



MACHINE LEARNING ALGORITHMS - A REVIEW

Tanuja Verma

ABSTRACT

In this paper, various machine learning algorithms have been talk about. These algorithms are used for various purposes like data mining, image processing, predictive analytics, etc. The main advantage of using machine learning is that, once an algorithm grasp what to do with data, it can do its work automatically

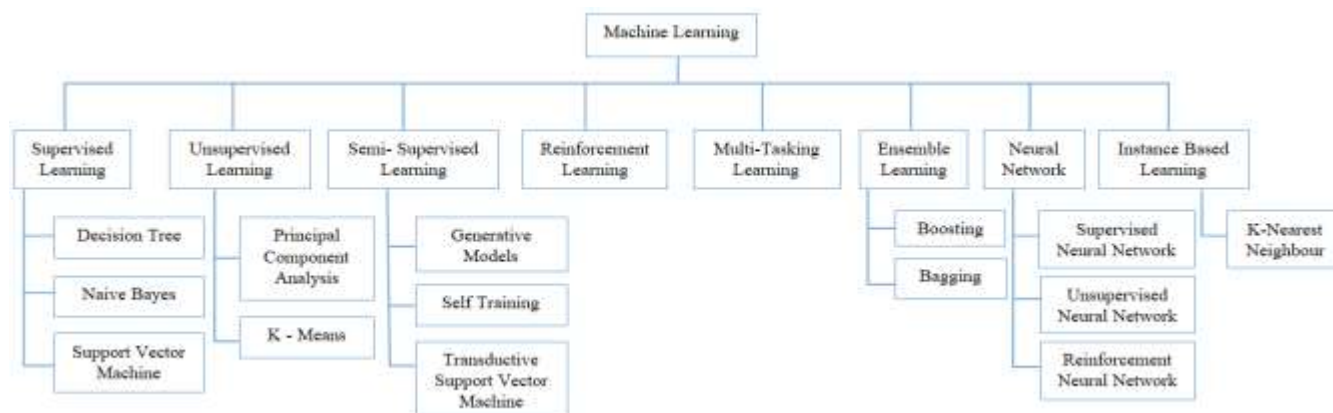
KEYWORDS: Algorithm, Machine Learning, Pseudo Code

I. INTRODUCTION

Machine learning is used to teach machines how to handle the data more efficiently. Sometimes after viewing the data, we cannot interpret the pattern or extract information from the data. In that case, we apply machine learning [1]. With the abundance of datasets available, the demand for machine learning is in rise.

Many industries from medicine to military apply machine learning to extract relevant information.

The purpose of machine learning is to learn from the data. Many studies have been done on how to make machines learn by themselves. Many mathematicians and programmers apply several approaches to find the solution of this problem. Some of them are demonstrated in below figure.



II. TYPES OF LEARNING

1. Supervised Learning

The supervised machine learning algorithms are those algorithms which needs external assistance. The input dataset is divided into train and test dataset. The train dataset has output variable which needs to be predicted or classified. All algorithms learn some kind of patterns from the training dataset and apply them to the test dataset for prediction or classification [2]. Three most famous supervised machine learning algorithms have been discussed here

1.1 Decision Tree: Decision trees are those types of trees which groups attributes by sorting them based on their values. Decision tree is used mainly for classification purpose. Each tree consists of nodes and branches. Each node represents attributes in a

group that is to be classified and each branch represents a value that the node can take [2].

1.2 Naive Bayes: Naïve Bayes mainly targets the text classification industry. It is mainly used for clustering and classification purpose [3]. The underlying architecture of Naïve Bayes depends on the conditional probability. It creates trees based on their probability of happening. These trees are also known as Bayesian Network.

1.3 Support Vector Machine: Another most widely used state-of-the-art machine learning technique is Support Vector Machine (SVM). It is mainly used for classification. SVM works on the principle of margin calculation. It basically, draws margins between the classes. The margins are drawn in such a fashion that the distance between the margin and the classes is maximum and hence, minimizing the classification error.



2. Unsupervised Learning

The unsupervised learning algorithm learns few features from the data. When new data is introduced, it uses the previously learned features to recognize the class of the data. It is mainly used for clustering and feature reduction.

The two main algorithms for clustering and dimensionality reduction techniques are discussed below.

2.1. Principal Component Analysis: In Principal Component Analysis or PCA, the dimension of the data is reduced to make the computations faster and easier. To understand how PCA works, let's take an example of 2D data. When the data is being plot in a graph, it will take up two axes. PCA is applied on the data, the data then will be 1D.

2.2. K-Mean Clustering: Clustering or grouping is a type of unsupervised learning technique that when initiates, creates groups automatically. The item which possesses similar characteristics are put in the same cluster. This algorithm is called k-means because it creates k distinct clusters. The mean of the values in a particular cluster is the center of that cluster [4].

3. Semi-Supervised Learning

Semi-Supervised learning algorithm is a technique which combines the power of both supervised and unsupervised learning. It can be fruit-full in those areas of machine learning and data mining where the unlabeled data is already present and getting the labeled data is a tedious process [5]. There are many categories of semi-supervised learning. Some of which are discussed below:

3.1. Generative Models: Generative models are one of the oldest semi-supervised learning method assumes a structure like $p(x,y) = p(y)p(x|y)$ where $p(x|y)$ is a mixed distribution e.g. Gaussian mixture models. Within the unlabeled data, the mixed components can be identifiable. One labeled example per component is enough to confirm the mixture distribution.

3.2. Self Training: In self-training, a classifier is trained with a portion of labeled data. The classifier is then fed with unlabeled data. The unlabeled points and the predicted labels are added together in the training set. This procedure is then repeated further. Since the classifier is learning itself, hence the name self-training.

3.3 Transductive SVM: Transductive support vector machine or TSVM is an extension of SVM. In TSVM, the labeled and unlabeled data both are considered. It is used to label the unlabeled data in such a way that the margin is maximum between the labeled and unlabeled data. Finding an exact solution by TSVM is a NP-hard problem.

4. Reinforcement Learning

Reinforcement learning is a type of learning which makes decisions based on which actions to take such that the outcome is more positive. The learner has no knowledge which actions to take until it's been given a situation. The action which is taken by the learner may affect situations and their actions in the future. Reinforcement learning solely depends on two criteria: trial and error search and delayed outcome [6].

5. Multi-Tasking Learning

Multitask learning has a simple goal of helping other learners to perform better. When multitask learning algorithms are applied on a task, it remembers the procedure how it solved the problem or how it reaches to the particular conclusion. The algorithm then uses these steps to find the solution of other similar problem or task. This helping of one algorithm to another can also be termed as inductive transfer mechanism. If the learners share their experience with each other, the learners can learn concurrently rather than individually and can be much faster [7].

6. Ensemble Learning

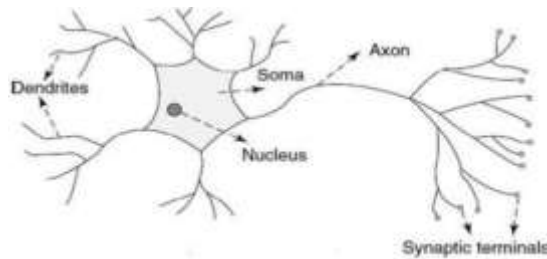
When various individual learners are combined to form only one learner then that particular type of learning is called ensemble learning. The individual learner may be Naïve Bayes, decision tree, neural network, etc. Ensemble learning is a hot topic since 1990s. It has been observed that, a collection of learners is almost always better at doing a particular job rather than individual learners [8]. Two popular Ensemble learning techniques are given below:

6.1. Boosting: Boosting is a technique in ensemble learning which is used to decrease bias and variance. Boosting creates a collection of weak learners and convert them to one strong learner. A weak learner is a classifier which is barely correlated with true classification. On the other hand, a strong learner is a type of classifier which is strongly correlated with true classification.

6.2. Bagging: Bagging or bootstrap aggregating is applied where the accuracy and stability of a machine learning algorithm needs to be increased. It is applicable in classification and regression. Bagging also decreases variance and helps in handling overfitting.

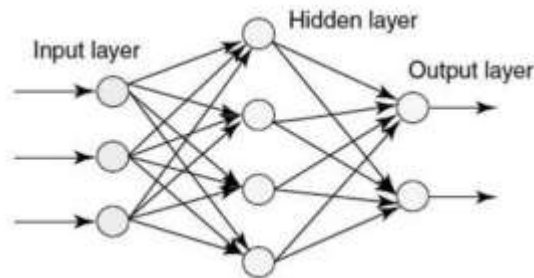
7. Neural Network Learning

The neural network (or artificial neural network or ANN) is derived from the biological concept of neurons. A neuron is a cell like structure in a brain. To understand neural network, one must understand how a neuron works. A neuron has mainly four parts (see Figure). They are dendrites, nucleus, soma and axon.



The dendrites receive electrical signals. Soma processes the electrical signal. The output of the process is carried by the axon to the dendrite terminals where the output is sent to next neuron.

The nucleus is the heart of the neuron. The inter-connection of neuron is called neural network where electrical impulses travel around the brain.



An artificial neural network behaves the same way. It works on three layers. The input layer takes input (much like dendrites). The hidden layer processes the input (like soma and axon). Finally, the output layer sends the calculated output (like dendrite terminals) [9]. There are basically three types of artificial neural network: supervised, unsupervised and reinforcement.

7.1. Supervised Neural Network: In the supervised neural network, the output of the input is already known. The predicted output of the neural network is compared with the actual output. Based on the error, the parameters are changed, and then fed into the neural network again. Supervised neural network is used in feed forward neural network.

7.2. Unsupervised Neural Network: Here, the neural network has no prior clue about the output the input. The main job of the network is to categorize the data according to some similarities. The neural network checks the correlation between various inputs and groups them.

7.3. Reinforcement Neural Network: In reinforced neural network, the network behaves as if a human communicates with the environment. From the environment, a feedback has been provided to the network acknowledging the fact that whether the decision taken by the network is right or wrong. If the decision is right, the connection which points to that particular output is strengthened. The connections are weakened otherwise. The network has no previous information about the output.

8. Instance Based Learning

In instance-based learning, the learner learns a particular type of pattern. It tries to apply the same pattern to the newly fed data. Hence the name instance-based. It is a type of lazy learner which waits for the test data to arrive and then act on it together with training data. The complexity of the learning algorithm increases with the size of the data. Given below is a well-known example of instance-based learning which is k-nearest neighbor.

8.1 K-Nearest Neighbor: In k-nearest neighbor (or KNN), the training data (which is well-labeled) is fed into the learner. When the test data is introduced to the learner, it compares both the data. K most correlated data is taken from training set. The majority of k is taken which serves as the new class for the test data [10].

III. CONCLUSION

This paper surveys various machine learning algorithms. Today each and every person is using machine learning knowingly or unknowingly. From getting a recommended product in online shopping to updating photos in social networking sites. This paper gives an introduction to most of the popular machine learning algorithms.

REFERENCES

1. W. Richert, L. P. Coelho, "Building Machine Learning Systems with Python", Packt Publishing Ltd., ISBN 978-1-78216-140-0
2. S.B. Kotsiantis, "Supervised Machine Learning: A Review of Classification Techniques", *Informatica* 31 (2007) 249-268
3. D. Lowd, P. Domingos, "Naïve Bayes Models for Probability Estimation"
4. S. S. Shwartz, Y. Singer, N. Srebro, "Pegasos: Primal Estimated sub - Gradient Solver for SVM", *Proceedings of the 24th International Conference on Machine Learning, Corvallis, OR, 2007*
5. X. Zhu, A. B. Goldberg, "Introduction to Semi - Supervised Learning", *Synthesis Lectures on Artificial Intelligence and Machine Learning*, 2009, Vol. 3, No. 1, Pages 1-130
6. R. S. Sutton, "Introduction: The Challenge of Reinforcement Learning", *Machine Learning*, 8, Page 225-227, Kluwer Academic Publishers, Boston, 1992
7. R. Caruana, "Multitask Learning", *Machine Learning*, 28, 41-75, Kluwer Academic Publishers, 1997



8. D. Opitz, R. Maclin, "Popular Ensemble Methods: An Empirical Study", *Journal of Artificial Intelligence Research*, 11, Pages 169- 198, 1999
9. V. Sharma, S. Rai, A. Dev, "A Comprehensive Study of Artificial Neural Networks", *International Journal of Advanced Research in Computer Science and Software Engineering*, ISSN 2277128X, Volume 2, Issue 10, October 2012
10. P. Harrington, "Machine Learning in Action", Manning Publications Co., Shelter Island, New York, ISBN 9781617290183, 2012