

(REVIEW ARTICLE)



Comprehensive review of supervised machine learning algorithms to identify the best and error free

Omankwu Obinnaya Chinecherem *, Ugwuja Nnenna Esther and Kanu Chigbundu

Department of Computer Sc. Michael Okpara University of Agriculture, Umudike, Nigeria.

International Journal of Scholarly Research in Engineering and Technology, 2023, 02(01), 013–019

Publication history: Received on 05 November 2022; revised on 24 December 2022; accepted on 26 December 2022

Article DOI: <https://doi.org/10.56781/ijret.2023.2.1.0028>

Abstract

Supervised classification is one of the tasks most frequently carried out by the intelligent systems. Supervised Machine Learning (SML) is the search for algorithms that reason from externally supplied instances to produce general hypotheses, which then make predictions about future instances. This paper; compares various supervised. Seven different machine learning algorithms were considered: Decision Table, Random Forest (RF) , Naïve Bayes (NB) , Support Vector Machine (SVM), Neural Networks (Perceptron), JRip and Decision Tree [J48]l. And also reviews various Supervised Machine Learning (ML) classification techniques with the aim of identifying the Best and Error free algorithm.

Keywords: Machine Learning; Classifiers; Data Mining Techniques; Data Analysis; Learning Algorithms; Supervised Machine Learning

1 Introduction

Machine learning is one of the fastest growing areas of computer science, with far-reaching applications. It refers to the automated detection of meaningful patterns in data. Machine learning tools are concerned with endowing programs with the ability to learn and adapt (Shai et al 2014).

Machine Learning has become one of the mainstays of Information Technology and with that, a rather central, albeit usually hidden, part of our life. With the ever increasing amounts of data becoming available there is a good reason to believe that smart data analysis will become even more pervasive as a necessary ingredient for technological progress.

There are several applications for Machine Learning (ML), the most significant of which is data mining. People are often prone to making mistakes during analyses or, possibly, when trying to establish relationships between multiple features (Lemnar, 2012).

Data Mining and Machine Learning are Siamese twins from which several insights can be derived through proper learning algorithms. There has been tremendous progress in data mining and machine learning as a result of evolution of smart and Nano technology which brought about curiosity in finding hidden patterns in data to derive value. The fusion of statistics, machine learning, information theory, and computing has created a solid science, with a firm mathematical base, and with very powerful tools.

Machine learning algorithms are organized into a taxonomy based on the desired outcome of the algorithm. Supervised learning generates a function that maps inputs to desired outputs.

*Corresponding author: Omankwu Obinnaya C

Unprecedented data generation has made machine learning techniques become sophisticated from time to time. This has called for utilization for several algorithms for both supervised and unsupervised machine learning. Supervised learning is fairly common in classification problems because the goal is often to get the computer to learn a classification system that we have created (Taiwo, O. A. (2010)..

ML is perfectly intended for accomplishing the accessibility hidden within Big Data. ML hand over's on the guarantee of extracting importance from big and distinct data sources through outlying less dependence scheduled on individual track as it is data determined and spurts at machine scale. Machine learning is fine suitable towards the intricacy of handling through dissimilar data origin and the vast range of variables as well as amount of data concerned where ML prospers on increasing datasets. The extra data supply into a ML structure, the more it be able to be trained and concern the consequences to superior value of insights. At the liberty from the confines of individual level thought and study, ML is clever to find out and show the patterns hidden in the data (Pradeep, 2017).

One standard formulation of the supervised learning task is the classification problem: The learner is required to learn (to approximate the behavior of) a function which maps a vector into one of several classes by looking at several input-output examples of the function. Inductive machine learning is the process of learning a set of rules from instances (examples in a training set), or more generally speaking, creating a classifier that can be used to generalize from new instances. The process of applying supervised ML to a real-world problem is described in Figure 1.

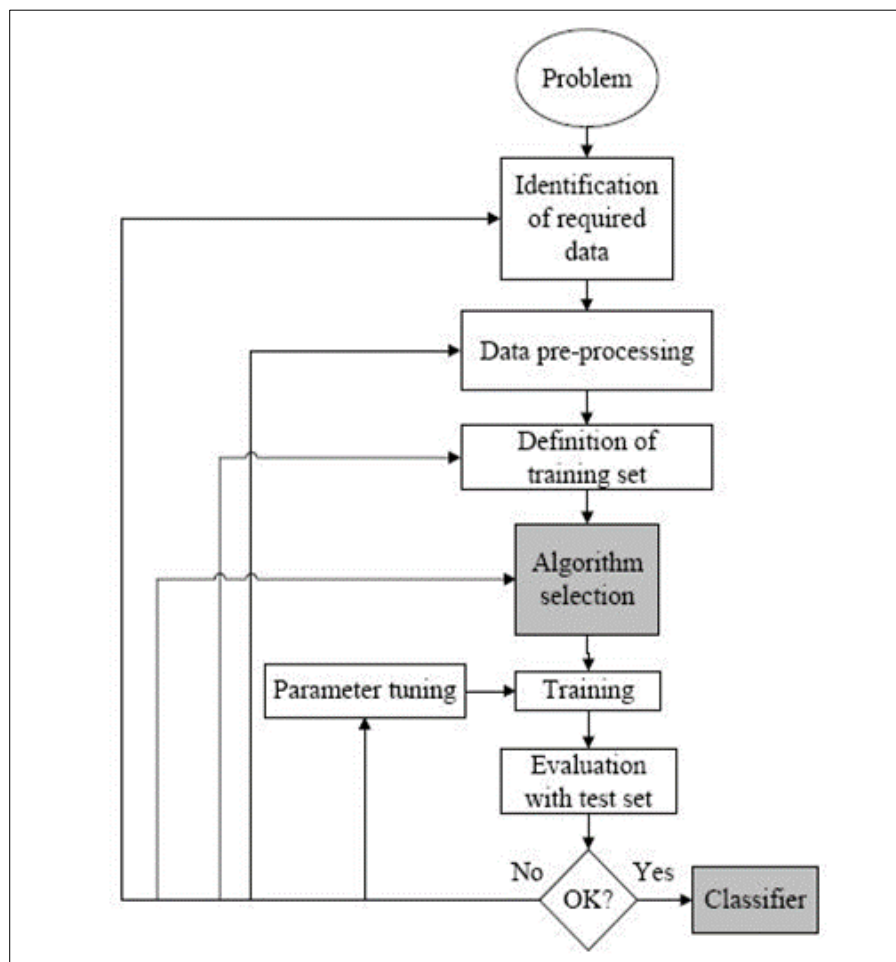


Figure 1 The Processes of Supervised Machine Learning (Osisanwo et al; 2019)

This work focuses on the classification of ML algorithms and determining the most efficient algorithm with highest accuracy and precision. As well as establishing the performance of different algorithms on large and smaller data sets with a view classify them correctly and give insight on how to build supervised machine learning models.

The remaining part of this work is arranged as follows: Section 2 presents the literature review discussing classification of different supervised learning algorithms; section 3 gives the conclusion and recommendation for further works.

2 Literature review

2.1 Classification of Supervised Learning Algorithms

According to Taiwo, 2010, the supervised machine learning algorithms which deals more with classification includes the following: Linear Classifiers, Logistic Regression, Naïve Bayes Classifier, Perceptron, Support Vector Machine; Quadratic Classifiers, K-Means Clustering, Boosting, Decision Tree, Random Forest (RF); Neural networks, Bayesian Networks and so on.

2.1.1 Linear Classifiers

Linear models for classification separate input vectors into classes using linear (hyperplane) decision boundaries (Good, 1951). The goal of classification in linear classifiers in machine learning is to group items that have similar feature values, into groups. Timothy, (2018) Stated that a linear classifier achieves this goal by making a classification decision based on the value of the linear combination of the features. A linear classifier is often used in situations where the speed of classification is an issue, since it is rated the fastest classifier Timothy, (2018) Also, linear classifiers often work very well when the number of dimensions is large, as in document classification, where each element is typically the number of counts of a word in a document. The rate of convergence among data set variables however depends on the margin. Roughly speaking, the margin quantifies how linearly separable a dataset is, and hence how easy it is to solve a given classification problem (Setiono et al, 2020),

2.1.2 Logistic regression

This is a classification function that uses class for building and uses a single multinomial logistic regression model with a single estimator. Logistic regression usually states where the boundary between the classes exists, also states the class probabilities depend on distance from the boundary, in a specific approach. This moves towards the extremes (0 and 1) more rapidly when data set is larger. These statements about probabilities which make logistic regression more than just a classifier. It makes stronger, more detailed predictions, and can be fit in a different way; but those strong predictions could be wrong. Logistic regression is an approach to prediction, like Ordinary Least Squares (OLS) regression. However, with logistic regression, prediction results in a dichotomous outcome. Logistic regression is one of the most commonly used tools for applied statistics and discrete data analysis. Logistic regression is linear interpolation (Newsom, 2015).

2.1.3 Naive Bayesian (NB) Networks

These are very simple Bayesian networks which are composed of directed acyclic graphs with only one parent (representing the unobserved node) and several children (corresponding to observed nodes) with a strong assumption of independence among child nodes in the context of their parent (Good, 1951). Thus, the independence model (Naive Bayes) is based on estimating (Nilsson, 1965). Bayes classifiers are usually less accurate than other more sophisticated learning algorithms (such as ANNs). However, Domingos, P. & Pazzani, M. (2020) performed a large-scale comparison of the naive Bayes classifier with state-of-the-art algorithms for decision tree induction, instance-based learning, and rule induction on standard benchmark datasets, and found it to be sometimes superior to the other learning schemes, even on datasets with substantial feature dependencies. Bayes classifier has attribute-independence problem which was addressed with Averaged One-Dependence Estimators (Hormozi et al, 2012).

2.1.4 Multi-layer Perceptron

This is a classifier in which the weights of the network are found by solving a quadratic programming problem with linear constraints, rather than by solving a non-convex, unconstrained minimization problem as in standard neural network training (Taiwo, 2010). Other well-known algorithms are based on the notion of perceptron (Rosenblatt, 2018). Perceptron algorithm is used for learning from a batch of training instances by running the algorithm repeatedly through the training set until it finds a prediction vector which is correct on all of the training set. This prediction rule is then used for predicting the labels on the test set (Kotsiantis, 2017).

2.1.5 Support Vector Machines (SVMs)

These are the most recent supervised machine learning technique [24]. Support Vector Machine (SVM) models are closely related to classical multilayer perceptron neural networks. SVMs revolve around the notion of a $-\text{margin}$ —either side of a hyperplane that separates two data classes. Maximizing the margin and thereby creating the largest possible distance between the separating hyperplane and the instances on either side of it has been proven to reduce an upper bound on the expected generalisation error (Vapnik, 1995).

2.1.6 *K-means*

According to Bishop, (1995) and Tapas Kanungo, (2002).K-means is one of the simplest unsupervised learning algorithms that solve the well-known clustering problem. The procedure follows a simple and easy way to classify a given data set through a certain number of clusters (assume k clusters) fixed a priori.K-Means algorithm is employed when labeled data is not available (Alex et al,2008).

2.1.7 *Decision Trees*

Decision Trees (DT) are trees that classify instances by sorting them based on feature values. Each node in a decision tree represents a feature in an instance to be classified, and each branch represents a value that the node can assume. Instances are classified starting at the root node and sorted based on their feature values (Kotsiantis, 2017).Decision tree learning, used in data mining and machine learning, uses a decision tree as a predictive model which maps observations about an item to conclusions about the item's target value. More descriptive names for such tree models are classification trees or regression trees (Hastie et al, 2001). Decision tree classifiers usually employ post-pruning techniques that evaluate the performance of decision trees, as they are pruned by using a validation set. Any node can be removed and assigned the most common class of the training instances that are sorted to it (Kotsiantis, 2017).

2.1.8 *Neural Networks*

[2]opined Neural Networks (NN) that can actually perform a number of regression and/or classification tasks at once,although commonly each network performs only one. In the vast majority of cases, therefore, the network will have a single output variable, although in the case of many-state classification problems, this may correspond to a number of output units (the post-processing stage takes care of the mapping from output units to output variables).Artificial Neural Network (ANN) depends upon three fundamental aspects, input and activation functions of the unit, network architecture and the weight of each input connection. Given that the first two aspects are fixed, the behavior of the ANN is defined by the current values of the weights. The weights of the net to be trained are initially set to random values, and then instances of the training set are repeatedly exposed to the net. The values for the input of an instance are placed on the input units and the output of the net is compared with the desired output for this instance. Then, all the weights in the net are adjusted slightly in the direction that would bring the output values of the net closer to the values for the desired output. There are several algorithms with which a network can be trained (Neocleous C. &Schizas C. (2002).

2.1.9 *Bayesian Network*

A Bayesian Network (BN) is a graphical model for probability relationships among a set of variables (features). Bayesian networks are the most well-known representative of statistical learning algorithms (Kotsiantis, 2017). The most interesting feature of BNs, compared to decision trees or neural networks, is most certainly the possibility of taking into account prior information about a given problem, in terms of structural relationships among its features (Kotsiantis, 2017). A problem of BN classifiers is that they are not suitable for datasets with many features (Cheng et al, 2002). This prior expertise, or domain knowledge, about the structure of a Bayesian network can take the following forms:

- Declaring that a node is a root node, i.e., it has no parents.
- Declaring that a node is a leaf node, i.e., it has no children.
- Declaring that a node is a direct cause or direct effect of another node.
- Declaring that a node is not directly connected to another node.
- Declaring that two nodes are independent, given a condition-set.

2.2 ***Features of Machine Learning Algorithms Supervised machine learning techniques are***

Applicable in numerous domains. A number of Machine Learning (ML) application oriented papers can be found in Setiono et al, (2020) and Witten et al, (2005).

Generally, SVMs and neural networks tend to perform much better when dealing with multi-dimensions and continuous features. On the other hand, logic-based systems tend to perform better when dealing with discrete/categorical features. For neural network models and SVMs, a large sample size is required in order to achieve its maximum prediction accuracy whereas NB may need a relatively small dataset.

There is general agreement that k-NN is very sensitive to irrelevant features: this characteristic can be explained by the way the algorithm works. Moreover, the presence of irrelevant features can make neural network training very inefficient, even impractical. Most decision tree algorithms cannot perform well with problems that require diagonal

partitioning. The division of the instance space is orthogonal to the axis of one variable and parallel to all other axes. Therefore, the resulting regions after partitioning are all hyperrectangles. The ANNs and the SVMs perform well when multi-collinearity is present and a nonlinear relationship exists between the input and output features.

Naive Bayes (NB) requires little storage space during both the training and classification stages: the strict minimum is the memory needed to store the prior and conditional probabilities. The basic kNN algorithm uses a great deal of storage space for the training phase, and its execution space is at least as big as its training space. On the contrary, for all non-lazy learners, execution space is usually much smaller than training space, since the resulting classifier is usually a highly condensed summary of the data. Moreover, Naive Bayes and the kNN can be easily used as incremental learners whereas rule algorithms cannot. Naive Bayes is naturally robust to missing values since these are simply ignored in computing probabilities and hence have no impact on the final decision. On the contrary, kNN and neural networks require complete records to do their work.

Finally, Decision Trees and NB generally have different operational profiles, when one is very accurate the other is not and vice versa. On the contrary, decision trees and rule classifiers have a similar operational profile. SVM and ANN have also a similar operational profile. No single learning algorithm can uniformly outperform other algorithms over all datasets.

Finally, Decision Trees and NB generally have different operational profiles, when one is very accurate the other is not and vice versa. On the contrary, decision trees and rule classifiers have a similar operational profile. SVM and ANN have also a similar operational profile. No single learning algorithm can uniformly outperform other algorithms over all datasets.

Different data sets with different kind of variables and the number of instances determine the type of algorithm that will perform well. There is no sing

3 Conclusion and Recommendation

ML classification requires thorough fine tuning of the parameters and at the same time sizeable number of instances for the data set. It is not a matter of time to build the model for the algorithm only but precision and correct classification. Therefore, the best learning algorithm for a particular data set, does not guarantee the precision and accuracy for another set of data whose attributes are logically different from the other.

However, the key question when dealing with ML classification is not whether a learning algorithm is superior to others, but under which conditions a particular method can significantly outperform others on a given application problem.

Meta-learning is moving in this direction, trying to find functions that map datasets to algorithm performance (Neocleous, 2002). To this end, meta-learning uses a set of attributes, called meta-attributes, to represent the characteristics of learning tasks, and searches for the correlations between these attributes and the performance of learning algorithms. Some characteristics of learning tasks are: the number of instances, the proportion of categorical attributes, the proportion of missing values, the entropy of classes, etc.

Brazdil, (2003) provided an extensive list of information and statistical measures for a dataset. After a better understanding of the strengths and limitations of each method, the possibility of integrating two or more algorithms together to solve a problem should be investigated. The objective is to utilize the strengths of one method to complement the weaknesses of another. If we are only interested in the best possible classification accuracy, it might be difficult or impossible to find a single classifier that performs as well as a good ensemble of classifiers. SVM, NB and RF machine learning algorithms can deliver high precision and accuracy regardless of the number of attributes and data instances. Therefore, ML algorithms require precision, accuracy and minimum error to have supervised predictive machine learning.

Compliance with ethical standards

Acknowledgments

I acknowledge the efforts of students of computer science, Michael Okpara University of Agriculture, Umudike and my Mentor Prof Anigbogu O.S whom I follow his footsteps in the field of AI- Machine Learning.

Disclosure of conflict of interest

Conflict of interest was that Ugwuja Nnenna and Chigbundu Kanu wanted the study to be carried out on a specific Machine Learning Algorithm, but after much deliberations, we finally agreed to work on the comprehensive review of the Machine Learning Algorithms

Statement of informed consent

Informed consent was obtained from all individual participants included in the study.

References

- [1] Alex S.&Vishwanathan, S.V.N. (2008). Introduction to Machine Learning. Published by the press syndicate of the University of Cambridge, Cambridge, United Kingdom. Copyright© Cambridge University Press 2008. ISBN: 0521-82583-0. Available at KTH website:<https://www.kth.se/social/upload/53a14887f276540ebc81aec3/online.pdf> Retrieved from website: <http://alex.smola.org/drafts/thebook.pdf>
- [2] Bishop, C. M. (1995). Neural Networks for Pattern Recognition. Clarendon Press, Oxford, England. 1995. Oxford University Press, Inc. New York, NY, USA ©1995 ISBN:0198538642 Available at: http://cs.du.edu/~mitchell/mario_books/Neural_Networks_f_or_Pattern_Recognition_-_Christopher_Bishop.pdf
- [3] Brazdil P., Soares C. &da Costa, J. (2003). Ranking Learning Algorithm using IBL and Meta – Learning on accuracy and Time Result: Machine Learning Vol. 50; Issue 3; 2003
- [4] Cheng, J., Greiner, R., Kelly, J., Bell, D.& Liu, W. (2002). Learning Bayesian networks from data: An information-theory based approach. Artificial Intelligence Volume 137. Copyright © 2002. Published by Elsevier Science B.V. All rights reserved pp. 43 – 90. Available at science Direct: <http://www.sciencedirect.com/science/article/pii/S0004370202001911>
- [5] Domingos, P. &Pazzani, M. (2020). On the optimality of the simple Bayesian classifier under zero-one loss. Machine Learning Volume 29, pp. 103–130 Copyright © 1997 Kluwer Academic Publishers. Manufactured in The Netherlands. Available at University of Trento website: <http://disi.unitn.it/~p2p/RelatedWork/Matching/domingos97optimality.pdf>
- [6] Elder, J. (n.d). Introduction to Machine Learning and Pattern Recognition. Available at LASSONDE University EECS Department York website: http://www.eecs.yorku.ca/course_archive/2011-12/F/4404-5327/lectures/01%20Introduction.pdf
- [7] Good, I.J. (1951). Probability and the Weighing of evidence Philosophy; volume 26, issue 97; Published by Charles Griffin and Company, London; 1950.
- [8] Hormozi, H., Hormozi, E. &Nohooji, H. R. (2012). The Classification of the Applicable Machine Learning Methods in Robot Manipulators. International Journal of Machine Learning and Computing (IJMLC), Vol. 2, No. 5, 2012 doi: 10.7763/IJMLC.2012.V2.189pp. 560 – 563. Available at IJMLC website: <http://www.ijmlc.org/papers/189-C00244-001.pdf>
- [9] Kotsiantis, S. B. (2017). Supervised Machine Learning: A Review of Classification Techniques. Informatica 31 (2007). Pp. 249 – 268. Retrieved from IJS website: <http://wen.ijs.si/ojs-2.4.3/index.php/informatica/article/download/148/140>.
- [10] Lemnaru C. (2012). Strategies for dealing with Real World Classification Problems, (Unpublished PhD thesis) Faculty of Computer Science and Automation; University Technica; Din Cluj-Napoca
- [11] Logistic Regression pp. 223 – 237. Available at: <https://www.stat.cmu.edu/~cshalizi/uADA/12/lectures/ch12.pdf>
- [12] Neocleous C. &Schizas C. (2002). Artificial Neural Network Learning: A Comparative Review. In: Vlahavas I.P., Spyropoulos C.D. (eds)Methods and Applications of Artificial Intelligence. Hellenic Conference on Artificial IntelligenceSETN 2002. Lecture Notes in Computer Science, Volume 2308. Springer, Berlin, Heidelberg, doi: 10.1007/3-540-46014-4_27 pp. 300-313. Available at: https://link.springer.com/chapter/10.1007/3-540-46014-4_27
- [13] Newsom,I.(2015).Data Analysis II Logistic Regression

- [14] Nilsson, N.J. (1965). Learning machines. New York: McGraw-Hill. Published in: Journal of IEEE Transactions on Information Theory Volume 12 Issue 3, 1966. doi: 10.1109/TIT.1966.1053912 pp. 407 – 407. Available at ACM digital library website: <http://dl.acm.org/citation.cfm?id=2267404>
- [15] Pradeep, K. R. & Naveen, N. C. (2017). A Collective Study of Machine Learning (ML) Algorithms with Big Data Analytics (BDA) for Healthcare Analytics (HcA). International Journal of Computer Trends and Technology (IJCTT) – Volume 47 Number 3, 2017. ISSN: 2231-2803, doi: 10.14445/22312803/IJCTT-V47P121, pp 149 – 155. Available from IJCTT website: <http://www.ijcttjournal.org/2017/Volume47/number-3/IJCTT-V47P121.pdf>
- [16] Rob Schapire (n.d) Machine Learning Algorithms for Classification.
- [17] Rosenblatt, F. (2018), Principles of Neurodynamics. Spartan, New York.
- [18] 18. Setiono R. and Loew, W. K. (2020), FERNN: An algorithm for fast extraction of rules from neural networks, Applied Intelligence.
- [19] ShaiShalev-Shwartz and Shai Ben-David (2014). Understanding Machine Learning From Theory to Algorithms
- [20] T. Hastie, R. Tibshirani, J. H. Friedman (2001) – The elements of statistical learning,|| Data mining, inference, and prediction, 2001, New York: Springer Verlag.
- [21] Taiwo, O. A. (2010). Types of Machine Learning Algorithms New Advances in Machine Learning, Yagang Zhang (Ed.), ISBN: 978-953-307-034-6, InTech, University of Portsmouth. United Kingdom. Pp 3 – 31. Available at InTech open website: <http://www.intechopen.com/books/new-advances-in-machine-learning/types-of-machine-learning-algorithms>.
- [22] Tapas Kanungo, D. M. (2002). A local search approximation algorithm for k-means clustering. Proceedings of the eighteenth annual symposium on Computational geometry Barcelona, Spain: ACM Press
- [23] Timothy Jason Shepard, P. J. (2018). Decision Fusion Using a Multi-Linear Classifier. In Proceedings of the International Conference on Multisource-Multisensor Information Fusion.
- [24] Vapnik, V. N. (1995). The Nature of Statistical Learning Theory (2nd Edition). Springer Verlag. Retrieved from <https://www.andrew.cmu.edu/user>
- [25] Witten, I. H. & Frank, E. (2005). Data Mining: Practical machine learning tools and techniques (2nd ed.), ISBN: 0-12-088407-0, Morgan Kaufmann Publishers, San Francisco, CA, U.S.A. © 2005 Elsevier Inc. Retrieved from website: [ftp://93.63.40.27/pub/manuela.sbarra/Data Mining Practical Machine Learning Tools and Techniques - WEKA.pdf](ftp://93.63.40.27/pub/manuela.sbarra/Data%20Mining%20Practical%20Machine%20Learning%20Tools%20and%20Techniques%20-%20WEKA.pdf)
- [26] Osisanwo F.Y., Akinsola J.E.T., Awodele O., Hinmikaiye J. O., Olakanmi O and Akinjobi J.(2017). Supervised Machine Learning Algorithms: Classification and Comparison. International Journal of Computer Trends and Technology (IJCTT) – Volume 48 Number 3 June 2017