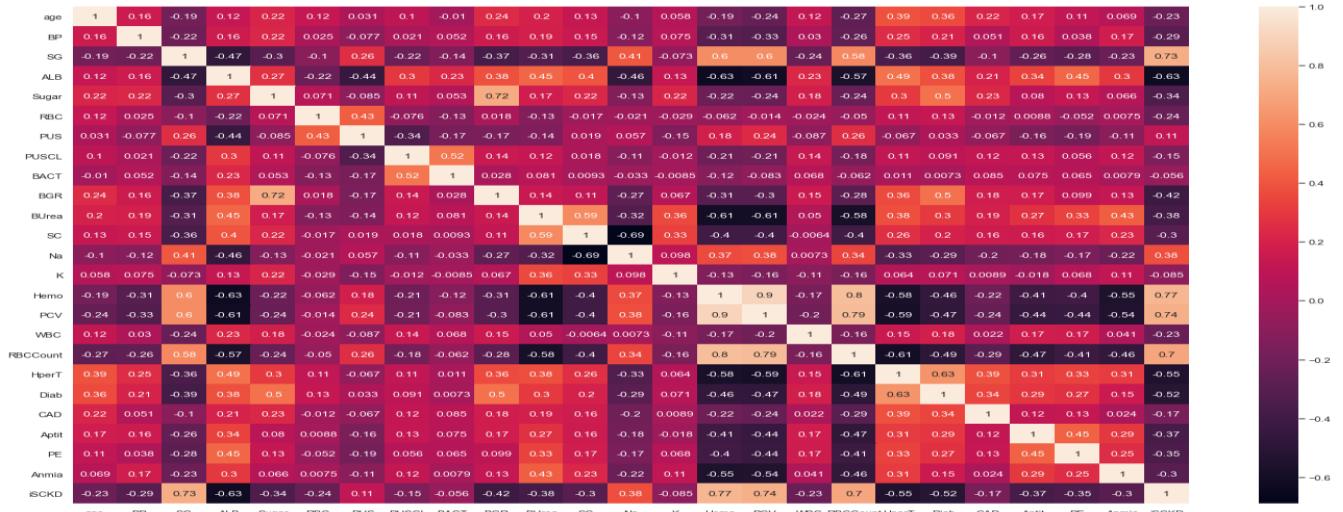


Fig. 4.2. Heat Maps of Correlation between Selected Features in case of Full feature selection



Fig. 4.3. Heat Maps of Correlation between Selected Features in case of CFS



Fig. 4.4. Heat Maps of Correlation between Selected Features in case of Wrapper FS

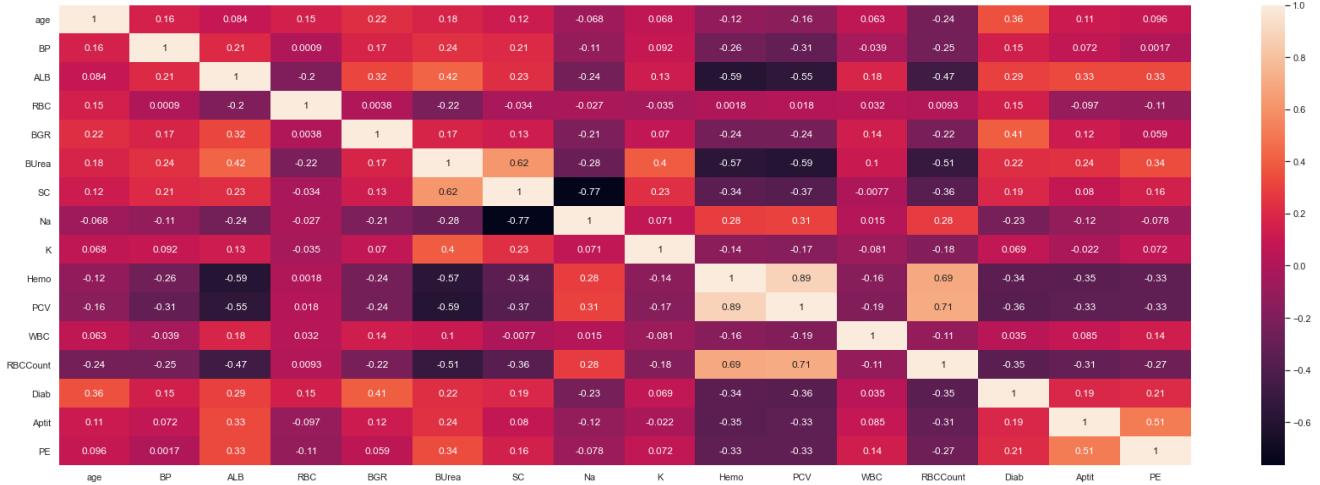


Fig. 4.5. Heat Maps of Correlation between Selected Features in case of LASSO FS



Fig. 4.6. Heat Maps of Correlation between Selected Features in case of Full features with SMOTE



Graphs of Accuracy, Precision, Recall:

Here are the plotted graphs of Accuracy, Precision, Recall respectively for the Machine Learning classifiers in case of Full feature selection, CFS, Wrapper FS, LASSO FS, Full features with SMOTE, LASSO with SMOTE respectively.

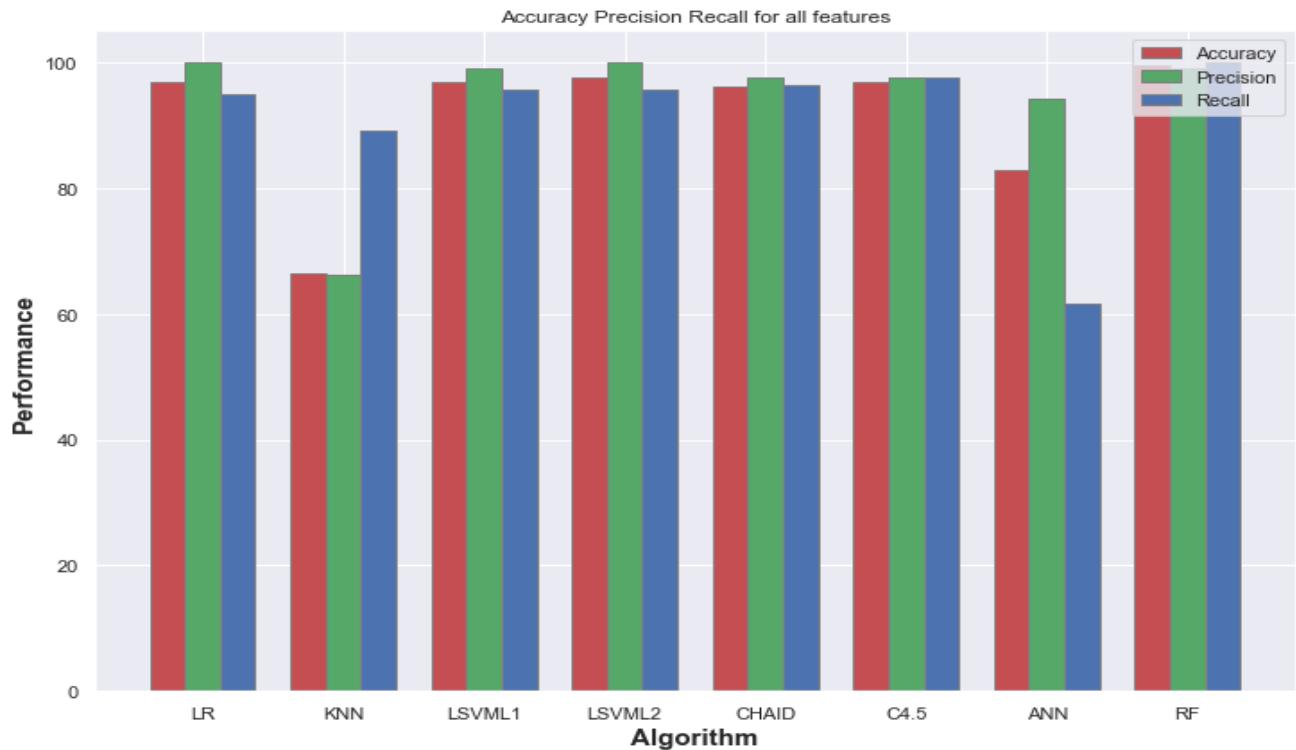


Fig. 4.8. Accuracy, Precision, Recall for ML classifiers in case of Full feature selection

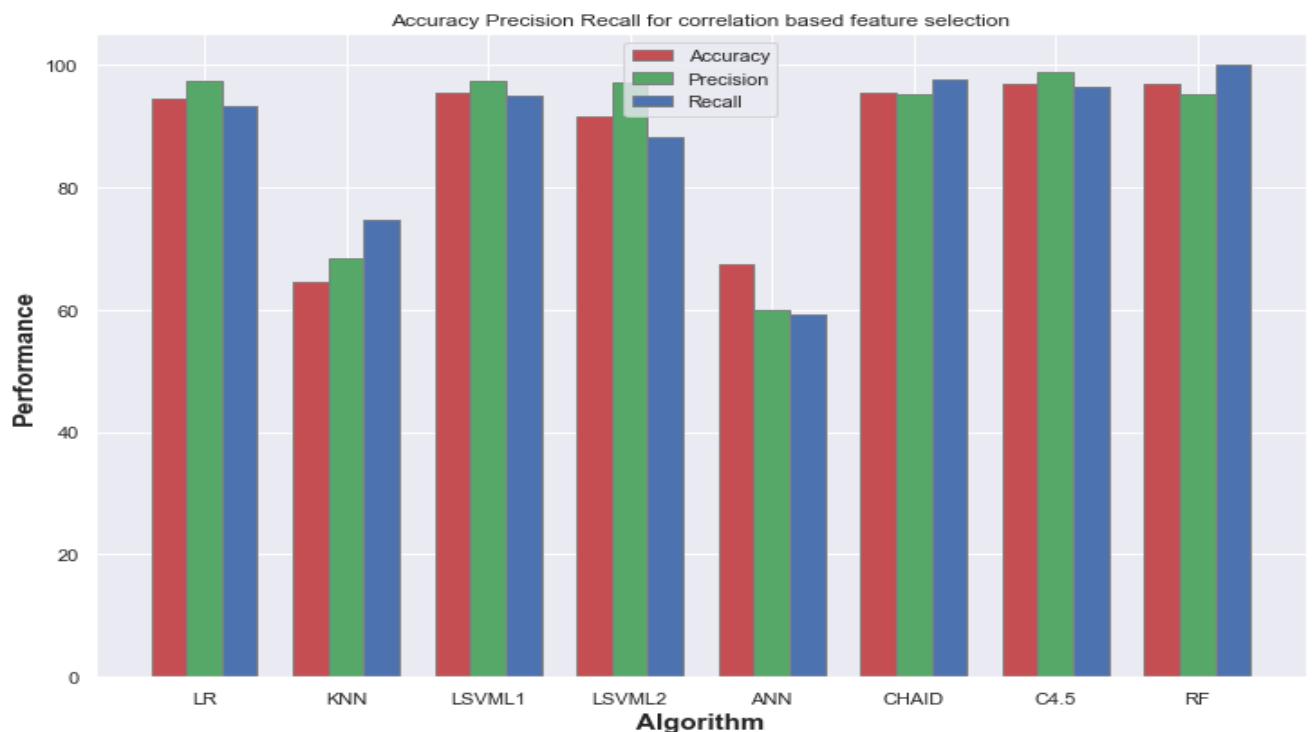


Fig. 4.9. Accuracy, Precision, Recall for ML classifiers in case of CFS

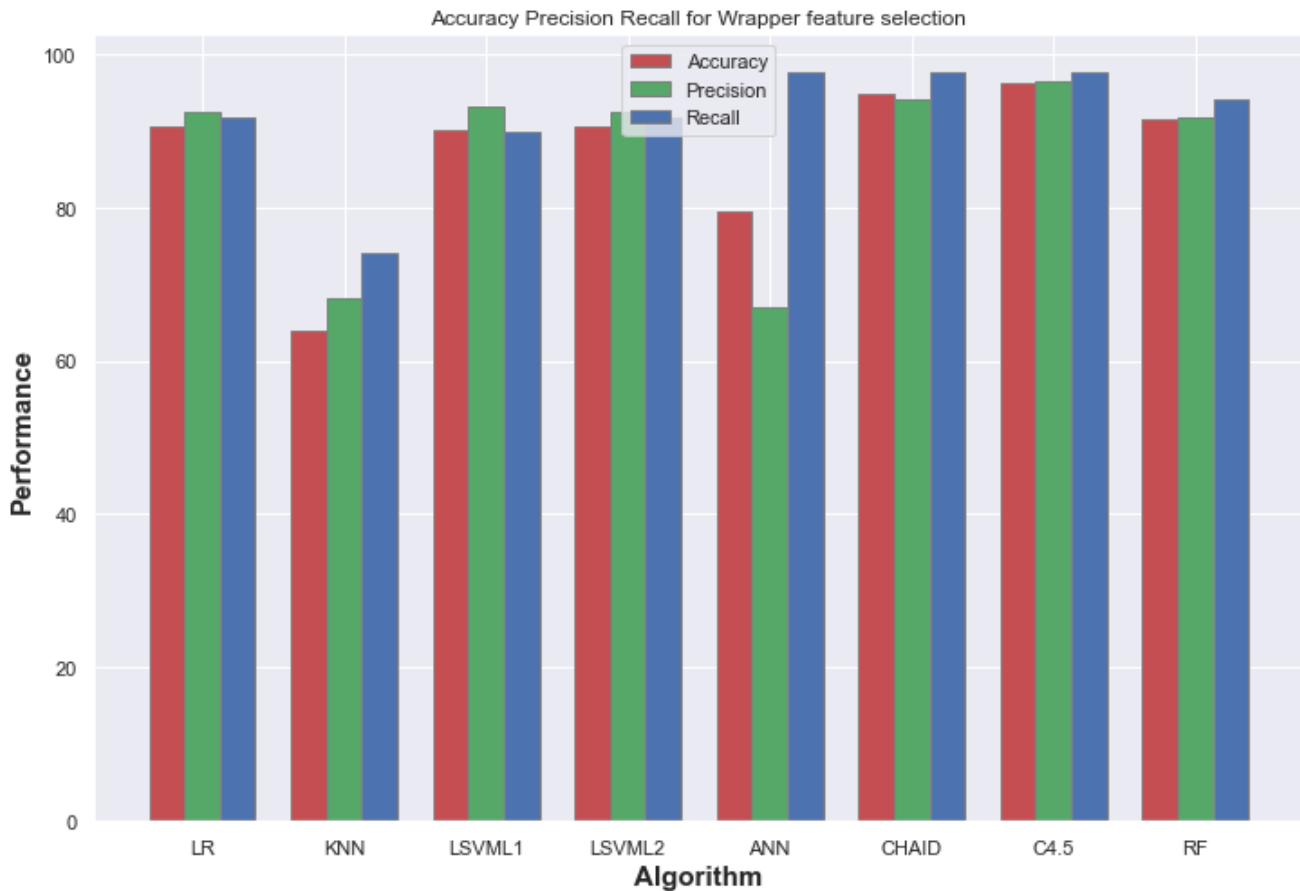


Fig. 4.10. Accuracy, Precision, Recall for ML classifiers in case of Wrapper FS

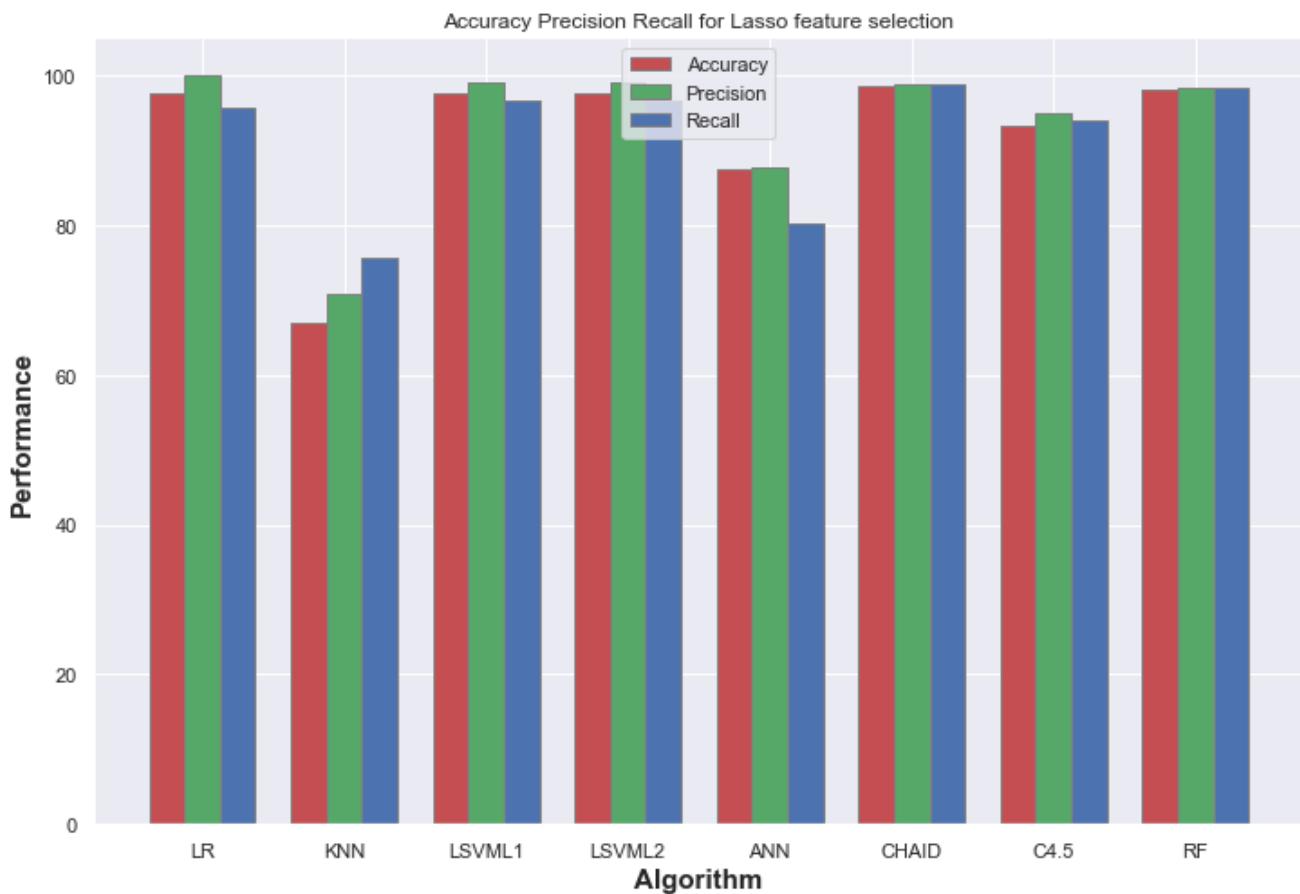


Fig. 4.11. Accuracy, Precision, Recall for ML classifiers in case of LASSO FS

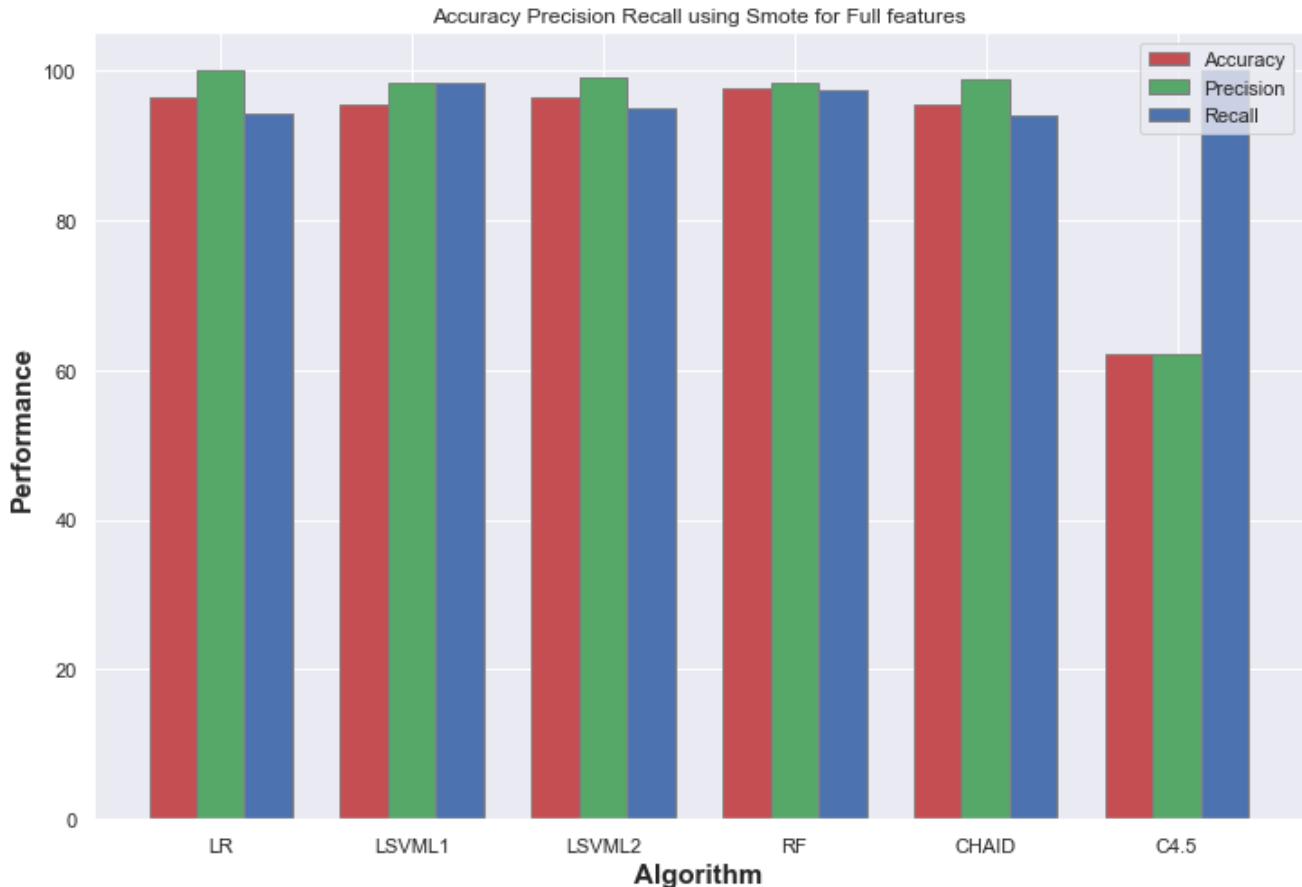


Fig. 4.12. Accuracy, Precision, Recall for ML classifiers in case of Full features with SMOTE

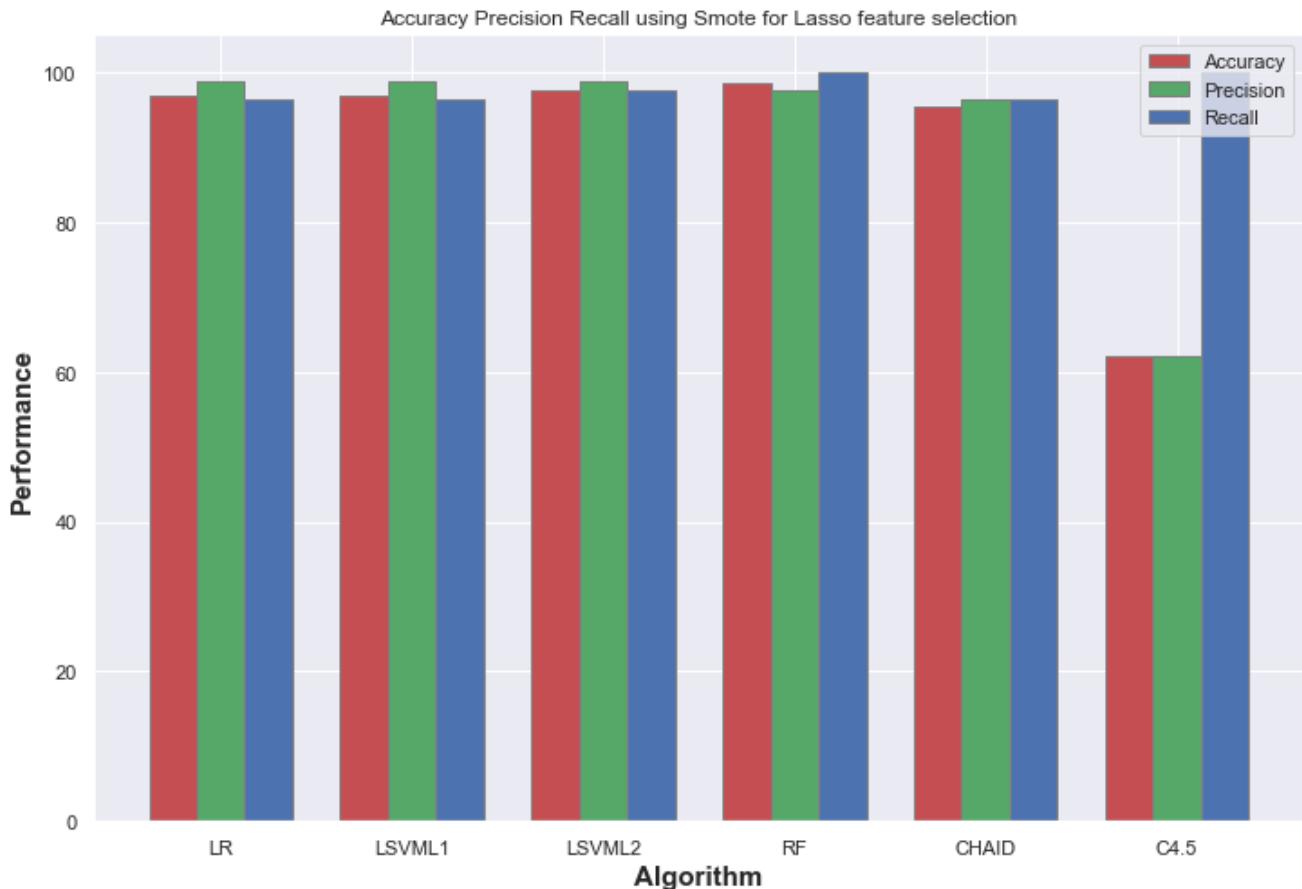


Fig. 4.13. Accuracy, Precision, Recall for ML classifiers in case of LASSO FS with SMOTE
Classification Error:

Here are the plotted graphs of Classification Error for the Machine Learning classifiers in case of Full feature selection, CFS, Wrapper FS, LASSO FS, Full features with SMOTE, LASSO with SMOTE respectively.

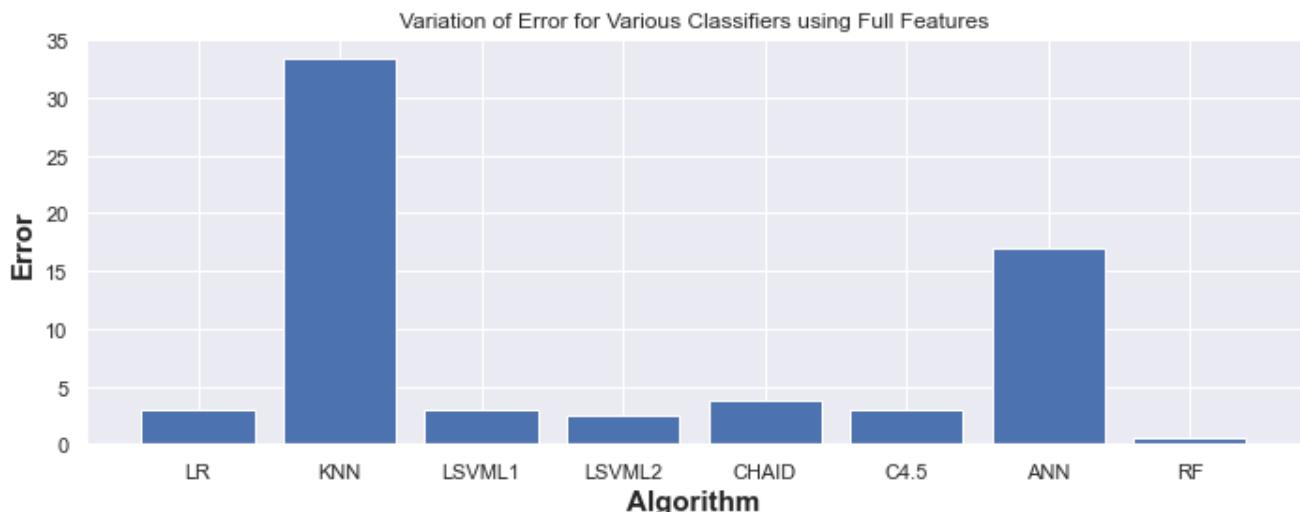


Fig. 4.14. Classification Error for ML classifiers in case of Full feature selection

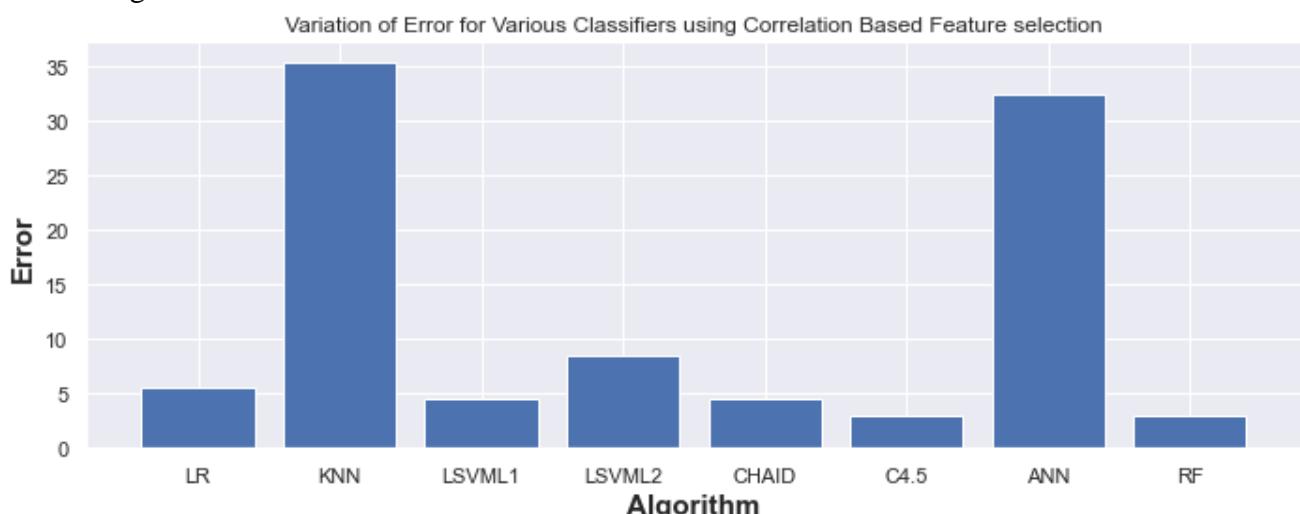


Fig. 4.15. Classification Error for ML classifiers in case of CFS

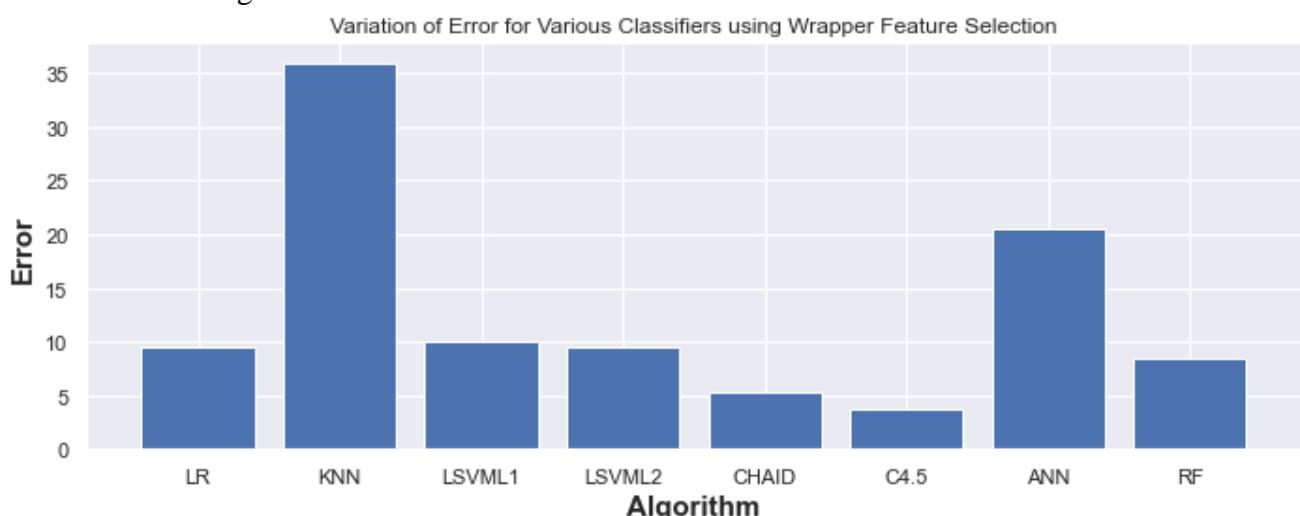


Fig. 4.16. Classification Error for ML classifiers in case of Wrapper FS

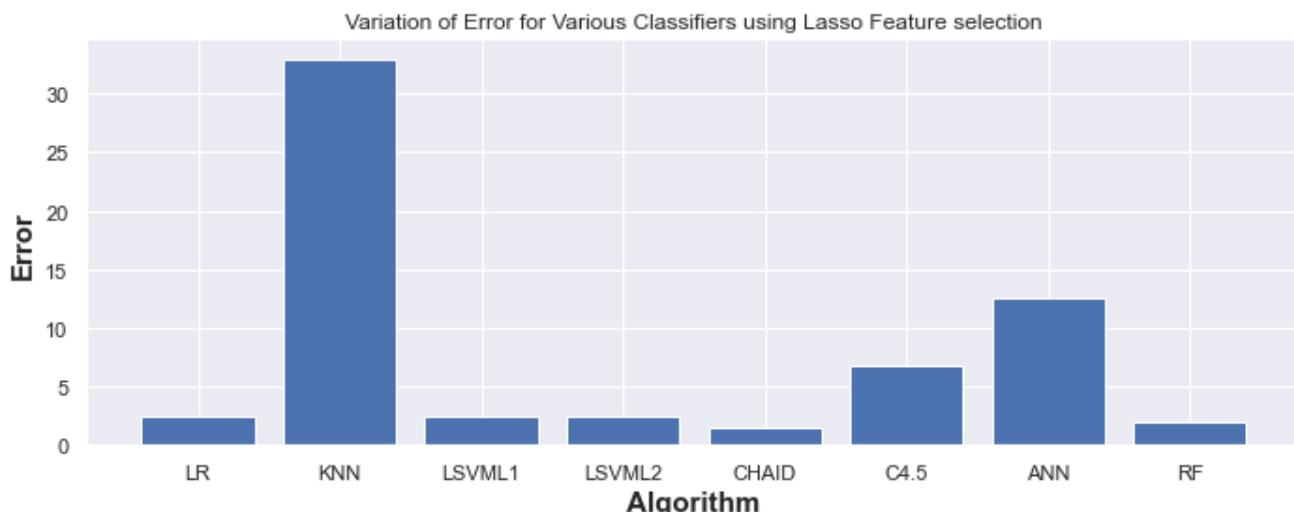


Fig. 4.17. Classification Error for ML classifiers in case of LASSO FS

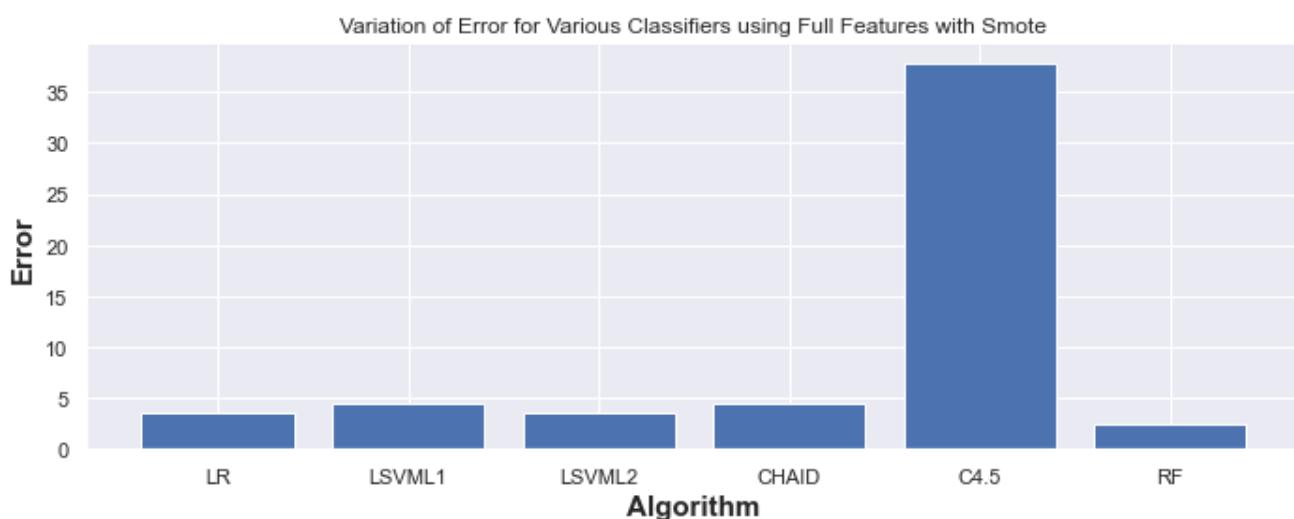


Fig. 4.18. Classification Error for ML classifiers in case of Full features with SMOTE

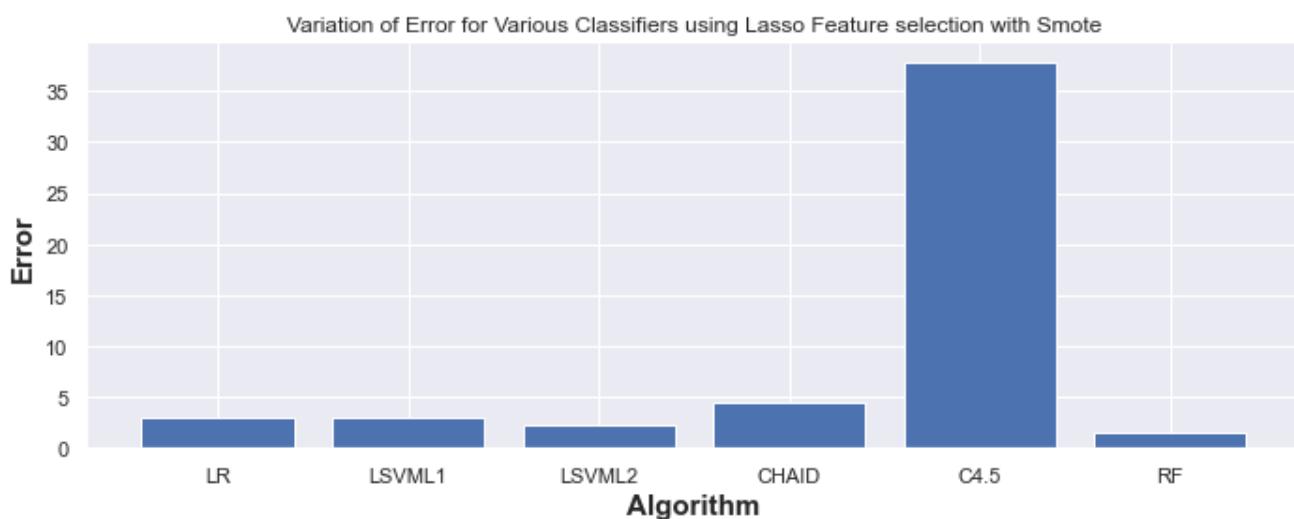


Fig. 4.19. Classification Error for ML classifiers in case of LASSO FS with SMOTE

F- Measure:

Here are the plotted graphs of F- Measure for the Machine Learning classifiers in case of Full feature selection, CFS, Wrapper FS, LASSO FS, Full features with SMOTE, LASSO with SMOTE respectively.

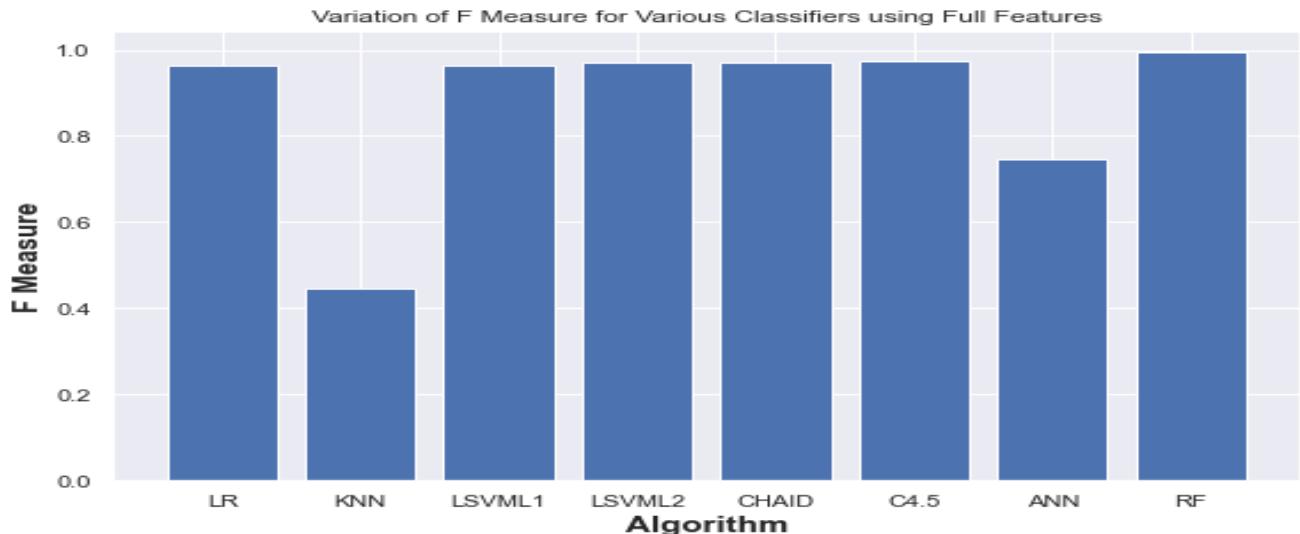


Fig. 4.20. F- Measure for ML classifiers in case of Full feature selection

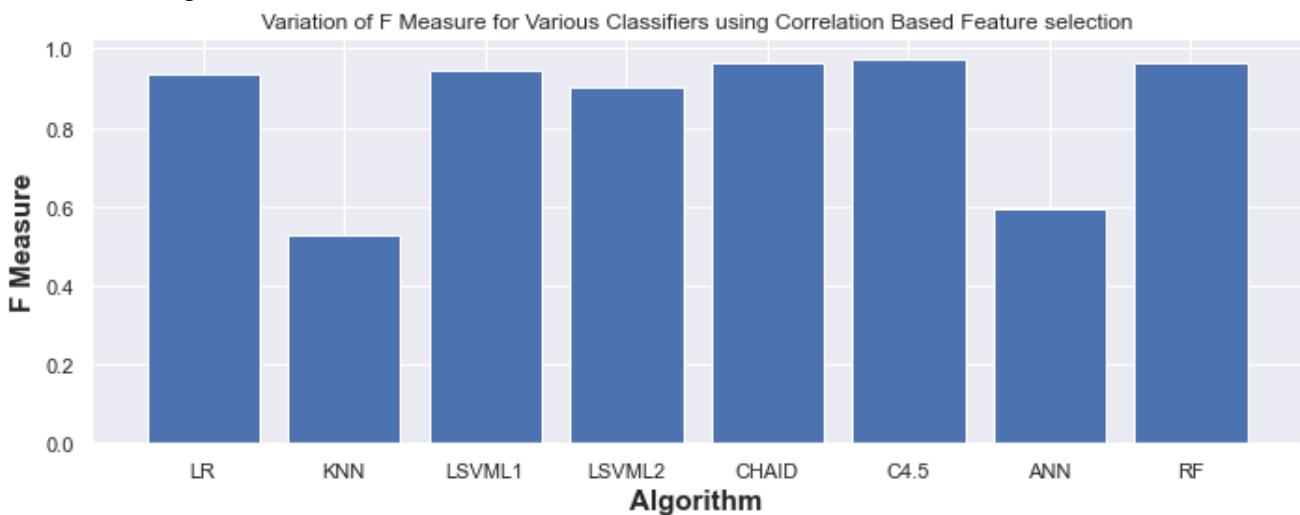


Fig. 4.21. F- Measure for ML classifiers in case of CFS

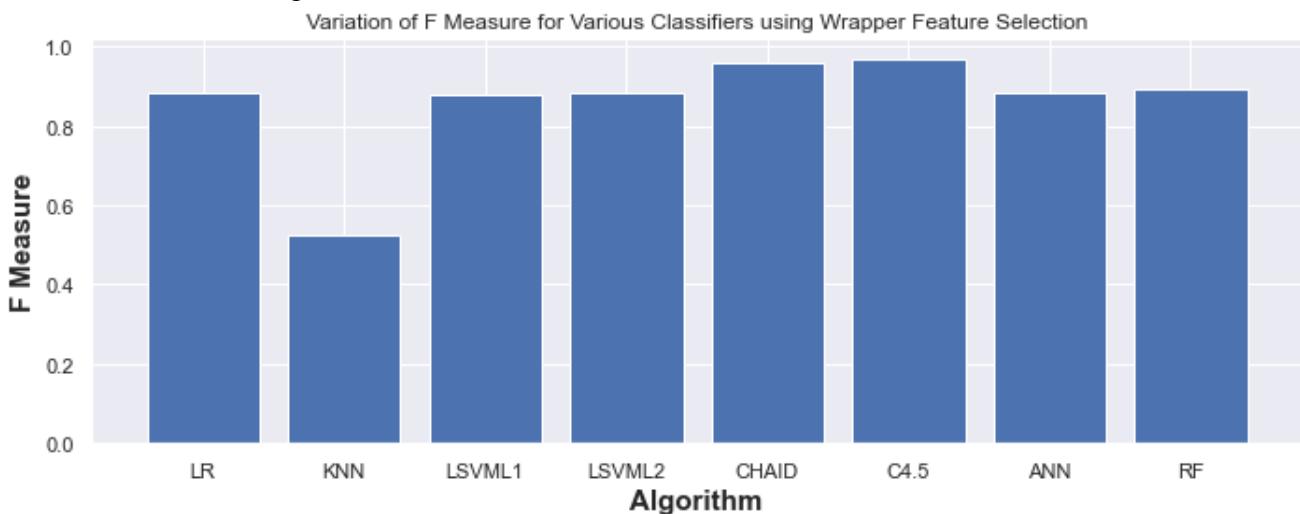


Fig. 4.22. F- Measure for ML classifiers in case of Wrapper FS

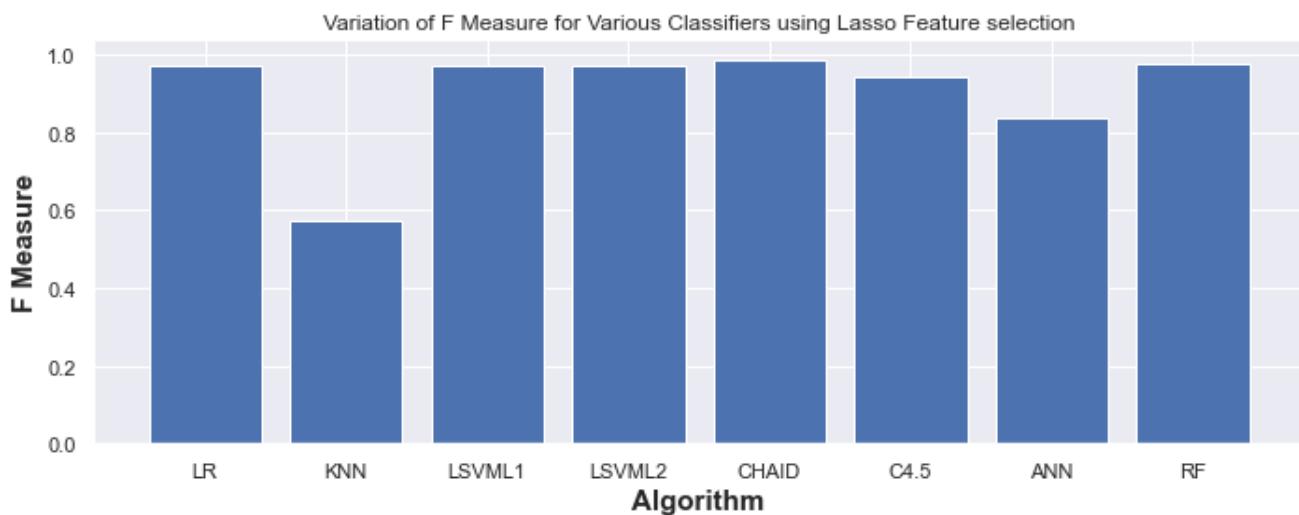


Fig. 4.23. F- Measure for ML classifiers in case of LASSO FS

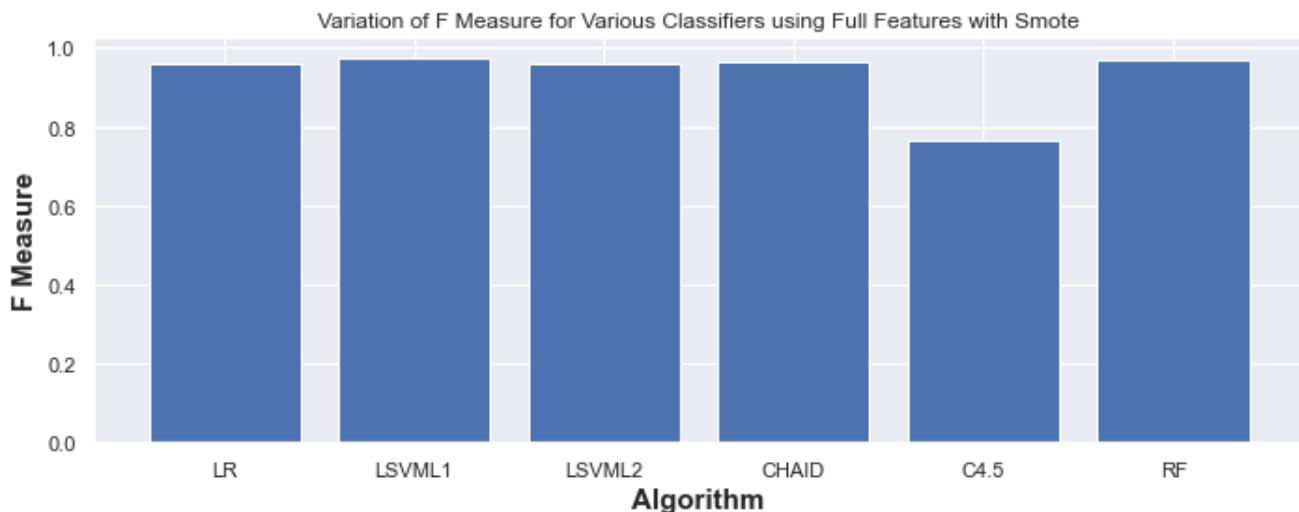


Fig. 4.24. F- Measure for ML classifiers in case of Full features with SMOTE

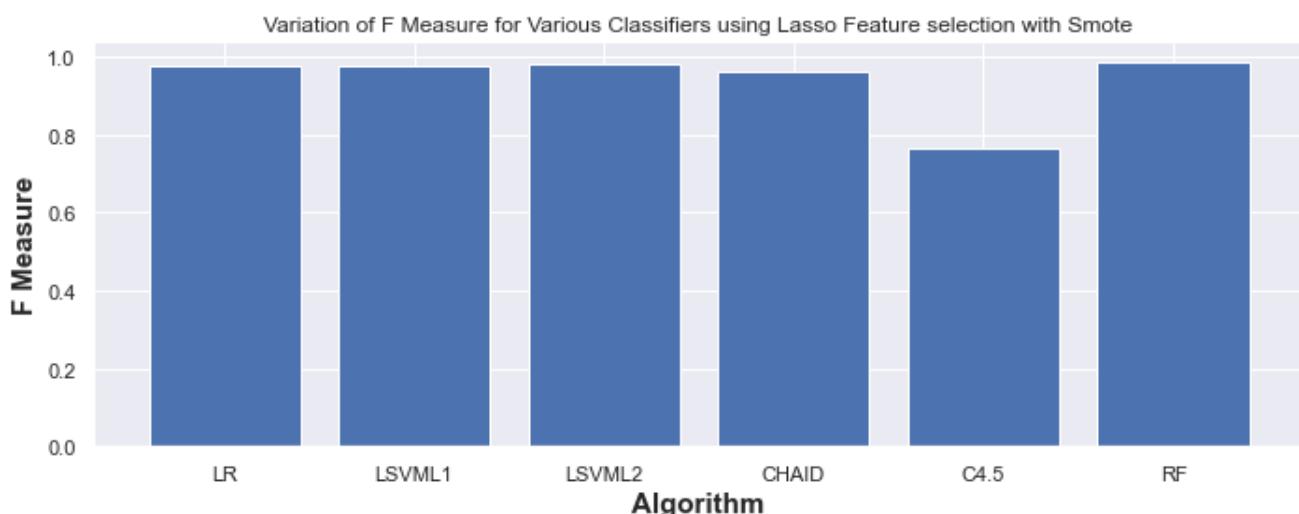


Fig. 4.25. F- Measure for ML classifiers in case of LASSO FS with SMOTE

AUC:

Here are the plotted graphs of AUC for the Machine Learning classifiers in case of Full feature selection, CFS, Wrapper FS, LASSO FS, Full features with SMOTE, LASSO with SMOTE respectively.

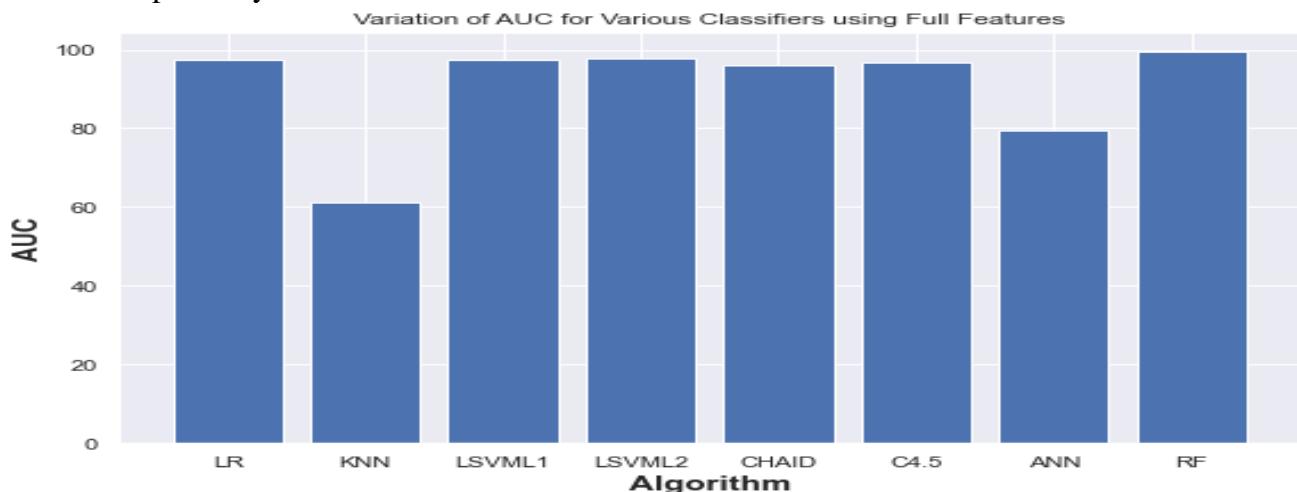


Fig. 4.26. AUC for ML classifiers in case of Full feature selection

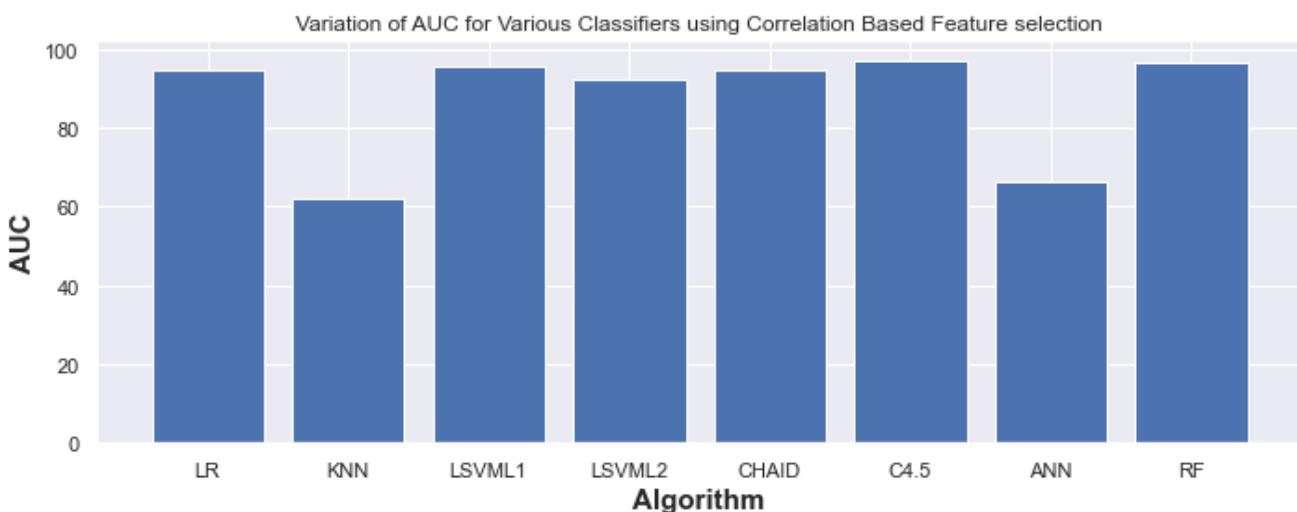


Fig. 4.27. AUC for ML classifiers in case of CFS

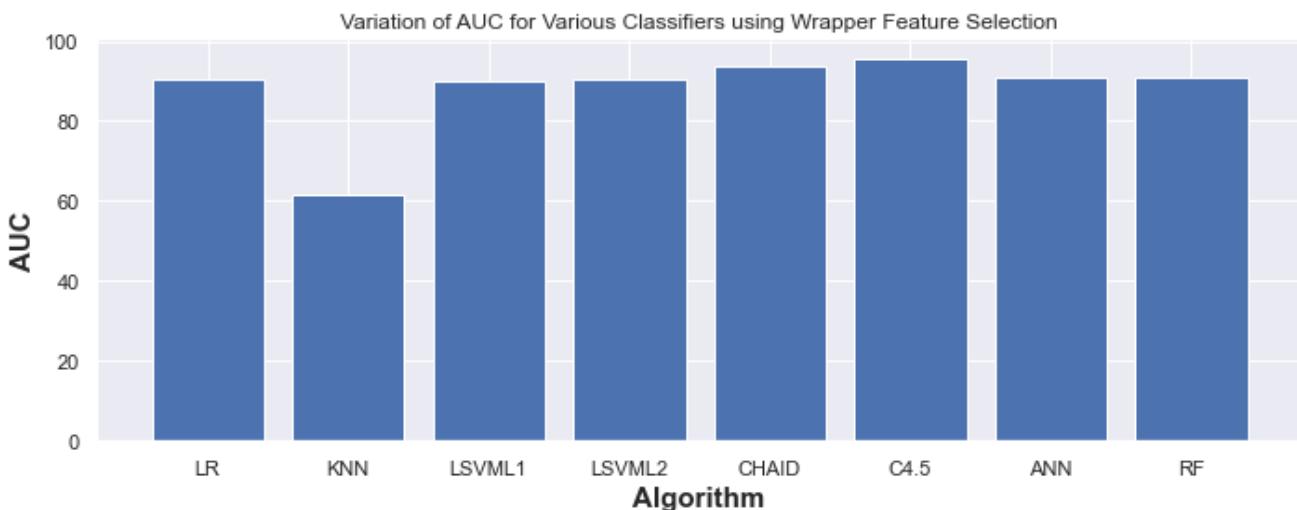


Fig. 4.28. AUC for ML classifiers in case of Wrapper FS

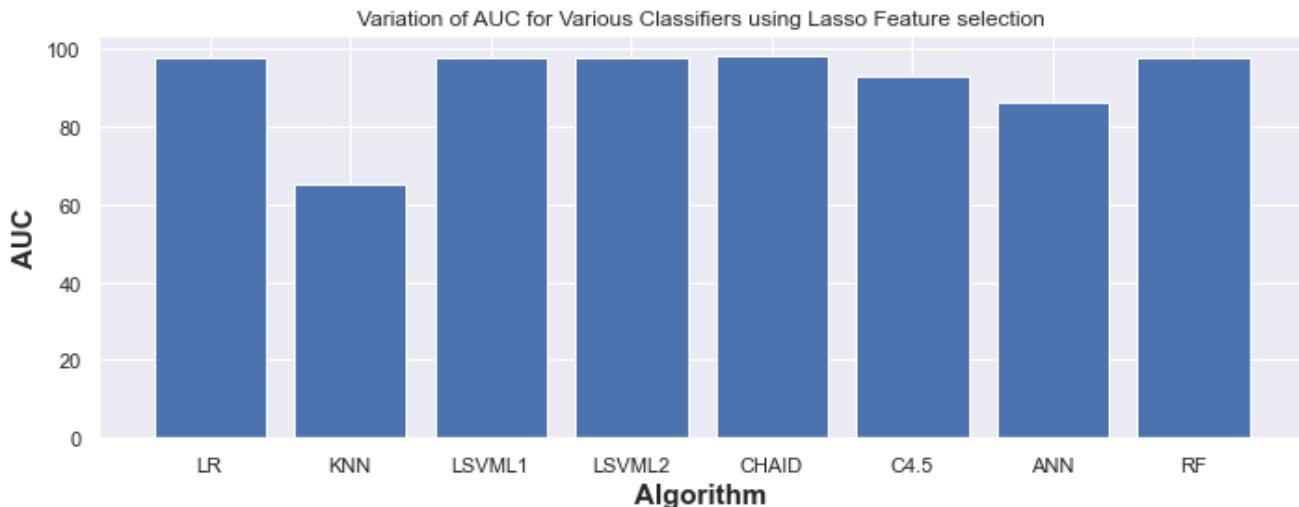


Fig. 4.29. AUC for ML classifiers in case of LASSO FS

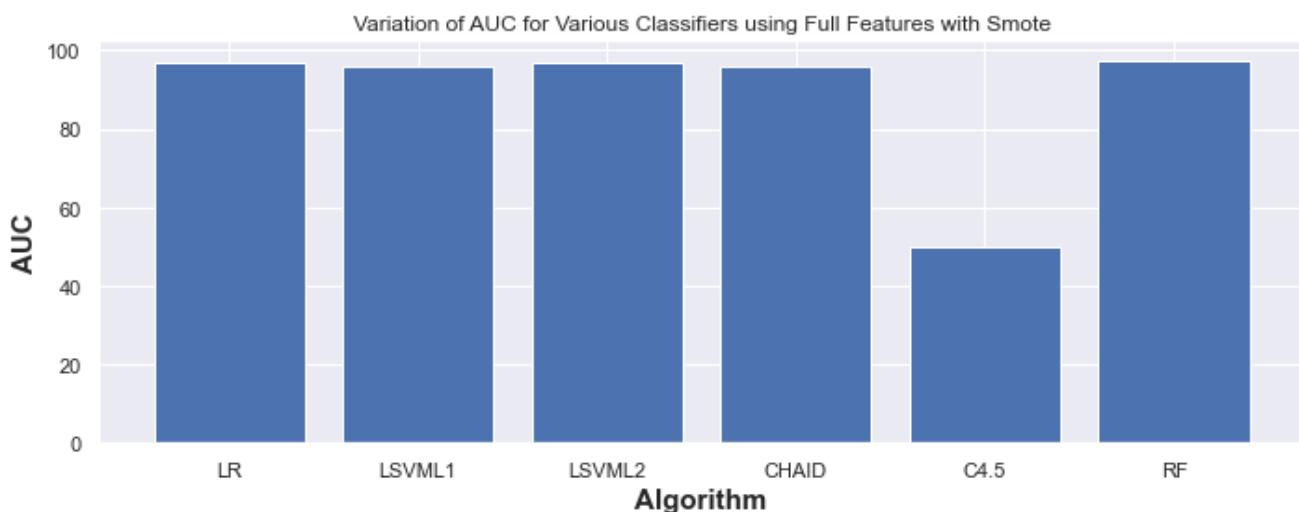


Fig. 4.30. AUC for ML classifiers in case of Full features with SMOTE

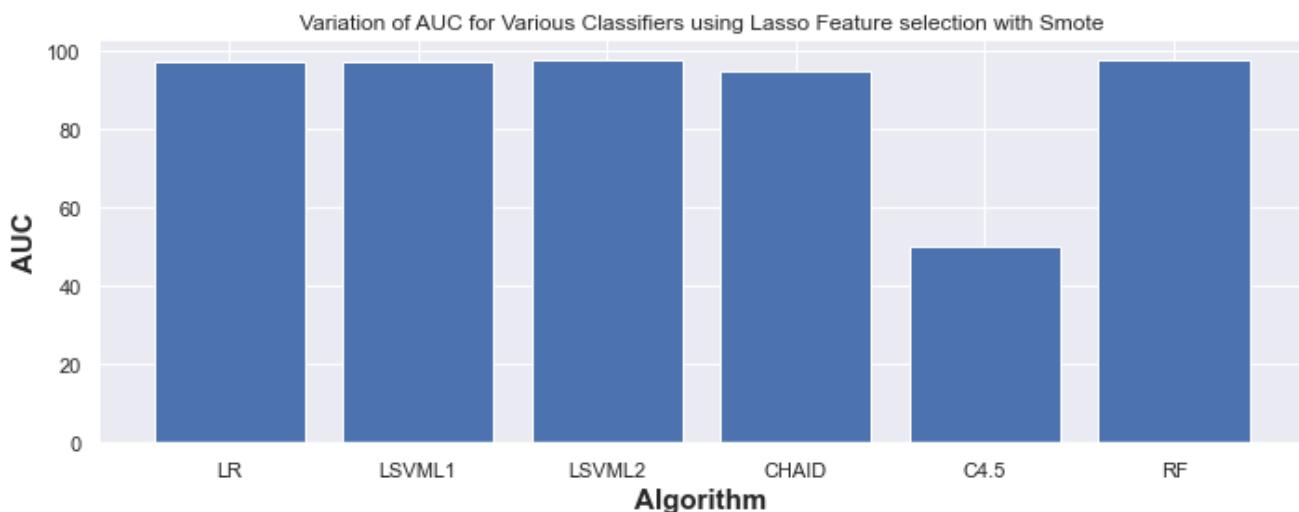


Fig. 4.31. AUC for ML classifiers in case of LASSO FS with SMOTE

GINI Index:

Here are the plotted graphs of GINI Index for the Machine Learning classifiers in case of Full feature selection, CFS, Wrapper FS, LASSO FS, Full features with SMOTE, LASSO with SMOTE respectively.

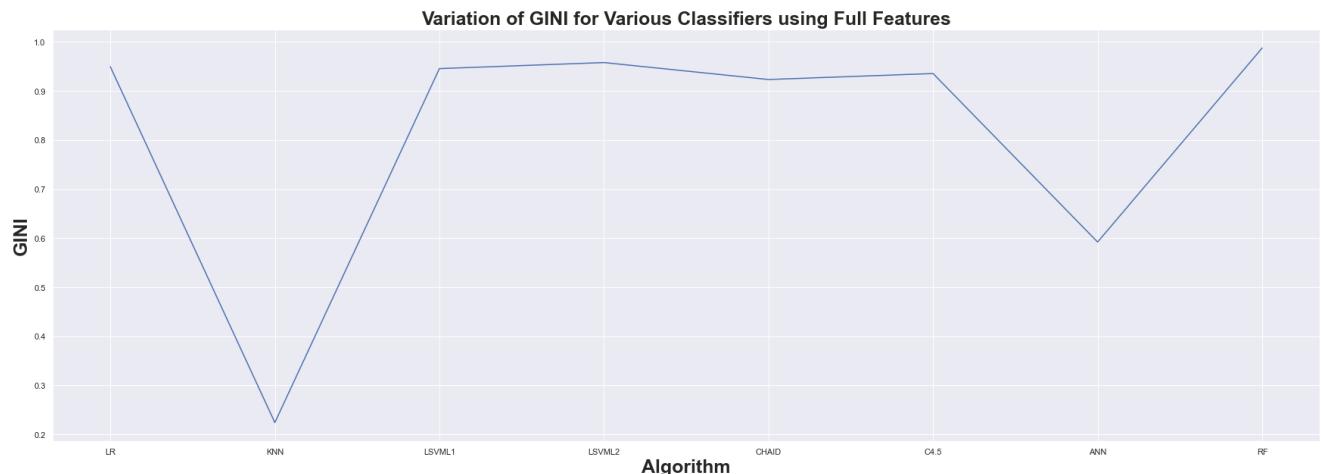


Fig. 4.32. GINI Index for ML classifiers in case of Full feature selection

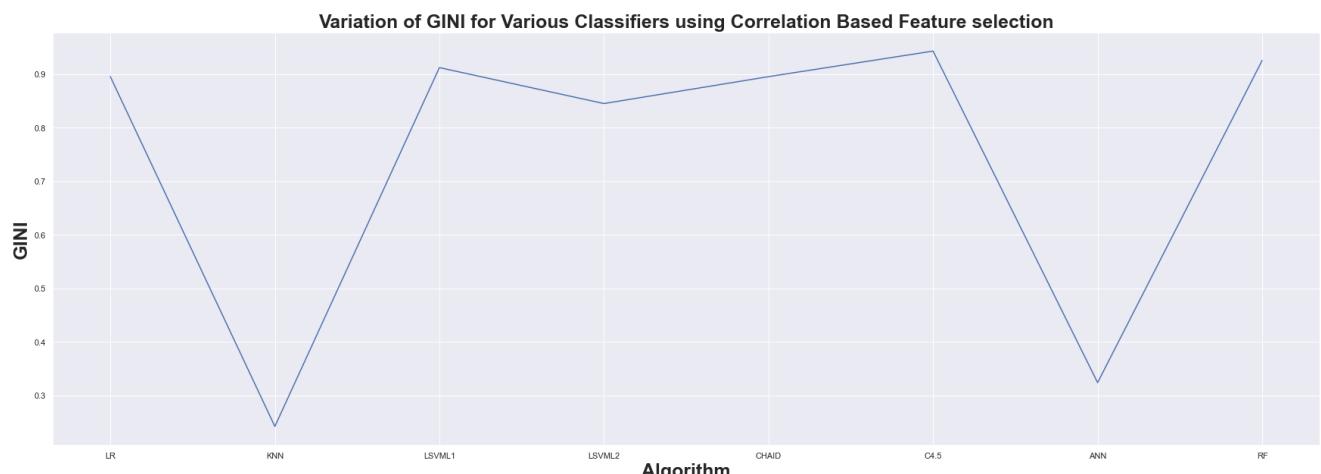


Fig. 4.33. GINI Index for ML classifiers in case of CFS

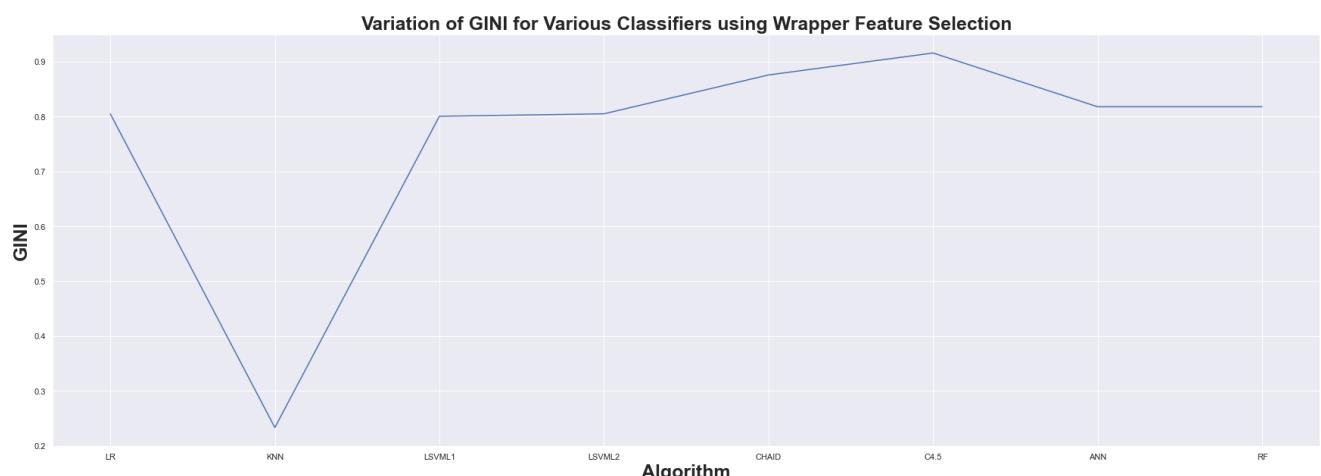


Fig. 4.34. GINI Index for ML classifiers in case of Wrapper FS

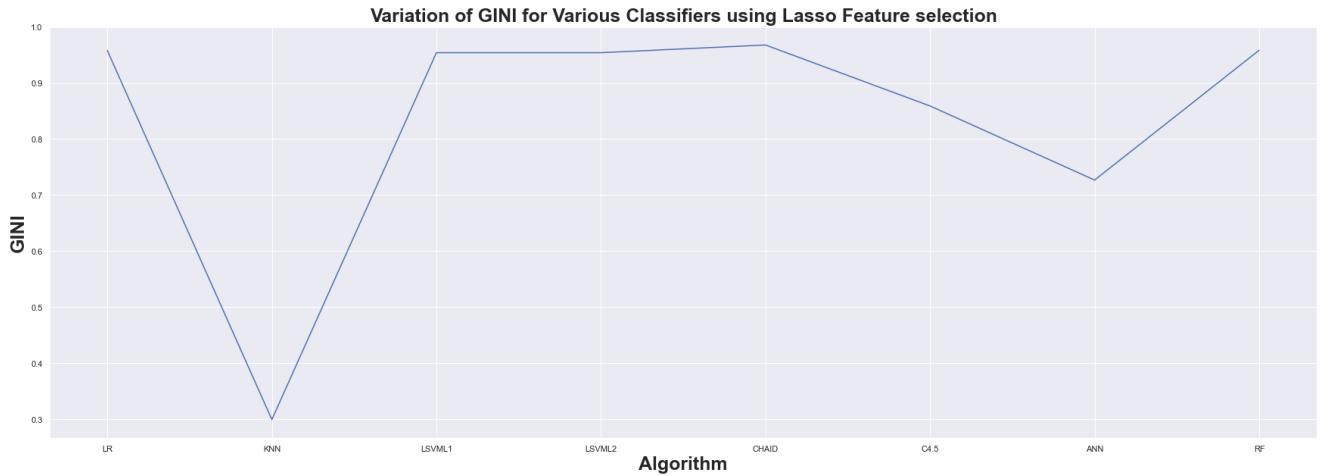


Fig. 4.35. GINI Index for ML classifiers in case of LASSO FS

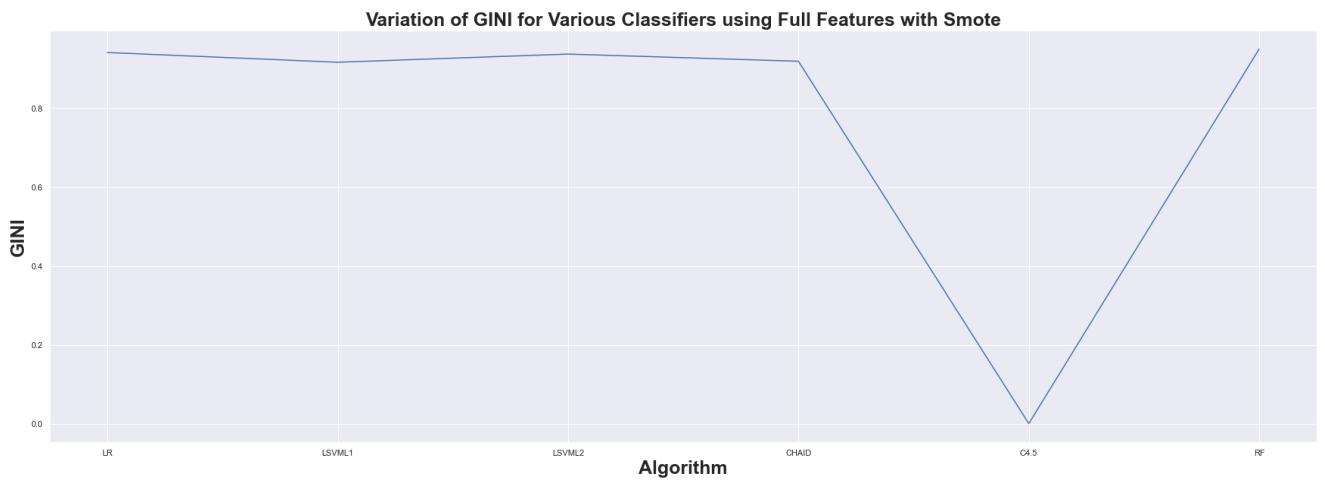


Fig. 4.36. GINI Index for ML classifiers in case of Full features with SMOTE

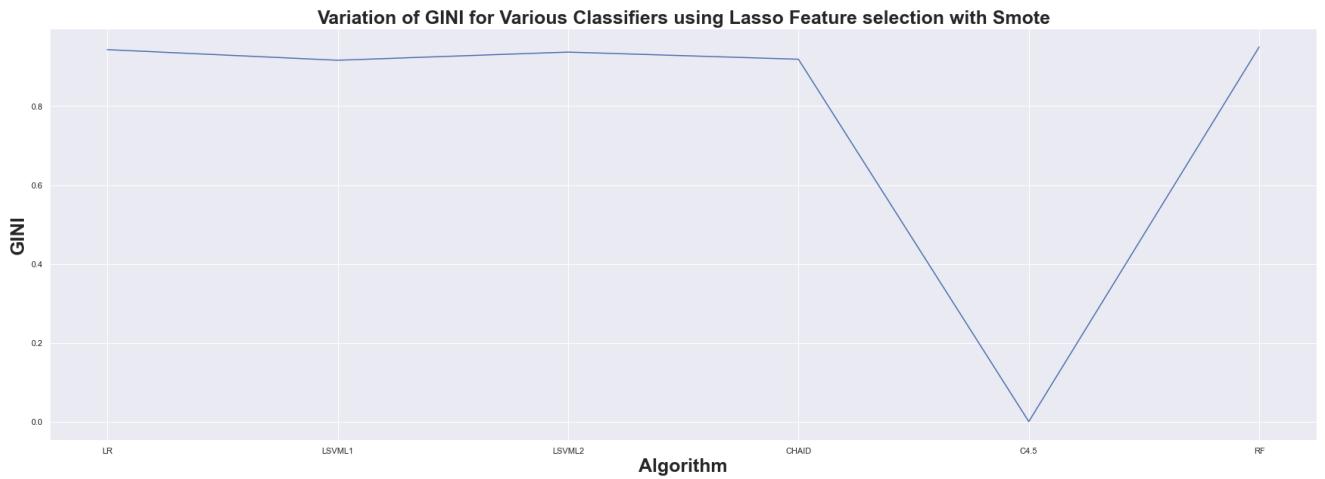


Fig. 4.37. GINI Index for ML classifiers in case of LASSO FS with SMOTE

CHAPTER 5

CONCLUSION AND FUTURE PLANS

5.1. Conclusion

The main topic of this research is the prediction of CKD using ML classifier techniques. We employed ML classifier methods as RF, CHAID, LSVM with penalty L1 and penalty L2, LR, KNN, LSVM with penalty L2, C 4.5, and ANN. For each classifier, the results for full feature, correlation-based, wrapper, and LASSO FS are computed.

LR, LSVM L1, LSVM L2, CHAID, C 4.5, and RF outperform Full Features and LASSO Feature Selection in terms of accuracy and other parameters.

SMOTE is used to calculate the outcomes of the classifiers and FS methods stated above. We did not employ the KNN and ANN algorithms for SMOTE since they did not yield results that were sufficient.

Random Forest outperforms LSVM penalty L2 with an accuracy of 98.49 percent and LASSO Feature Selection and SMOTE with an accuracy of 97.73 percent.

The findings led to the conclusion that Random Forest with LASSO Feature Selection with SMOTE yields better outcomes for predicting CKD.

5.2. Future Plans

We oversampled our data set with SMOTE because it had an average amount of instances. The same method can be used to perform undersampling in the future when a data collection with many instances is being considered.

CHAPTER 6

REFERENCES

- [1] Pankaj Chittora, Sandeep Chaurasia, Prasun Chakrabarti, Gaurav Kumawat1, Tulika Chakrabarti, Zbigniew Leonowicz, Michał Jasiński, Lukasz Jasiński, Radomir Gono, Elżbieta Jasińska, and Vadim Bolshev , (2021) “prediction of chronic kidney disease - a machine learning perspective”, *ieee access*.2021.3053763.
- [2] J. Qin, L. Chen, Y. Liu, C. Liu, C. Feng, and B. Chen, “A machine learning methodology for diagnosing chronic kidney disease,” *IEEE Access*, vol. 8, pp. 20991–21002, 2020.
- [3] L. Kilvia De Almeida, L. Lessa, A. Peixoto, R. Gomes, and J. Celestino, “Kidney failure detection using machine learning techniques,” in *Proc. 8th Int. Workshop ADVANCEs ICT Infrastructures Services*, 2020, pp. 1–8.
- [4] S. Shankar, S. Verma, S. Elavarthy, T. Kiran, and P. Ghuli, “Analysis and prediction of chronic kidney disease,” *Int. Res. J. Eng. Technol.*, vol. 7, no. 5, May 2020, pp. 4536–4541.
- [5] G. R. Vasquez-Morales, S. M. Martinez-Monterrubio, P. Moreno-Ger, and J. A. Recio-Garcia, “Explainable prediction of chronic renal disease in the colombian population using neural networks and case-based reasoning,” *IEEE Access*, vol. 7, pp. 152900–152910, 2019.
- [6] M. Almasoud and T. E. Ward, “Detection of chronic kidney disease using machine learning algorithms with least number of predictors,” *Int. J. Adv. Comput. Sci. Appl.*, vol. 10, no. 8, pp. 89–96, 2019.
- [7] Y. Amirkaliyev, S. Shamiluulu, and A. Serek, “Analysis of chronic kidney disease dataset by applying machine learning methods,” in *Proc. IEEE 12th Int. Conf. Appl. Inf. Commun. Technol. (AICT)*, Oct. 2018, pp. 1–4.
- [8] W. Gunarathne, K. D. M Perera, and K. A. D. C. P Kahandawaarachchi, “Performance evaluation on machine learning classification techniques for disease classification and forecasting through data analytics for chronic kidney disease (CKD),” in *Proc. IEEE 17th Int. Conf. Bioinf. Bioeng. (BIBE)*, Oct.2017, pp.291–296.
- [9] K. R. A. Padmanaban and G. Parthiban, “Applying machine learning techniques for predicting the risk of chronic kidney disease,” *Indian J. Sci. Technol.*, vol. 9, no. 29, Aug. 2016.

CHAPTER 6

APPENDIX -BASE PAPER

The screenshot shows a search result for a paper on the IEEE Xplore platform. The paper is titled "Prediction of Chronic Kidney Disease - A Machine Learning Perspective". It is published by IEEE and has 3530 full-text views. The abstract discusses the use of machine learning techniques for diagnosing Chronic Kidney Disease. The page includes sections for Document Sections, Authors, Figures, References, Citations, and Keywords. To the right, there is a sidebar for "Need Full-Text" with a "CONTACT IEEE TO SUBSCRIBE" button, and a "More Like This" section with related papers.

Fig. 6.1. Screenshot of Indexing

PRIMARY SOURCES

- | | | |
|---|---|----|
| 1 | Pankaj Chittora, Sandeep Chaurasia, Prasun Chakrabarti, Gaurav Kumawat et al.
"Prediction of Chronic Kidney Disease - A Machine Learning Perspective", IEEE Access, 2021
Publication | 2% |
| 2 | Olayinka Ayodele Jongbo, Adebayo Olusola Adetunmbi, Roseline Bosede Ogunrinde, Bukola Badeji-Ajisafe. "Development of an Ensemble Approach to Chronic Kidney Disease Diagnosis", Scientific African, 2020
Publication | 1% |
| 3 | Samrat Kumar Dey, Khandaker Mohammad Mohi Uddin, Hafiz Md. Hasan Babu, Md. Mahbubur Rahman, Arpita Howlader, K.M. Aslam Uddin. "Chi2-MI: A hybrid feature selection based machine learning approach in diagnosis of chronic kidney disease", Intelligent Systems with Applications, 2022
Publication | 1% |
-

- 4 "Decision Sciences for COVID-19", Springer Science and Business Media LLC, 2022 Publication 1 %
- 5 Novia Arinda Pradisty, A. Aldrie Amir, Martin Zimmer. "Plant species- and stage-specific differences in microbial decay of mangrove leaf litter: the older the better?", Oecologia, 2021 Publication 1 %
- 6 core.ac.uk Internet Source 1 %
- 7 Wendy D. Fife. "Comparison Between Malignant and Nonmalignant Splenic Masses in Dogs using Contrast-Enhanced Computed Tomography", Veterinary Radiology & Ultrasound, 7/2004 Publication <1 %
- 8 www.biorxiv.org Internet Source <1 %
- 9 www.imedpub.com Internet Source <1 %
- 10 nceca.in Internet Source <1 %
- 11 ebin.pub Internet Source <1 %

- 12 W.H.S.D Gunarathne, K.D.M Perera, K.A.D.C.P Kahandawaarachchi. "Performance Evaluation on Machine Learning Classification Techniques for Disease Classification and Forecasting through Data Analytics for Chronic Kidney Disease (CKD)", 2017 IEEE 17th International Conference on Bioinformatics and Bioengineering (BIBE), 2017
Publication <1 %
-
- 13 Yan Chai Hum. "Segmentation of Hand Bone for Bone Age Assessment", Springer Science and Business Media LLC, 2013
Publication <1 %
-
- 14 conf.kln.ac.lk <1 %
Internet Source
-
- 15 Muhammet Sinan Basarslan, Fatih Kayaalp. "Performance Analysis Of Fuzzy Rough Set-Based And Correlation-Based Attribute Selection Methods On Detection Of Chronic Kidney Disease With Various Classifiers", 2019 Scientific Meeting on Electrical-Electronics & Biomedical Engineering and Computer Science (EBBT), 2019
Publication <1 %
-
- 16 d.researchbib.com <1 %
Internet Source
-
- 17 vbn.aau.dk <1 %
Internet Source

<1 %

-
- 18 Nikhila. "Chronic Kidney Disease Prediction using Machine Learning Ensemble Algorithm", 2021 International Conference on Computing, Communication, and Intelligent Systems (ICCCIS), 2021

Publication

Exclude quotes Off

Exclude bibliography Off

Exclude matches Off