

# Comparative Analysis for the Prediction of major Diseases using Supervised and Hybrid Machine Learning Algorithms

Dillip Narayan Sahu<sup>1</sup>, Vijay Pal Singh<sup>2\*</sup>

<sup>1</sup>Research Scholar, Department of Computer Science, OPJS University, Rajasthan, India

<sup>2</sup>Associate Professor, Department of Computer Science, OPJS University, Rajasthan, India

**Abstract:** In the few recent years, common major diseases have emerged together of the foremost common causes of deaths worldwide. As the changing in lifestyle, food habits, working cultures etc, has significantly contributed to the present issues related to health across world-wide including the developed, underdeveloped and developing countries, a challenge to the medical science to overcome from this situation[1]. As per as WHO (World Health Organization) Global Health Estimates report is concerned, an estimated 74% of all deaths were noncommunicable diseases globally, 3 out of 10 major diseases are communicable. In this paper, we have taken 5 major diseases (Ischaemic Heart Disease (IHD) with Stroke, Chronic Kidney Disease (CKD), Diabetes Mellitus (DM) including BP, Chronic Liver, and Cancer) among the top 10 deadliest diseases[2]. All these major diseases can be curable with proper diagnosis and early detection. The purpose of this paper is to establish some Machine Learning supervised algorithms with some hybrid approach for better comparative analysis and predict for the particular disease at an early stage with a greater accuracy level. The outcome of this paper also justify that the hybrid algorithm model has better processing, performance with more accuracy outputs so as to help the medical and healthcare sector in the early stage disease prediction.

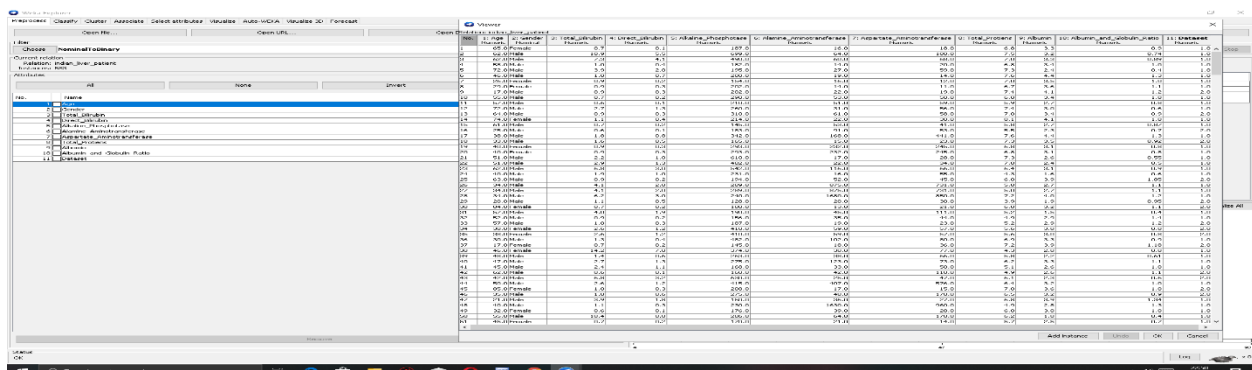
**Keywords:** Algorithm, Classifier, Machine Learning, Predictive Analysis.

## INTRODUCTION

The use of machine learning with deep learning in medical and health care sector is increasing day by day. The Machine Learning is a powerful field of Artificial Intelligence where we can do a better analysis, processing, visualizing, cleaning, classification, clustering, etc. with the help of different machine learning algorithms which are exist as well as some hybrid methods to predict and to obtain better results [3][4]. In this paper, we have taken 5 different dataset based on 5 different major diseases as per the report of World Health Organization (WHO). We investigated a total of 5 different major diseases based on the WHO report on the major diseases with different experiments and observations by the help of Weka data analysis for machine learning tool.

People of many developing countries are still living in poorness with dysfunctional health care systems and very restricted access to basic treatment [5]. The information needed for the analysis of health state of affairs and health issues comprise of mortality, morbidity, demographic condition, socioeconomic factors etc [6].

Following are the number of the main diseases in today's world that cause health issues and deaths in several countries in the world [7]. The five major diseases are- Cardio Vascular Disease (Heart), Chronic Kidney Disease, Chronic Liver Disease, Cancer Disease, Diabetes Mellitus.



Age	Sex	BMI	Glucose	Blood Pressure	Heart Rate	Cholesterol	Triglycerides	HDL	LDL	Liver Disease
39	M	29.2	100	130	75	200	100	50	150	0
41	F	26.7	110	120	80	180	90	40	140	0
35	M	31.5	120	140	90	220	110	60	160	0
45	F	28.8	130	150	100	240	120	70	170	0
38	M	27.1	115	135	85	190	95	45	145	0
42	F	30.4	125	145	95	210	105	55	155	0
36	M	29.9	118	138	88	205	102	52	152	0
44	F	27.6	122	142	92	215	108	58	158	0
37	M	32.1	135	155	105	250	130	75	175	0
43	F	28.3	128	148	98	225	115	65	165	0
34	M	26.5	105	125	78	175	85	35	135	0
46	F	31.8	138	158	108	235	125	75	175	0
33	M	25.9	102	122	75	170	82	32	132	0
47	F	32.5	142	162	112	245	135	80	180	0
32	M	25.2	100	120	72	165	80	30	130	0
48	F	33.2	145	165	115	250	140	85	185	0
31	M	24.8	98	118	70	160	78	28	128	0
49	F	34.0	148	168	118	255	145	90	190	0
30	M	24.1	95	115	68	155	75	26	125	0
50	F	34.8	150	170	120	260	150	95	195	0
29	M	23.5	92	112	65	150	72	24	122	0
51	F	35.5	152	172	122	265	155	100	200	0
28	M	22.9	90	110	62	145	70	22	120	0
52	F	36.2	155	175	125	270	160	105	205	0
27	M	22.2	88	108	60	140	68	20	118	0
53	F	37.0	158	178	128	275	165	110	210	0
26	M	21.6	85	105	58	135	65	18	115	0
54	F	37.8	160	180	130	280	170	115	215	0
25	M	21.0	82	102	55	130	62	16	112	0
55	F	38.5	162	182	132	285	175	120	220	0
24	M	20.4	80	100	52	125	60	14	110	0
56	F	39.2	165	185	135	290	180	125	225	0
23	M	19.8	78	98	50	120	58	12	108	0
57	F	40.0	168	188	138	295	185	130	230	0
22	M	19.2	75	95	48	115	55	10	105	0
58	F	40.8	170	190	140	300	190	135	235	0
21	M	18.6	72	92	45	110	52	8	102	0
59	F	41.5	172	192	142	305	195	140	240	0
20	M	18.0	70	90	42	105	50	6	100	0
60	F	42.2	175	195	145	310	200	145	245	0
19	M	17.4	68	88	40	100	48	4	98	0
61	F	43.0	178	198	148	315	205	150	250	0
18	M	16.8	65	85	38	95	45	2	95	0
62	F	43.8	180	200	150	320	210	155	255	0
17	M	16.2	62	82	35	90	42	0	92	0
63	F	44.5	182	202	152	325	215	160	260	0
16	M	15.6	60	80	32	85	40	0	90	0
64	F	45.2	185	205	155	330	220	165	265	0
15	M	15.0	58	78	30	80	38	0	88	0
65	F	46.0	188	208	158	335	225	170	270	0
14	M	14.4	55	75	28	75	35	0	85	0
66	F	46.8	190	210	160	340	230	175	275	0
13	M	13.8	52	72	25	70	32	0	82	0
67	F	47.5	192	212	162	345	235	180	280	0
12	M	13.2	50	70	22	65	30	0	80	0
68	F	48.2	195	215	165	350	240	185	285	0
11	M	12.6	48	68	20	60	28	0	78	0
69	F	49.0	198	218	168	355	245	190	290	0
10	M	12.0	45	65	18	55	25	0	75	0
70	F	49.8	200	220	170	360	250	195	295	0
9	M	11.4	42	62	15	50	22	0	72	0
71	F	50.5	202	222	172	365	255	200	300	0
8	M	10.8	40	60	12	45	20	0	70	0
72	F	51.2	205	225	175	370	260	205	305	0
7	M	10.2	38	58	10	40	18	0	68	0
73	F	52.0	208	228	178	375	265	210	310	0
6	M	9.6	35	55	8	35	15	0	65	0
74	F	52.8	210	230	180	380	270	215	315	0
5	M	9.0	32	52	5	30	12	0	62	0
75	F	53.5	212	232	182	385	275	220	320	0
4	M	8.4	30	50	3	25	10	0	60	0
76	F	54.2	215	235	185	390	280	225	325	0
3	M	7.8	28	48	1	20	8	0	58	0
77	F	55.0	218	238	188	395	285	230	330	0
2	M	7.2	25	45	0	15	5	0	55	0
78	F	55.8	220	240	190	400	290	235	335	0
1	M	6.6	22	42	0	10	3	0	52	0
79	F	56.5	222	242	192	405	295	240	340	0
0	M	6.0	20	40	0	5	1	0	50	0
80	F	57.2	225	245	195	410	300	245	345	0

Fig.1 Chronic Liver Dataset

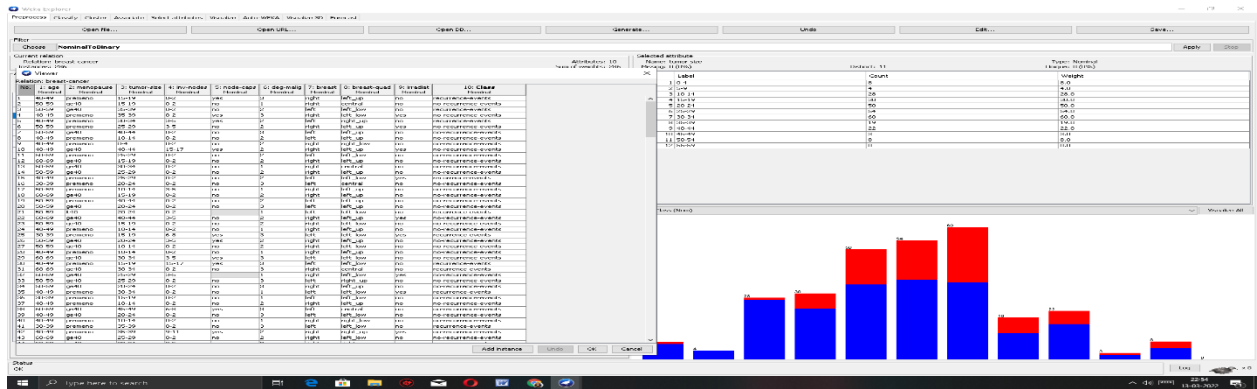


Fig.2 Cancer Dataset

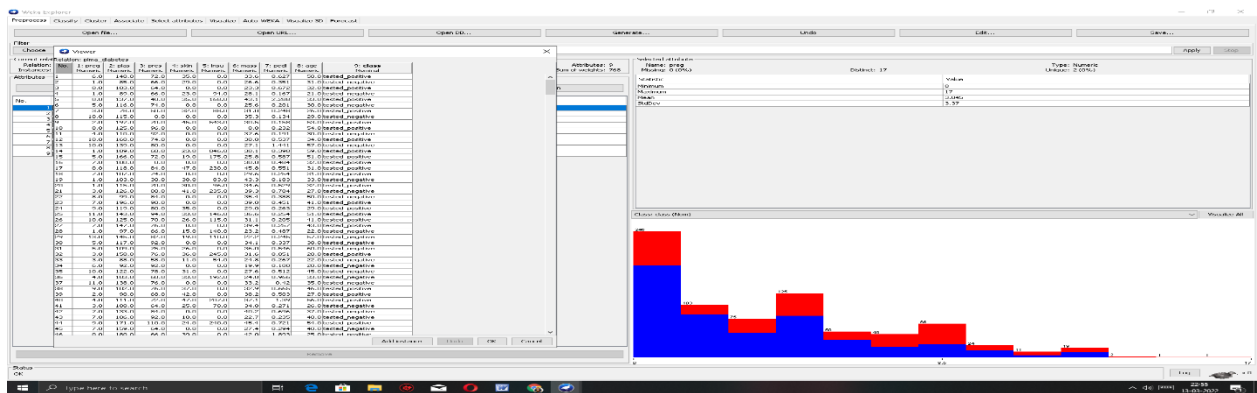


Fig.3 Diabetes Mellitus Dataset

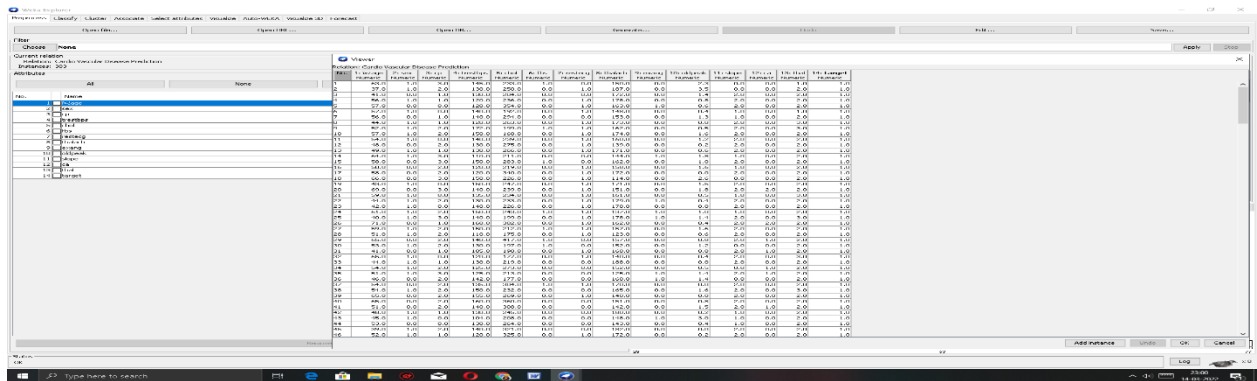


Fig.4 Cardio Vascular Disease Dataset

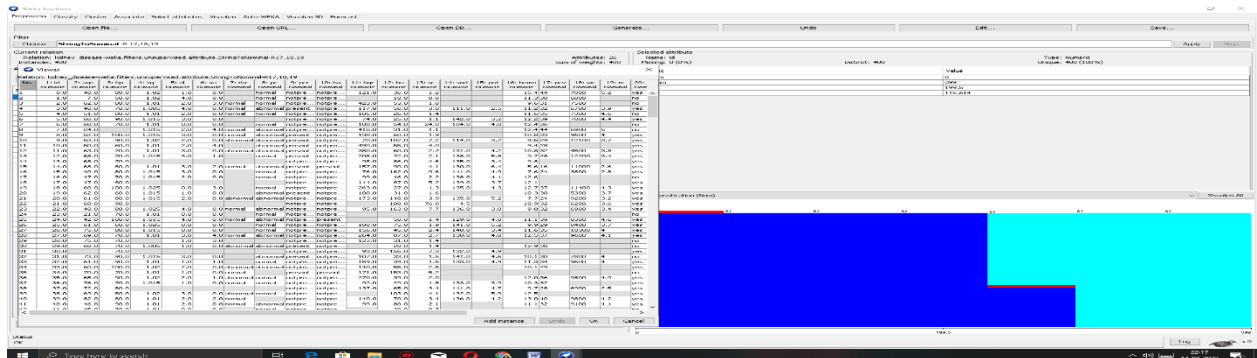


Fig.5 Chronic Kidney Disease Dataset

**Experiments and Observations-1**

Classifier Algorithm=BayesNet on CKD dataset

10 fold cross validation

Classifier Output- === Run information ===

Scheme: weka.classifiers.bayes.BayesNet -D -Q weka.classifiers.bayes.net.search.local.K2 -- -P 1 -S BAYES -E weka.classifiers.bayes.net.estimate.SimpleEstimator -- -A 0.5

Relation: kidney\_disease-weka.filters.unsupervised.attribute.StringToNominal-R17,18,19

Instances: 400 Attributes: 26

Id	age	bp	sg	al	su	rbc	pc	pcc	ba	bgr	bu	sc
Sod	pot	hemo	pcv	wc	rc	htn	dm	cad	appet	pe	ane	

classification

Test mode: 10-fold cross-validation

=== Classifier model (full training set) ===

Bayes Network Classifier

#attributes=26 #classindex=25

=== Stratified cross-validation ===== Summary ===

Correctly Classified Instances 391 97.75 %

Incorrectly Classified Instances 9 2.25 %

Kappa statistic 0.9534

Mean absolute error 0.0167

Root mean squared error 0.1039

Relative absolute error 5.2518 %

Root relative squared error 26.1069 %

Total Number of Instances 400

=== Detailed Accuracy By Class ===

	TP Rate	FP Rate	Precision	Recall	F-Measure	MCC	ROC Area	PRC Area	Class
	0.972	0.013	0.992	0.972	0.982	0.953	0.990	0.987	ckd
	0.000	0.018	0.000	0.000	0.000	-0.009	0.536	0.008	ckd
	1.000	0.000	1.000	1.000	1.000	1.000	1.000	1.000	notckd

Weighted Avg. 0.978 0.008 0.990 0.978 0.984 0.966 0.991 0.987

=== Confusion Matrix ===

a	b	c	<-- classified as
241	7	0	a = ckd
2	0	0	b = ckd
0	0	150	c = notckd

**Experiments and Observations-2**

Classifier Algorithm=NaiveBayes on CKD dataset

10 fold cross validation

Classifier Output- === Run information ===

Scheme: weka.classifiers.bayes.NaiveBayes

Relation: kidney\_disease-weka.filters.unsupervised.attribute.StringToNominal-R17,18,19

Instances: 400

Attributes: 26

Test mode: 10-fold cross-validation

=== Classifier model (full training set) ===

Naive Bayes Classifier

Class

Attribute	ckd	ckd	notckd
	(0.62)	(0.01)	(0.37)

=== Stratified cross-validation ===== Summary ===

Correctly Classified Instances 392 98 %

Incorrectly Classified Instances 8 2 %

Kappa statistic 0.958

Mean absolute error 0.0153

Root mean squared error 0.1131

Relative absolute error 4.8193 %

Root relative squared error 28.4297 %



Total Number of Instances 400

==== Detailed Accuracy By Class ====

TP Rate	FP Rate	Precision	Recall	F-Measure	MCC	ROC Area	PRC Area	Class
0.976	0.013	0.992	0.976	0.984	0.958	0.990	0.990	ckd
0.000	0.003	0.000	0.000	0.000	-0.004	0.099	0.005	ckd
1.000	0.020	0.968	1.000	0.984	0.974	1.000	1.000	notckd

Weighted Avg. 0.980 0.016 0.978 0.980 0.979 0.959 0.990 0.989

==== Confusion Matrix ====

```

a b c <-- classified as
242 1 5 | a = ckd
2 0 0 | b = ckd
0 0 150 | c = notckd
    
```

**Experiments and Observations-3**

Classifier Algorithm=meta.Bagging on CKD dataset

10 fold cross validation

==== Run information ====

Scheme: weka.classifiers.meta.Bagging -P 100 -S 1 -num-slots 1 -I 10 -W weka.classifiers.trees.REPTree -- -M 2 -V 0.001 -N 3 -S 1 -L -1 -I 0.0

Relation: kidney\_disease-weka.filters.unsupervised.attribute.StringToNominal-R17,18,19

Instances: 400 Attributes: 26

Test mode: 10-fold cross-validation

==== Classifier model (full training set) ====

Bagging with 10 iterations and base learner

weka.classifiers.trees.REPTree -M 2 -V 0.001 -N 3 -S 1 -L -1 -I 0.0

==== Stratified cross-validation ===== Summary =====

Correctly Classified Instances 397 99.25 %

Incorrectly Classified Instances 3 0.75 %

Kappa statistic 0.9841

Mean absolute error 0.009

Root mean squared error 0.0663

Relative absolute error 2.8347 %

Root relative squared error 16.6627 %

Total Number of Instances 400

==== Detailed Accuracy By Class ====

TP Rate	FP Rate	Precision	Recall	F-Measure	MCC	ROC Area	PRC Area	Class
0.996	0.013	0.992	0.996	0.994	0.984	0.988	0.980	ckd
0.000	0.000	?	0.000	?	0.440	0.006	0.006	ckd
1.000	0.004	0.993	1.000	0.997	0.995	1.000	1.000	notckd

Weighted Avg. 0.993 0.010 ? 0.993 ? ? 0.990 0.983

==== Confusion Matrix =====

```

a b c <-- classified as
247 0 1 | a = ckd
2 0 0 | b = ckd
0 0 150 | c = notckd
    
```

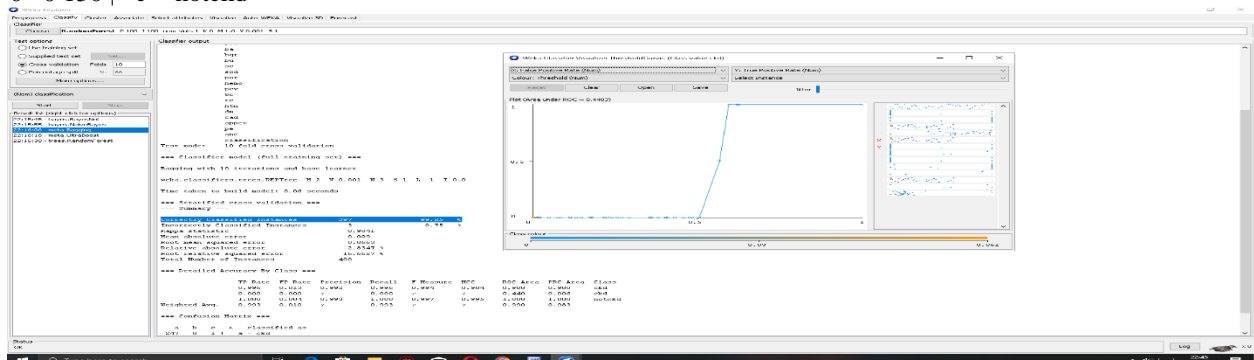


Fig.6 Bagging Classifier with accuracy of 99.25%

**Experiments and Observations-4**

Classifier Algorithm=meta.UltraBoost on CKD dataset

10 fold cross validation

=== Run information ===

```
Scheme:          weka.classifiers.meta.UltraBoost -S 1 -B "weka.classifiers.meta.FilteredClassifier -F
\"weka.filters.unsupervised.attribute.RemoveType -V -T nominal\" -S 1 -W weka.classifiers.bayes.NaiveBayes" -B
"weka.classifiers.meta.FilteredClassifier -F \"weka.filters.unsupervised.attribute.RemoveType -V -T numeric\" -S 1 -W
weka.classifiers.functions.Logistic -- -R 1.0E-8 -M -1 -num-decimal-places 4"
```

Relation: kidney\_disease-weka.filters.unsupervised.attribute.StringToNominal-R17,18,19

Instances: 400

Attributes: 26

Test mode: 10-fold cross-validation

=== Classifier model (full training set) ===

UltraBoost Base classifiers

```
FilteredClassifier using weka.classifiers.bayes.NaiveBayes on data filtered through
weka.filters.unsupervised.attribute.RemoveType -V -T nominal
```

Filtered Header

```
@relation 'kidney_disease-weka.filters.unsupervised.attribute.StringToNominal-R17,18,19-
weka.filters.unsupervised.attribute.RemoveType-V-Tnominal'
```

=== Stratified cross-validation ===== Summary ===

Correctly Classified Instances 397 99.25 %

Incorrectly Classified Instances 3 0.75 %

Kappa statistic 0.9842

Mean absolute error 0.044

Root mean squared error 0.1245

Relative absolute error 13.8347 %

Root relative squared error 31.2791 %

Total Number of Instances 400

=== Detailed Accuracy By Class ===

	TP Rate	FP Rate	Precision	Recall	F-Measure	MCC	ROC Area	PRC Area	Class
	0.996	0.013	0.992	0.996	0.994	0.984	0.988	0.981	ckd
	0.000	0.003	0.000	0.000	0.000	-0.004	0.606	0.009	ckd
	1.000	0.000	1.000	1.000	1.000	1.000	1.000	1.000	notckd
Weighted Avg.	0.993	0.008	0.990	0.993	0.991	0.985	0.991	0.983	

=== Confusion Matrix ===

```
a b c <-- classified as
247 1 0 | a = ckd
2 0 0 | b = ckd
0 0 150 | c = notckd
```

**Experiments and Observations-5**

Classifier Algorithm=RandomForest on CKD dataset

10 fold cross validation

=== Run information ===

Scheme: weka.classifiers.trees.RandomForest -P 100 -I 100 -num-slots 1 -K 0 -M 1.0 -V 0.001 -S 1

Relation: kidney\_disease-weka.filters.unsupervised.attribute.StringToNominal-R17,18,19

Instances: 400 Attributes: 26

=== Stratified cross-validation ===== Summary ===

Correctly Classified Instances 398 99.5 %

Incorrectly Classified Instances 2 0.5 %

Kappa statistic 0.9894

Mean absolute error 0.04

Root mean squared error 0.0813

Relative absolute error 12.5748 %

Root relative squared error 20.4263 %

Total Number of Instances 400

=== Detailed Accuracy By Class ===

	TP Rate	FP Rate	Precision	Recall	F-Measure	MCC	ROC Area	PRC Area	Class
	0.996	0.013	0.992	0.996	0.994	0.984	0.988	0.981	ckd
	0.000	0.003	0.000	0.000	0.000	-0.004	0.606	0.009	ckd
	1.000	0.000	1.000	1.000	1.000	1.000	1.000	1.000	notckd
Weighted Avg.	0.993	0.008	0.990	0.993	0.991	0.985	0.991	0.983	

	1.000	0.013	0.992	1.000	0.996	0.989	0.989	0.983	ckd
	0.000	0.000	?	0.000	?	?	0.298	0.005	ckd
	1.000	0.000	1.000	1.000	1.000	1.000	1.000	1.000	notckd
Weighted Avg.	0.995	0.008	?	0.995	?	?	0.989	0.985	

=== Confusion Matrix ===

a b c ← classified as  
 248 0 0 | a = ckd  
 2 0 0 | b = ckd  
 0 0 150 | c = notckd

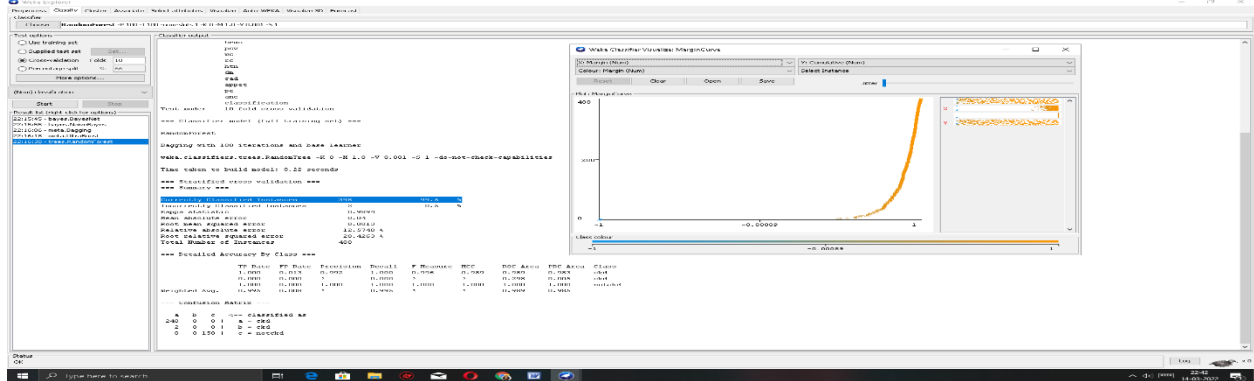


Fig.7 Random Forest Classifier with accuracy of 99.5%

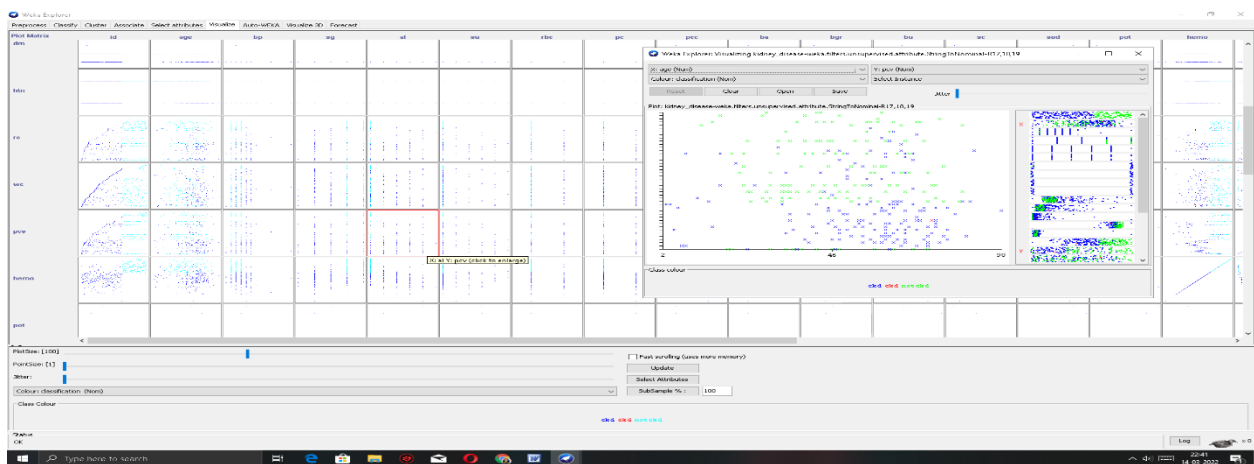


Fig.8 Visualize Curve for Classification

### CONCLUSION

Though it is not an easy task to clearly classify errors and make a predictive analysis from the real time disease dataset from the medical and healthcare domain, but it is very important to apply the Artificial Intelligence - Machine Learning model to clearly analyze, visualize, preprocess, predict, classify based on different symptoms of the patients for a better early case diagnosis with a good accuracy level which in turn a necessity for the medical and health care sector and for the health of most people in this world who suffered from these kind of major diseases. As we can see, we have taken 5 major disease dataset as input and process for a better result by applying different existing as well as some hybrid methods for an early case diagnosis of the disease. Hence Machine Learning hybrid classification algorithms will be an option in the medical and healthcare field which will be helpful for early case diagnosis of the major diseases based on different cases.

### REFERENCES

[1] Y Khourdifi, M Bahaj, “Heart disease prediction and classification using machine learning algorithms optimized by particle swarm optimization and ant colony optimization”, Int. J. Intell. Eng. Syst, volume 12, issue 1, 2019  
 [2] S Mohan, C Thirumalai, G Srivastava, “Effective heart disease prediction using hybrid machine learning techniques”, IEEE Access, volume 7, 2019





- [3] M Chen, Y Hao, K Hwang, L Wang, L Wang, "Disease prediction by machine learning over big data from healthcare communities", Ieee Access, volume 5, 2017
- [4] S Chae, S Kwon, D Lee, "Predicting infectious disease using deep learning and big data", International journal of environmental research and public health, volume 15, issue 8, 2018
- [5] A U Haq, J P Li, M H Memon, S Nazir, R Sun, "A hybrid intelligent system framework for the prediction of heart disease using machine learning algorithms", Mobile Information Systems, 2018
- [6] M Maniruzzaman, M J Rahman, B Ahammed, M M Abedin, "Classification and prediction of diabetes disease using machine learning paradigm", Health Information Science and Systems, volume 8, issue 1, 2020
- [7] Uddin, S., Khan, A., Hossain, M. et al., "Comparing different supervised machine learning algorithms for disease prediction", BMC Med Inform Decis Mak 19, 281, 2019
- [8] D. Dahiwade, G. Patle and E. Meshram, "Designing Disease Prediction Model Using Machine Learning Approach", 3rd International Conference on Computing Methodologies and Communication (ICCMC), 2019, pp. 1211-1215, 2019
- [9] M. Chen, Y. Hao, K. Hwang, L. Wang and L. Wang, "Disease prediction by machine learning over big data from healthcare communities", IEEE Access, vol. 5, no. 1, pp. 8869-8879, 2017.
- [10] Ajinkya Kunjir, Harshal Sawant and Nuzhat F. Shaikh, "Data Mining and Visualization for prediction of Multiple Diseases in Healthcare", IEEE big data analytics and computational intelligence, pp. 2325, Oct 2017.
- [11] B. Nithya and V. Ilango, "Predictive Analytics in Health Care Using Machine Learning Tools and Techniques", International Conference on Intelligent Computing and Control Systems, 2017.
- [12] S. Leoni Sharmila, C. Dharuman and P. Venkatesan, "Disease Classification Using Machine Learning Algorithms - A Comparative Study", International Journal of Pure and Applied Mathematics, vol. 114, no. 6, pp. 1-10, 2017.
- [13] Allen Daniel Sunny, Sajal Kulshreshtha, Satyam Singh, Srinabh, Mohan Ba and H Sarojadevi, "Disease Diagnosis System By Exploring Machine Learning Algorithms", International Journal of Innovations in Engineering and Technology (IJJET), vol. 10, no. 2, May 2018.
- [14] Shraddha Subhash Shirsath, "Disease Prediction Using Machine Learning Over Big Data", International Journal of Innovative Research in Science, vol. 7, no. 6, June 2018.
- [15] Deepthi, Y., Kalyan, K.P., Vyas, M., Radhika, K., Babu, D.K., Krishna Rao, N.V., "Disease Prediction Based on Symptoms Using Machine Learning.", Energy Systems, Drives and Automations. Lecture Notes in Electrical Engineering, vol 664. Springer, Singapore, 2018
- [16] Mir A, Dhage SN, "Diabetes disease prediction using machine learning on big data of healthcare", Fourth international conference on computing communication control and automation (ICCUBEA), 2018
- [17] Ray S, "A quick review of machine learning algorithms", International conference on machine learning, big data, cloud and parallel computing (Com-IT-Con), India, 14th-16th Feb 2019
- [18] Shetty D, Rit K, Shaikh S, Patil N, "Diabetes disease prediction using data mining", International conference on innovations in information, embedded and communication systems (ICIIECS), 2017
- [19] Thirunavukkarasu K, Singh AS, Irfan M, Chowdhury A, "Prediction of liver disease using classification algorithms", 4th international conference on computing communication and automation (ICCCA), 2018
- [20] SanthanaKrishnan J, Geetha S, "Prediction of heart disease using machine learning algorithms", 1st international conference on innovations in information and communication technology (ICIICT), 2019