

Accuracy of Prediction by Machine Learning Algorithms



Srinivasulu Reddy, Damodar

Abstract— Over the past few decades, Machine Learning (ML) has evolved from the endeavor of few computer enthusiasts exploiting the possibility of computers learning to play games, and a part of Mathematics (Statistics) that seldom considered computational approaches, to an independent research discipline that has not only provided the necessary base for statistical-computational principles of learning procedures, but also has developed various algorithms that are regularly used for text interpretation, design recognition, and a many other commercial purposes and has led to a separate research interest in data mining to identify hidden regularities or irregularities in social data that growing by second. This paper efforts on explaining the concept and evolution of Machine Learning, some of the popular Machine Learning algorithms and try to compare three most popular algorithms based on some basic notions. Sentiment140 dataset was used and performance of each algorithm in terms of training time, prediction time and accuracy of prediction have been documented and compared.

KEYWORDS: Machine Learning, Algorithm, Data, Training, accuracy.

I. INTRODUCTION

Machine learning is a paradigm that may refer to learning from past experience (which in this case is previous data) to improve future performance. The sole focus of this field is automatic learning methods. Learning refers to modification or improvement of algorithm based on past “experiences” automatically without any external assistance from human.

A. While designing a machine (a software system), the programmer always has a specific purpose in mind.

For instance, consider J. K. Rowling’s Harry Potter Series and Robert Galbraith’s Cormoran Strike Series. To confirm the claim that it was indeed Rowling who had written those books under the name Galbraith, two experts were engaged by The London Sunday Times and using Forensic Machine Learning they were able to prove that the claim was true. They develop a machine learning algorithm and “trained” it with Rowling’s as well as other writers writing examples to seek and learn the underlying designs and then “test” the books by Galbraith. The algorithm

concluded that Rowling’s and Galbraith’s writing matched the most in several aspects.

So instead of designing an algorithm to address the problematic directly, using Machine Learning, a researcher seek an approach through which the machine, i.e., the algorithm will come up with its own solution based on the example or training data set provided to it initially.

A. MACHINE LEARNING: INTERSECTION OF STATISTICS AND COMPUTER SCIENCE

Machine Learning was the phenomenal out come when Computer Science and Statistics joined forces. Computer Science efforts on building machines that solve particular problematic, and tries to identify if problematic are solvable at all. The main approach that Statistics fundamentally employs is data inference, modelling hypotheses and measuring reliability of the conclusions.

The defining idea of Machine Learning is a little different but partially dependent on both nonetheless. Whereas Computer Science concentrate on manually programming computers, ML addresses the problematic of getting computers to re-program themselves whenever exposed to new data based on some initial learning strategies provided. On the other hand, Statistics efforts on data inference and probability, Machine Learning includes additional concerns about the feasibility and effectiveness of architectures and algorithms to process those data, compounding several learning tasks into a compact one and performance measures.

B. MACHINE LEARNING AND HUMAN LEARNING

A third research area closely related to Machine Learning is the study of human and animal brain in Neuroscience, Psychology, and related fields. The researchers proposed that how a machine could learn from experience most probably would not be significantly different than how an animal or a human mind learn with time and experience. However, the research concentrated on solving machine learning problematic using learning methods of human brain did not yield much promising output so far than the researches concerned with statistical - computational approach. This might be due to the fact that human or animal psychology remains not fully understandable to date. Regardless of these difficulties, collaboration between human learning and machine learning is increasing for machine learning is being used to explain several learning techniques seeing in human or animals. For example, machine learning method of temporal difference was proposed to explain neural signals in animal learning. It is fairly expected that this collaboration is to grow considerably in coming years.

Manuscript published on 30 September 2019.

* Correspondence Author (s)

Srinivasulu Reddy, Asst. Professor, Dept. of CSE., Malla Reddy Engineering College for Women, Hyderabad, India.

Damodar, Asst. Professor, Dept. of CSE., Malla Reddy Engineering College for Women, Hyderabad, India.

© The Authors. Published by Blue Eyes Intelligence Engineering and Sciences Publication (BEIESP). This is an [open access](https://creativecommons.org/licenses/by-nc-nd/4.0/) article under the CC-BY-NC-ND license <http://creativecommons.org/licenses/by-nc-nd/4.0/>.

C. DATA MINING, ARTIFICIAL INTELLIGENCE AND MACHINE LEARNING

In practice, these three disciplines are so intertwined and overlapping that it's almost to draw a boundary or hierarchy among the three. To put it in other words, these three fields are symbiotically related and a combination of these approaches may be used as a tactic to produce more efficient and sensitive outputs.

Roughly, Data mining is basically about interpreting any kind of data, but it lays the foundation for both artificial intelligence and machine learning. In practice, it not only sample information from various sources but it analyses and recognises design and correlations that exists in those information that would have been difficult to interpret manually. Hence, data mining is not a mere method to prove a hypothesis but method for drawing relevant hypotheses. That mined data and the corresponding designs and hypotheses may be utilised the basis for both machine learning and artificial intelligence. Artificial intelligence may be broadly defined as machines those having the ability to solve a given problematic on their own without any human intervention. The solutions are not programmed directly into the system but the necessary data and the AI interpreting that data produce a solution by itself. The interpretation that goes underneath is nothing but a data mining algorithm. Machine learning takes promote the approach to an advanced level by providing the data essential for a machine to train and modify suitably when exposed to new data. This is known as "training". It efforts on extracting information from considerably large sets of data, and then detects and identifies underlying designs using various statistical measures to improve its ability to interpret new data and produce more effective outputs. Evidently, some parameters should be "tuned" at the incipient level for better productivity. Machine learning is the foothold of artificial intelligence. It is improbable to design any machine having abilities associated with intelligence, like language or vision, to get there at once. That task would have been almost impossible to solve. Moreover, a system cannot be considered completely intelligent if it lacked the ability to learn and improve from its previous exposures.

II. PRESENT RESEARCH QUESTIONS & RELATED WORK

The Several applications mentioned earlier suggests considerable advancement so far in ML algorithms and their fundamental theory. The discipline is divulging in several direction, probing a range of learning problematic. ML is a vast discipline and over past few decades numerous researchers have added their works in this field. The enumeration of these works are countably infinite and mentioning every work is out of the scope of this paper. However this paper describes the main research questions that are being pursued at present and provide references to some of the recent notable works on that task.

A. USING UNLABELLED DATA IN SUPERVISED LEARNING

Supervised learning algorithms approximate the relation between features and labels by defining an estimator $f: X \rightarrow Y$ for a particular group of pre-labeled training data $\{x_i, y_i\}$. The main challenge in this approach is pre-labeled

data is not always readily available. So before applying Supervised Classification, data need to be preprocessed, filtered and labeled using unsupervised learning, feature extraction, dimensionality reduction etc. there by adding to the total cost. This hike in cost can be reduced effectively if the Supervised algorithm can make use of unlabelled data (e.g., images) as well. Interestingly, in many special instances of learning problematic with additional assumptions, unlabelled data can indeed be warranted to improve the expected accuracy of supervised learning. Like, consider classifying web pages or detecting spam emails. Currently active researchers are seriously taking into account new algorithms or new learning problematic to exploit unlabelled data efficiently.

B. TRANSFERRING THE LEARNING EXPERIENCE

In many real life problematic, the supervised algorithm may involve learning a family of related functions (e.g., diagnosis functions for hospitals across the globe) rather than a single function. Even if the diagnosis functions for different cities (e.g., Kolkata and London) are presumed to be relatively different, some commonalities are anticipated as well. ML algorithms like hierarchical Bayesian methods give one approach that assumes the learning parameters of both the functions, say for Kolkata and London respectively, have some common prior probabilities, and allows the data from different city hospitals to overrule relevant priors as fitting. The subtlety further increases when the transfer among the functions are compounded.

C. LINKING DIFFERENT ML ALGORITHMS

Various ML algorithms have been introduced and experimented on in a number of domains. One trail of research aims to discover the possible correlations among the existing ML algorithms, and appropriate case or scenarios to use a particular algorithm. Consider, these two supervised classification algorithms, Naive Bayes and Logistic Regression. Both of them approach many data sets distinctly, but their equivalence can be demonstrated when implemented to specific types of training data (i.e., when the criteria of Naive Bayes classifier are fulfilled, and the number of examples in training set tends to infinity). In general, the conceptual understanding of ML algorithms, their convergence features, and their respective effectiveness and limitations to date remain a radical research concern.

D. BEST STRATEGICAL APPROACH FOR LEARNERS WHICH COLLECTS THEIR OWN DATA

A border research discipline efforts on learning systems that instead of mechanically using data collected by some other means, actively collects data for its own processing and learning. The research is devoted into finding the most effective strategy to completely hand over the control to the learning algorithm. For example consider a drug testing system which try to learn the success of the drug while monitoring the exposed patients for possible unknown side effects and try to in turn minimising them.

E. PRIVACY PRESERVING DATA MINING

This approach involves successfully applying data mining and obtaining outputs without exploiting the underlying information is attracting variety of research communities and beyond. Consider, a medical diagnosis routine trained with data from hospitals all over the world. But due to privacy concerns, this kind of applications is not largely pursued. Even if this presents a cross road between data mining and data privacy, ongoing research says a system can have both. One proposed solution of the above problematic is to develop a shared learning algorithm instead of a central database. Each of the hospitals will only be allowed to employ the algorithm under pre-defined restrictions to protect the privacy of the patients and then hand it over to the next. This is an booming research domain, combining statistical exploitation of data and recent cryptographic techniques to ensure data privacy.

F. NEVER-ENDING LEARNERS

Most of the machine learning tasks entails training the learner using certain data sets, then setting aside the learner and utilise the output. Whereas, learning in humans and other animals learn continuously, adapting different skills in succession with experience, and use these learnings and abilities in a thoroughly synergistic way. Despite of sizeable commercial applications of ML algorithms, learning in machines (computers) to date has remained strikingly lacking compared to learning in human or animal. An alternative approach that more diligently capture the multiplicity, adeptness and accumulating character of learning in human, is named as never- ending learning. For instance, the Never Ending Language Learner (NELL)^[8] is a learner whose function is learning to read webpages and has been reported to read the world wide web every hour since January 2010. NELL has obtained almost 80 million confidence- weighted opinions (Example, served With (tea, biscuits)) and has been able to learn million pairs of features and parameters that capacitate it to acquire these beliefs. Furthermore, it has become competent in reading (extracting) more beliefs, and overthrow old inaccurate ones, adding to a collection of confidence and provenance for each belief and there by improving each day than the last.

III. CATEGORISATION OF ML ALGORITHMS

An overwhelming number of ML algorithm have been designed and introduced over past years. Not everyone of them are widely known. Some of them did not satisfy or solve the problematic, so another was introduced in its place. Here the algorithms are broadly grouped into two category and those two groups are further sub-divided. This section try to name most popular ML algorithms and the next section compares three most widely used ML algorithms.

IV. MEASURING AND COMPARING PERFORMANCES OF POPULAR ML ALGORITHMS

Though various researchers have contributed to ML and numerous algorithms and techniques have been introduced as mentioned earlier, if it is closely studied most of the

practical ML approach includes three main supervised algorithm or their variant. These three are namely, Naive Bayes, Support Vector Machine and Decision Tree. Majority of researchers have utilised the concept of these three, be it directly or with a boosting algorithm to enhance the efficiency further. These three algorithms are discussed briefly in the following section.

A. NAIVE BAYES CLASSIFIER

It is a supervised classification method developed using Bayes' Theorem of conditional probability with a 'Naive' assumption that every pair of feature is mutually independent. That is, in simpler words, presence of a feature is not effected by presence of another by any means. Irrespective of this over-simplified assumption, NB classifiers performed quite well in many practical situations, like in text classification and spam detection. Only a small amount of training data is needed to estimate certain parameters. Beside, NB classifiers have considerably outperformed even highly advanced classification techniques.

B. SUPPORT VECTOR MACHINE

SVM, another supervised classification algorithm proposed by Vapnik in 1960s have recently attracted a major attention of researchers. The simple geometrical explanation of this approach involves determining an optimal separating plane or hyperplane that separates the two classes or clusters of data points justly and is equidistant from both of them. SVM was defined at first for linear distribution of data points. Later, the kernel function was introduced to tackle non- linear data as well.

C. DECISION TREE

A classification tree, popularly known as decision tree is one of the most successful supervised learning algorithm. It constructs a graph or tree that employs branching technique to demonstrate every probable output of a decision. In a decision tree representation, every internal node tests a feature, each branch corresponds to outcome of the parent node and every leaf finally assigns the class label. To classify an instance, a top-down approach is applied starting at the root of the tree. For a certain feature or node, the branch concurring to the value of the data point for that attribute is considered till a leaf is reached or a label is decided. Now, the performances of these three were roughly compared using a set of tweets with labels positive, negative and neutral. The raw tweets were taken from Sentiment140 data set. Then those are pre-processed and labeled using a python program. Each of these classifier were exposed to same data. Same algorithm of feature selection, dimensionality reduction and k-fold validation were employed in each cases. The algorithms were compared based on the training time, prediction time and accuracy of the prediction. The experimental output is given below.

V. FUTURE SCOPE

Machine learning is research area that has attracted a lot of brilliant minds and it has the potential to divulge further. But the three most important future sub-problematic are chosen to be discussed here.

A. EXPLAINING HUMAN LEARNING

A mentioned earlier, machine learning theories have been perceived fitting to comprehend features of learning in humans and animals. Reinforcement learning algorithms estimate the dopaminergic neurones induced activities in animals during reward-based learning with surprising accuracy. ML algorithms for uncovering sporadic delineations of naturally appearing images predict visual features detected in animals' initial visual cortex. Nevertheless, the important drivers in human or animal learning like stimulation, horror, urgency, hunger, instinctive actions and learning by trial and error over numerous time scales, are not yet taken into account in ML algorithms. This a potential opportunity to discover a more generalised concept of learning that entails both animals and machine.

B. PROGRAMMING LANGUAGES CONTAINING MACHINE LEARNING PRIMITIVES

In majority of applications, ML algorithms are incorporated with manually coded programs as part of an application software. The need of a new programming language that is self-sufficient to support manually written subroutines as well as those defined as "to be learned." It could enable the coder to define a set of inputs-outputs of every "to be learned" program and opt for an algorithm from the group of basic learning methods already imparted in the language. Programming languages like Python (Sckit-learn), R etc. already making use of this concept in smaller scope. But a fascinating new question is raised as to develop a model to define relevant learning experience for each subroutines tagged as "to be learned", timing, and security in case of any unforeseen modification to the program's function.

C. PERCEPTION

A generalised concept of computer perception that can link ML algorithms which are used in numerous form of computer perception today including but not limited to highly advanced vision, speech recognition etc., is another potential research area. One thought-provoking problematic is the integration of different senses (e.g., sight, hear, touch etc) to prepare a system which employ self-supervised learning to estimate one sensory knowledge using the others. Researches in developmental psychology have noted more effective learning in humans when various input modalities are supplied, and studies on co-training methods insinuates similar outputs.

VI. CONCLUSION

The foremost target of ML researchers is to design more efficient (in terms of both time and space) and practical general purpose learning methods that can perform better over a widespread domain. In the context of ML, the efficiency with which a method utilises data resources that is also an important performance paradigm along with time and space complexity. Higher accuracy of prediction and humanly interpretable prediction rules are also of high importance. Being completely data-driven and having the ability to examine a large amount of data in smaller intervals of time, ML algorithms has an edge over manual or direct programming. Also they are often more accurate and not prone to human bias. Consider the following scenarios:

Development of a software to solve perception tasks using sensors, like speech recognition, computer vision etc. It is easy for anyone to label an image of a letter by the alphabet it denotes, but designing an algorithm to perform this task is difficult. Customisation of a software according to the environment it is deployed to. Consider, speech recognition softwares that has to be customised according to the needs of the customer. Like e-commerce sites that customises the products displayed according to customers or email reader that enables spam detection as per user preferences. Direct programming lacks the ability to adapt when exposed to different environment. ML provides a software the flexibility and adaptability when necessary. In spite of some application (e.g., to write matrix multiplication programs) where ML may fail to be beneficial, with increase of data resources and increasing demand in personalised customisable software, ML will thrive in near future. Besides software development, ML will probably help reform the general outlook of Computer Science. By changing the defining question from "how to program a computer" to "how to empower it to program itself," ML pries the development of devices that are self-monitoring, self-diagnosing and self-repairing, and the utilises of the data flow available within the program rather than just processing it. Likewise, it will help reform Statistical rules, by providing more computational stance. Obviously, both Statistics and Computer Science will also embellish ML as they develop and contribute more advanced theories to modify the way of learning.

REFERENCES

1. T. M. Mitchell, Machine Learning, McGraw-Hill International, 1997.
2. T.M. Mitchell, The Discipline of Machine Learning, CMU-ML-06-108, 2006
3. N. Cristianini and J. Shawe-Taylor. An Introduction to Support Vector Machines. Cambridge University Press, 2000.
4. E. Osuna, R. Freund, and F. Girosi. Support vector machines: training and applications. AI Memo 1602, MIT, May 1997.
5. V. Vapnik. Statistical Learning Theory. John Wiley & Sons, 1998.
6. C.J.C. Burges. A tutorial on support vector machines for design recognition. Data Mining and Knowledge Discovery, 2(2):1-47, 1998.
7. Taiwo Oladipupo Ayodele, Types of Machine Learning Algorithms, New Advances in Machine Learning, Yagang Zhang (Ed.), InTech, 2010
8. T. Mitchell, W. Cohen, E. Hruschka, P. Talukdar, J. Betteridge, A. Carlson, B. Dalvi, M. Gardner, B. Kisiel, J. Krishnamurthy, N. Lao, K. Mazaitis, T. Mohamed, N. Nakashole, E. Platanios, A. Ritter, M. Samadi, B. Settles, R. Wang, D. Wijaya, A. Gupta, X. Chen, A. Saparov, M. Greaves, J. Welling, Never-Ending Learning, Proceedings of the Twenty-Ninth AAAI Conference on Artificial Intelligence, 2014
9. Pedregosa et al., Scikit-learn: Machine Learning in Python, JMLR 12, pp. 2825-2830, 2011.
10. Wang, J. and Jebara, T. and Chang, S.-F. Semi-supervised learning using greedy max-cut. Journal of Machine Learning Research, Volume 14(1), 771-800 2013
11. Chapelle, O. and Sindhvani, V. and Keerthi, S. S. Optimization Techniques for Semi-Supervised Support Vector Machines, Journal of Machine Learning Research, Volume 9, 203-233, 2013
12. J. Baxter. A model of inductive bias learning. Journal of Artificial Intelligence Research, 12:149-198, 2000.
13. S. Ben-David and R. Schuller. Exploiting task relatedness for multiple task learning. In Conference on Learning Theory, 2003.
14. W. Dai, G. Xue, Q. Yang, and Y. Yu, Transferring Naive Bayes classifiers for text classification. AAAI Conference on Artificial Intelligence, 2007.

15. H. Hlynsson. Transfer learning using the minimum description length principle with a decision tree application. Master's thesis, University of Amsterdam, 2007.
16. Z. Marx, M. Rosenstein, L. Kaelbling, and T. Dietterich. Transfer learning with an ensemble of background tasks. In NIPS Workshop on Transfer Learning, 2005.
17. R Conway and D Strip. Selective partial access to a database. In Proceedings of ACM Annual Conference, 85 - 89, 1976
18. P D Stachour and B M Thuraisingham Design of LDV A multilevel secure relational databasemanagement system, IEEE Trans. Knowledge and Data Eng., Volume 2, Issue 2, 190 - 209, 1990
19. R Oppliger, Internet security: Firewalls and beyond, Comm. ACM, Volume 40, Issue 5, 92 -102, 1997