

# Evaluation of Machine Learning Algorithms in Glass Information Mining

Charitha E<sup>1</sup>, Anjan Babu G<sup>2</sup>

Dept of Computer Science, SV University, Tirupati

**Abstract**— This paper means to bunch the glass dataset with respect to their metric data by using the best four AI portrayal computations like Decision Tree, Naive Bayes, KNN and SVM to find which estimation will have the alternative to bring to the table all the really testing accuracy. This paper zeroed in on contrasting the four AI models are on the Glass dataset and assess different execution measurements. Our trial results shows that tests display that Decision Tree and SVM performed feasibly in distinctive glass arrangement.

## I. INTRODUCTION

Procedure is the way toward tracking down a model or a cutoff that portrays and sees information classes and thoughts, to utilize the model to expect the classes of things whose class mark isn't known. Information deals can be seen as a two-stage measure: learning step in which a classifier is made portraying an appointed arrangement of classes or assessments by separating the status set contained educational summary tuples and their connected names. In the resulting improvement model is utilized for demand by first studying the reasonable precision of classifier worked during the mysterious new development [2]. It is finished utilizing the test information. The precision of classifier on a given test set tuples is level of tuples that are irrefutably referred to by the classifier. In the event that the exactness is over some worthy level, the classifier can be utilized to expect future tuples whose class mark isn't known.

Depiction is a sort of information assessment that can be utilized to make models depicting goliath information classes. Approach is an information mining hypothesis used to expect pack pay for information models [5]. It is one of the central constructions in information mining and is utilized in different applications, for example, plan interest, torture affirmation, client relationship the pioneers, and given out showing up. The objective of the depiction assessments is to amass a model from an immense heap of preparation information whose target class names are known and similarly this model is utilized to pack covered cases [6].

## II. MACHINE LEARNING (ML)

IDS engineers use various strategies for interference revelation. One of these methodologies relies upon ML. ML methodologies can predict and recognize risks before they achieve critical security events [5].

ML, a piece of man-made thinking, is a coherent request stressed over the arrangement and progression of computations that grant PCs to propel rehearses reliant upon observational data, for instance, from sensor data or data bases. A critical point of convergence of ML research is to this sort out some way to see complex models and make shrewd decisions subject to data [2][4]. ML has a wide extent of employments, including web crawlers, clinical end, text and handwriting affirmation, picture screening, load deciding, advancing and bargains finding, and so forth

Simulated intelligence procedures can be used to find and bring information by the techniques for models which can't be recognized adequately by human discernment. These segments are classifiers which portray the association data drawing closer into the system to pick whether the development is an attack or some normal activity.

The model can be farsighted to make assumptions later on, or clear to get data from data. To play out a farsighted or illustrative task, ML overall use two essential techniques: Classification and Clustering. All together, the program should predict the most probable characterization, class or name for novel insight into one or various predefined classes or name while gathering, the classes are not predefined during the learning cycle. In any case if the justification the IDS is to isolate between common or interference traffic, game plan is endorsed and if we hope to recognize the sort of interference, gathering can be more valuable [5][6]. To further develop the interference revelation structure and diminishing the fake negative and sham positive, which can be attempted by the usage of different computations? In this paper, Naive Bayes Classifier, Decision Tree Classifier, KNN Classifier and SVM Classifier are used for planning data and testing it.

### III. METHODOLOGY

In this research, we will use four machine learning algorithms for glass data detection.

#### 3.1 Decision Tree

Choice tree learning is perhaps the best strategy for directed plan learning. Choice trees are a direct recursive design for conveying a progressive portrayal measure in which a case, portrayed by a ton of characteristics, is assigned to one of a disjoint game plan of classes [5][6]. A Decision tree is a tree structure which bunches a data test into one of its possible classes. Choice trees are used to isolate data by making decision rules from the gigantic proportion of open information. A Decision tree classifier has a direct design which can be insignificantly taken care of and that capably orchestrates new data.

Choice trees contain center points and leaves. Each center point in the tree incorporates testing a particular quality and each leaf of the tree implies a class. Generally speaking, the test differentiates a quality worth and a consistent. Leaf centers give a gathering that applies to all models that show up at the leaf, or a great deal of portrayals, or a probability course over each possible game plan. To bunch a dark event, it is coordinated down the tree as demonstrated by the assessments of the characteristics attempted in reformist centers, and when a leaf is reached, the case is portrayed by the class dispensed to the leaf.

#### 3.2 Naive Bayes

Naive Bayes is truly outstanding and useful gathering computations. Gullible Bayes Classifier that is the probabilistic classifier reliant upon the Bayes Theorem. Gullible Bayes classifier expects that the effect of the characteristics regard on a given class is free on the assessment of various features [6]. The classifier simply picks the imprint with the most raised probability, given the information features. The blameless portion of the classifier is that it acknowledges a strong self-sufficiency between attributes; fundamentally it expects the probabilities for all of the information features are self-governing of each other.

Let  $H$  be a hypothesis and  $X$  is a data living in a particular  $C$  class. By then  $P(H/X)$  is known as the back probability that imparts our conviction level on a theory  $H$  after  $X$  data is given.  $P(H)$  addresses the  $H$  prior probability for all model data.  $P(H/X)$  is without a doubt more illuminating than  $P(H)$ . Bayes' speculation depicts the association between  $P(H/X)$ ,  $P(H)$ , and  $P(X)$  is showed up on condition 1 as follow:

$$P(H/X) = P(X/H) * P(H) / P(X)$$

#### 3.3 Support Vector Machine (SVM)

The SVM is another kind of ML method subject to quantifiable learning speculation. Because of incredible progression and a higher exactness, SVM has become the assessment point of convergence of the ML social class. SVMs are set of related managed learning systems used for portrayal and backslide [8][9]. A couple of late assessments have definite that the SVM all around are good for passing on better similarly as gathering precision than the other data plan estimations. SVM depends on genuine learning speculation by Vapnik et al proposed another learning method, which depends on a set number of tests in the information contained in the current getting ready content to get the best request results.

An uncommon property of SVM can't avoid being, SVM simultaneously limit the exploratory gathering screw up and support the numerical edge. So SVM called Maximum Margin Classifiers. SVM relies upon the Structural peril Minimization. SVM map input vector to a higher dimensional space where a maximal detaching hyperplane is assembled. Two equivalent hyperplanes are based on each side of the hyperplane that distinctive the data. The disengaging hyperplane is the hyperplane that expand the partition between the two equivalent hyperplanes. An assumption that is made that the greater the edge or detachment between these equivalent hyperplanes the better the hypothesis misstep of the classifier [1] [3].

#### 3.4 K nearest neighbor (KNN)

KNN is a popular portrayal computation showing extraordinary execution characteristics and a short time frame of getting ready time. KNN is clear, for the most part notable, significantly capable and convincing estimation for plan affirmation. KNN is a straight forward classifier, where tests are requested ward on the class of their nearest neighbor [2][5].

The KNN is a non-parametric request technique, which is clear yet amazing a significant part of the time [6]. For a data record  $d$  to be gathered, its  $K$  nearest neighbors are recuperated, and these constructions a space of  $d$ . Overwhelming majority projecting a polling form among the data records in the space is regularly used to pick the gathering for ' $d$ ' with or without

considered partition based weighting. In any case, to apply KNN we need to pick a reasonable impetus for K, and the accomplishment of request is a great deal of subject to this value. It could be said, the KNN technique is uneven by K. There are various strategies for picking the K regard, yet a fundamental one is to run the estimation conventionally with different K regards and pick the one with the best show.

#### IV. EXPLORATORY RESULTS

This part will give a format over the refined outcomes, the pre-owned information and the appraisal joint effort to orchestrate. We have considered the Glass information from UCI Machine Learning Repository dataset [7]. The assessments have been driven by utilizing Python Programming. It gives different information mining assessments and depiction instruments for information appraisal and farsighted outlining, with graphical UIs that assists client with effectively running these calculations on datasets. Python Sklearn two or three standard information mining undertakings, that are, information preprocessing, depiction, fall away from the faith, gathering, highlight choice. The Glass dataset contains 214 instances and 11 attributes and contains six class labels as shown in the figure-1. The detailed analysis of the dataset is shown in the figure-2 and figure-3 through density and histogram plots.

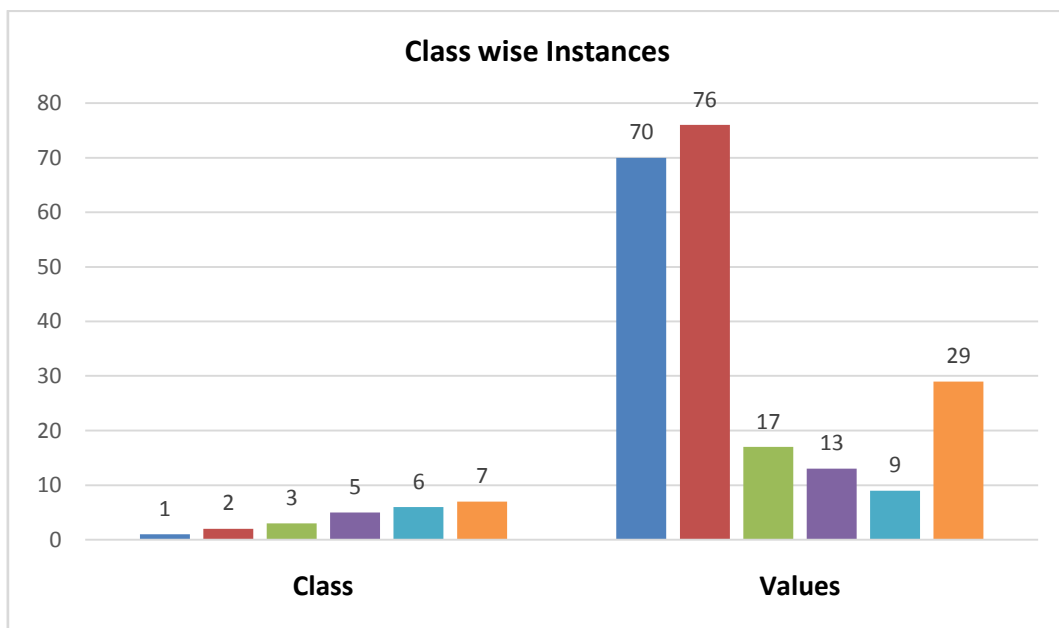


FIGURE 1: Class distribution instances

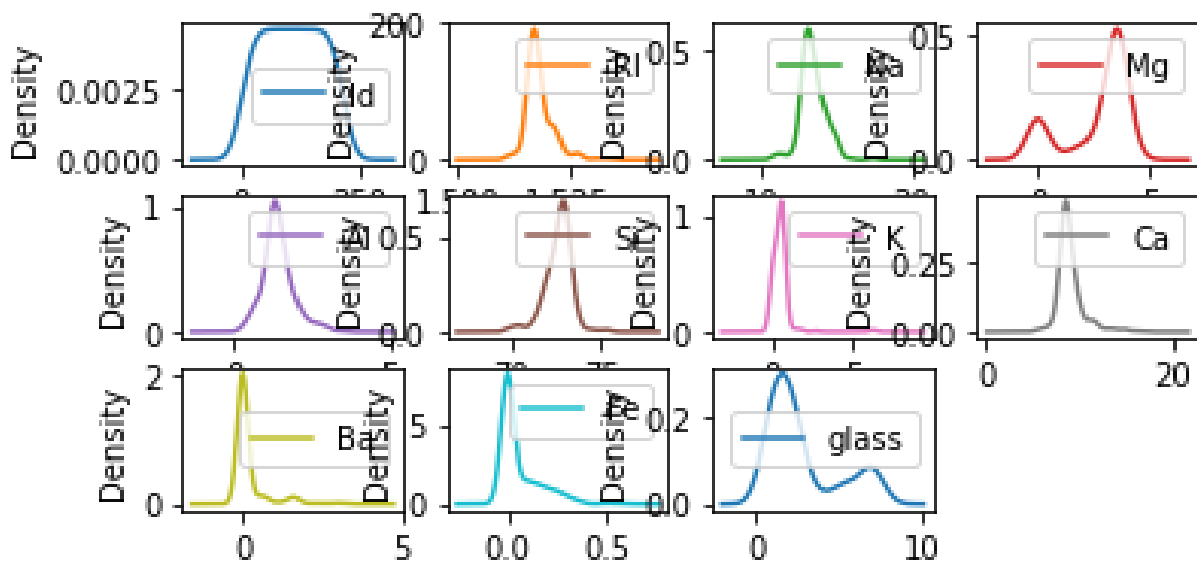


FIGURE 2: Density plot of glass dataset

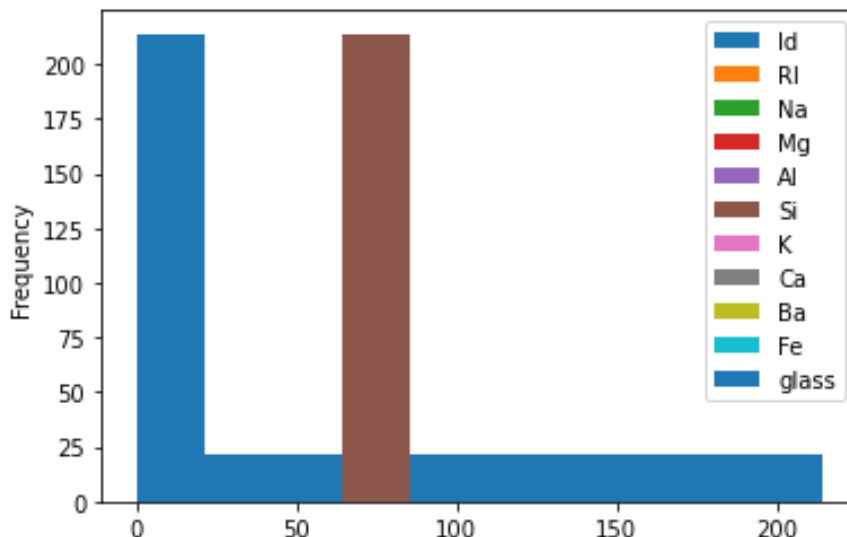


FIGURE 3: Histogram plot for entire glass dataset

#### 4.1 Results

So as to approve the forecast consequences of the examination of the four well known information mining methods and the 10-crease hybrid approval is utilized. The k-overlap hybrid approval is normally used to decrease the blunder came about because of arbitrary examining in the correlation of the accuracies of various forecast models. The whole arrangement of information is arbitrarily separated into k folds with similar number of cases in each crease. The preparation and testing are performed for k times and one overlay is chosen for additional testing while the rest are chosen for additional preparation. The current investigation partitioned the information into 10 folds where 1 overlap was for trying and 9 folds were for preparing for the 10-crease hybrid approval. We evaluate our four models using different execution estimations like exactness, Precision and Recall, the Experimental results are showed up in the table-1 same showed up in the figure-4.

TABLE 1  
PERFORMANCE OF ML ALGORITHMS

Algorithm	Accuracy	Precision	Recall
Decision Tree	96.9	97	97
Naïve Bayes	84.6	87	85
SVM	96.9	97	97
KNN	81.5	82	82

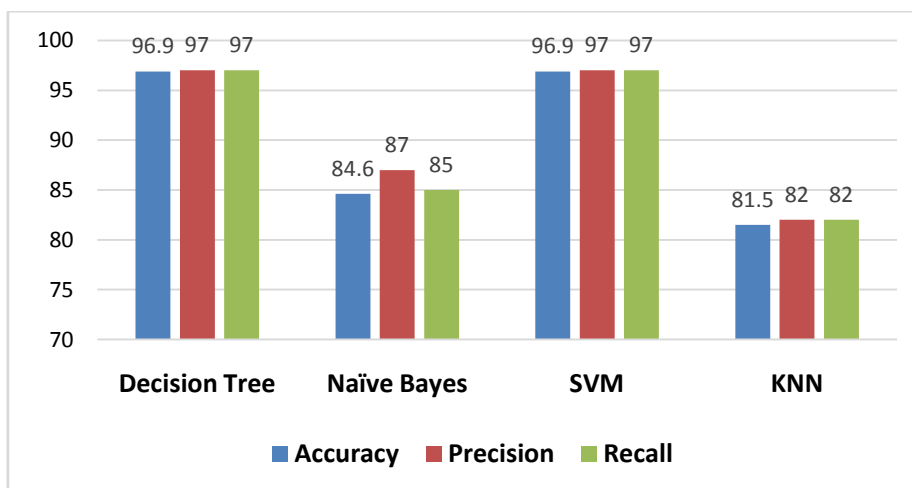


FIGURE 4: Performance of ML algorithms

We find in the Figure-4 for our preparation information, the introduction of the Decision Tree estimation has accomplished 96.9% Accuracy, Naïve Bayes has 84.6%, KNN has accomplished 81.5% and SVM model has achieved 96.9%.

As the result from assessment among the four figuring's, we find that most vital precision of Classification models are Decision Tree and SVM for both have same accuracy (96.9%). Precisely when veered from exactness and survey are also higher in the Decision Tree and SVM model when appeared differently in relation to other two models.

The experimental results screen shots are shown from the figure-5 to figure-8.

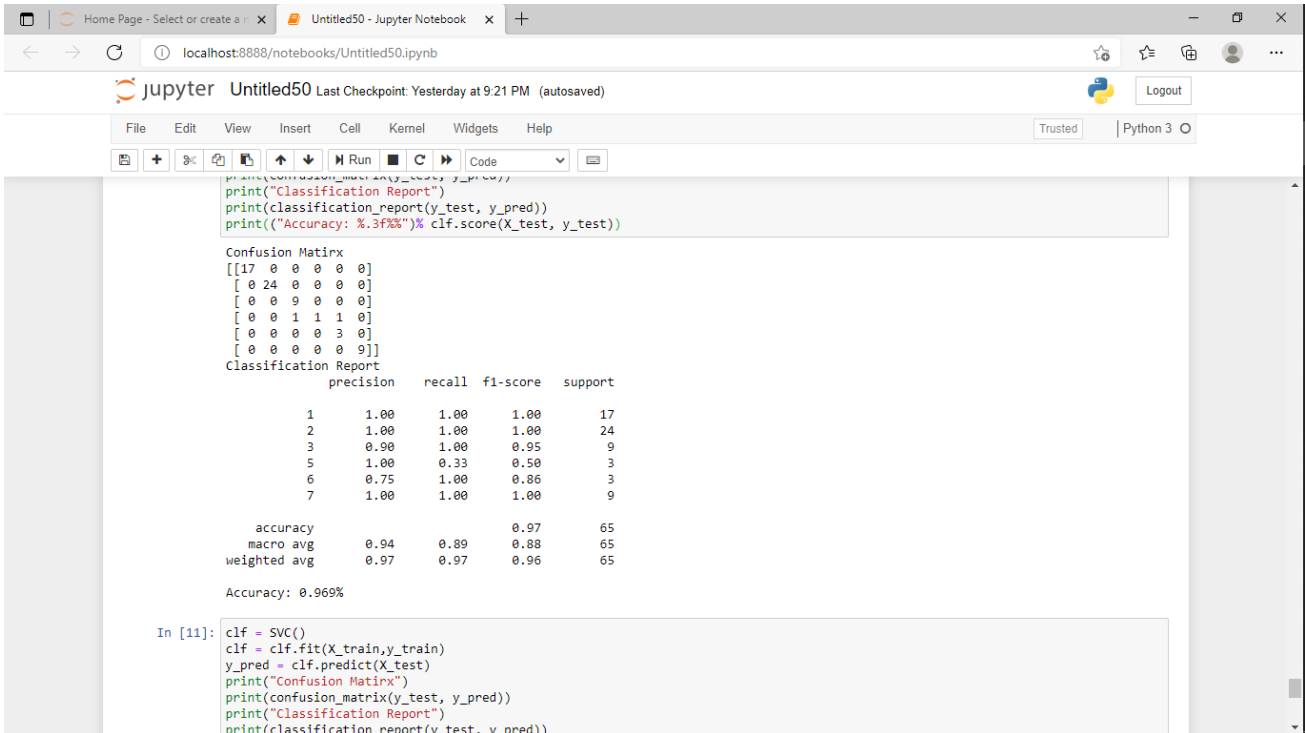


FIGURE-5: Screen shots of experimental results

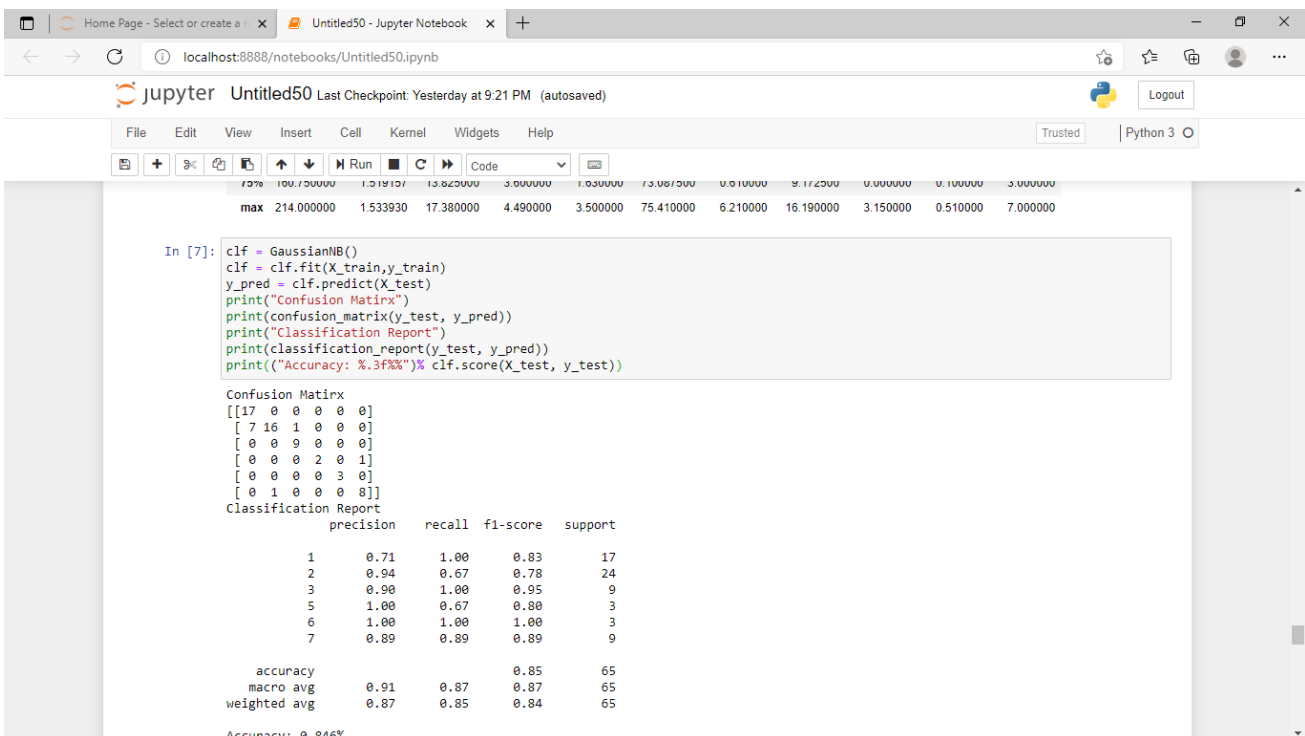
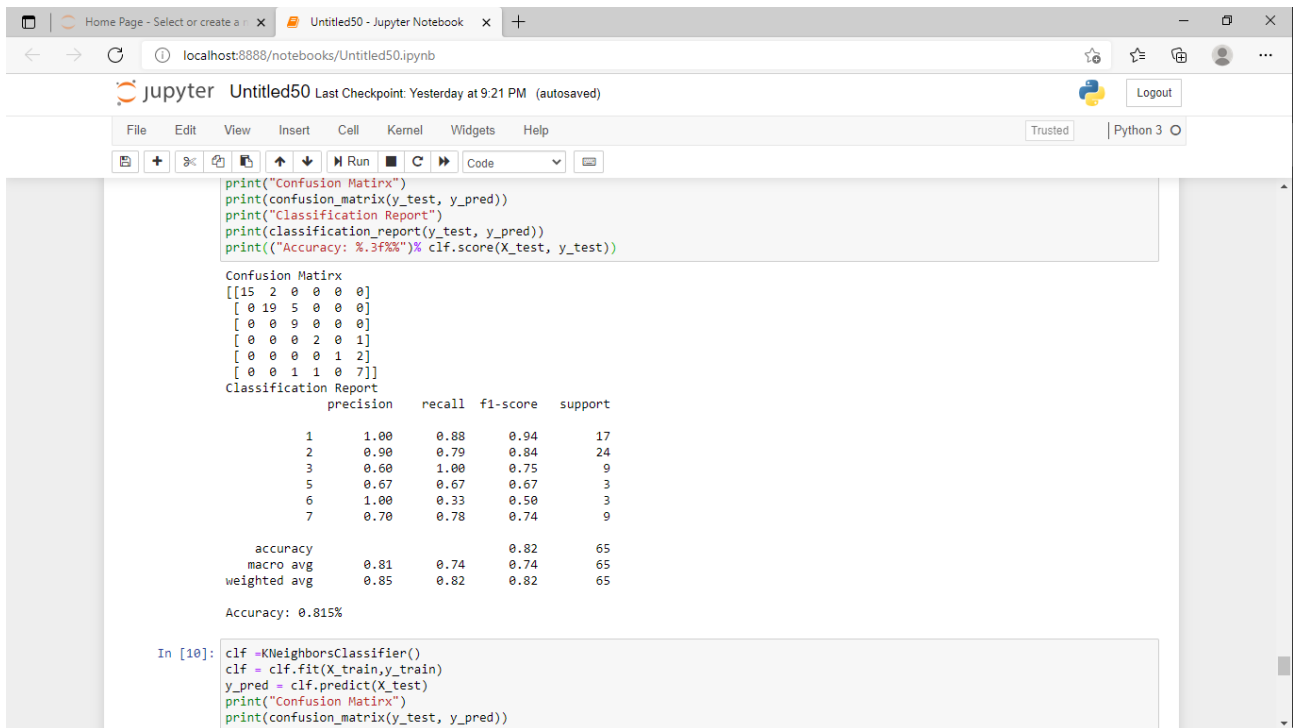
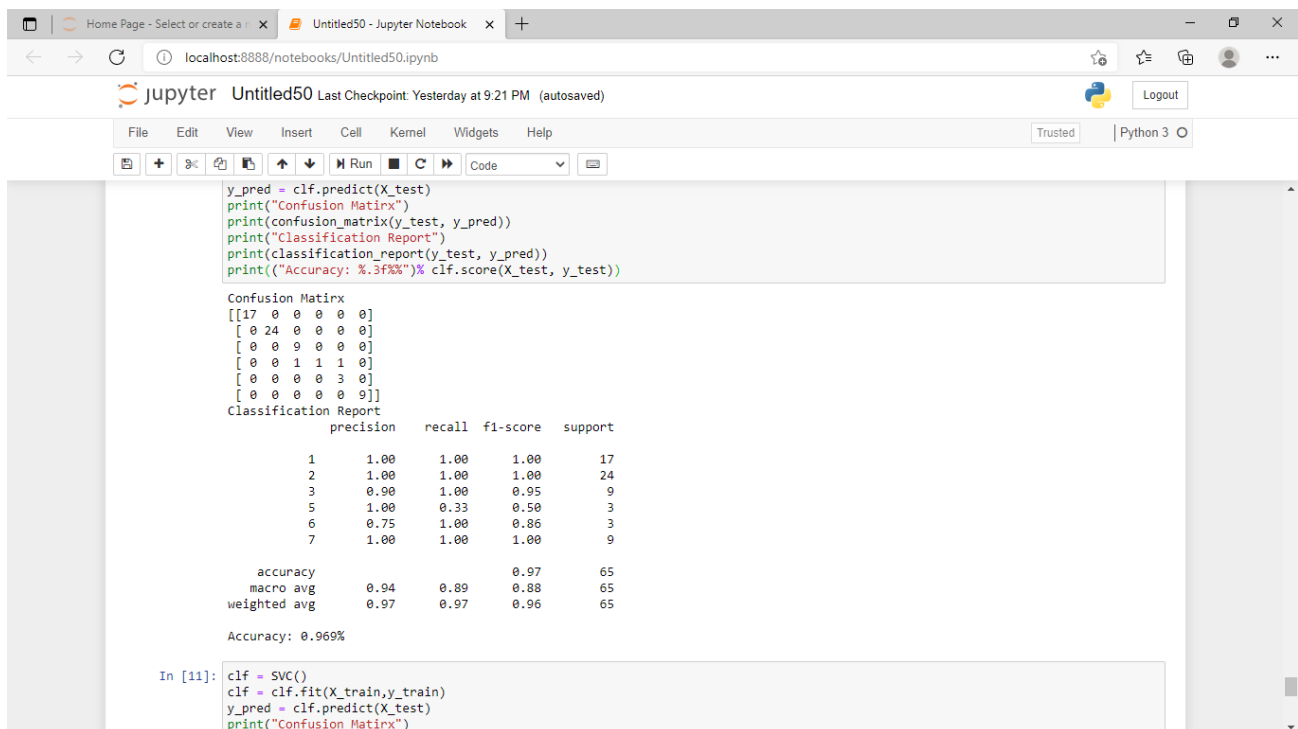


FIGURE 6: Screen shots of experimental results



**FIGURE 7: Screen shots of experimental results**



**FIGURE 8: Screen shots of experimental results**

## V. CONCLUSION

In this paper, four unique sorts of ML models were applied in particular Decision Tree, Naïve-Bayes, K Nearest Neighbors and Support Vector Machine for the glass identification framework. The exhibition of all these ML models were watched and looked at dependent on changed standard assessment boundaries, for example, Accuracy, Precision and Recall of the test information. Our test has been done with four distinctive grouping calculations for the dataset and in that decision tree and SVM both shows a high exactness contrasted with every other calculation. It was seen that the Decision Tree Classifier and SVM calculations performed in a way that is better than different models creating an exactness of 96.9%.

## REFERENCES

- [1] Bernhard Schölkopf and Alex Smola, Learning with kernels. MIT Press, Cambridge, MA, 2002.
- [2] D. Hand, H. Mannila, P. Smyth.: Principles of Data Mining. The MIT Press. (2001)
- [3] G. Bo and H. Xianwu, "SVM multi-class classification," Journal of Data Acquisition & Processing, vol. 21, pp. 334-339, 2006.
- [4] G. Ravi Kumar, K. Nagamani and G. Anjan Babu, "A Framework of Dimensionality Reduction Utilizing PCA for Neural Network Prediction", Lecture Notes on Data Engineering and Communications Technologies, ISBN 978-981-15-0977-3, Volume 37, PP:173-180, 2020
- [5] Han J, Kamber M. Data Mining: Concepts and Techniques [J]. Data Mining Concepts Models Methods & Algorithms Second Edition, 2011, 5 (4)
- [6] Ian H. Witten and Eibe Frank. Data Mining: Practical machine learning tools and techniques.2nd ed. San Francisco: Morgan Kaufmann, 2005.
- [7] UCI machine learning repository. <http://archive.ics.uci.edu/ml/>.
- [8] Vapnik, V.N. Statistical Learning Theory. John Wiley and Sons, New York, USA, 1998.
- [9] Vapnik, V.N. The Natural of Statistical Learning theory. Springer-Verleg, New York, USA 1995.