

ROLE OF MATHEMATICS IN MACHINE LEARNING

Dr. Suresh Dara^{*1}, Subham Surmara^{*2}, Sai Kiran Reddy^{*3},

Kamala Vani^{*4}, Aditya Sai Varma^{*5}

^{*1}Professor, Computer Science And Engineering, B.V Raju Institute Of Technology, Narsapur,
Telangana, India.

^{*2,3,4,5}Student, Computer Science And Engineering, B.V Raju Institute Of Technology, Narsapur,
Telangana, India.

ABSTRACT

All Machine Learning algorithms are built on a mathematical foundation. Because Deep Learning is a subset of Machine Learning, the above holds true for Deep Learning, Shallow Learning, Optimization, and all other Data Science methods. These algorithms assist us in extracting information from the data. We write these algorithms in a programming language (typically libraries are available), and the computer machine that executes them on the data set seems to be intelligent, thus the title Artificial Intelligence. Machine learning provides a model that can learn from data and make predictions using an algorithm. It's used to figure out how something works and why one model is superior to another. Machine learning comes with a built-in mathematical stipulation. It is a field that combines probability, statistics, linear algebra, computer science, and algorithms to develop intelligent software. These programs can extract relevant and insightful information from data in order to arrive at business insights. Because machine learning is based on the study and application of algorithms, a strong foundation in mathematics is required.

Keywords: Mathematics Of Machine Learning, Statistics, Calculus, Linear Algebra, Probability, Computer Science, Deep Learning, Artificial Intelligence.

I. INTRODUCTION

Many people are aiming to transfer to the AI/ML/Data Science area these days, which is quite encouraging and in line with the world's changing speed. However, these individuals are perplexed by questions such as: I want to be a machine learning expert without having to learn a lot of math. Is that possible? What role does mathematics play in data science and AI/ML? As previously said, there are numerous libraries available to conduct various machine learning tasks, making it simple to ignore the mathematical aspects of the topic. Various issues, such as computer games, self-driving automobiles, and object recognition, are difficult to tackle with traditional programming methods. Machine learning is one approach to tell computers how to learn from data. Machine Learning is used to assist Amazon propose products to you, YouTube recommend videos, and spam mail be classified, among other things. To achieve this, we use a combination of mathematics and a lot of programming. The goal of machine learning is to create algorithms that can learn from data and generate predictions [1]. Machine learning is predicated on mathematical foundations. Mathematics is required to complete the Data Science project and to solve the Deep Learning use cases. Mathematics clarifies the basic principle of the algorithms and explains why one is superior to the other. You can develop models even if you don't understand the logic behind how algorithms function, but wait... What would you do if you didn't know which one was best and when to utilise it? To work as a data scientist, you must be familiar with the mathematics that underpin machine learning techniques. It's unavoidable. Every recruiter and seasoned machine learning specialist would attest to the fact that it is a crucial aspect of a data scientist's job.

II. MACHINE LEARNING PROCESS

Machine learning is the process of creating systems that are particularly engineered to learn and develop on their own. Machine learning's goal is to create algorithms that automatically assist a system in gathering data and using that data to learn more. Systems are required to search for patterns in the data they acquire and use those patterns to make critical decisions for themselves. Machine learning, in general, is the process of teaching computers to think and act like people, to demonstrate human-like intelligence, and to give them a brain. The goal of imbuing computers with intelligence appears overwhelming and unachievable. FIGURE 1 illustrates the major phases followed for proposing a model for a particular problem via ML [2].

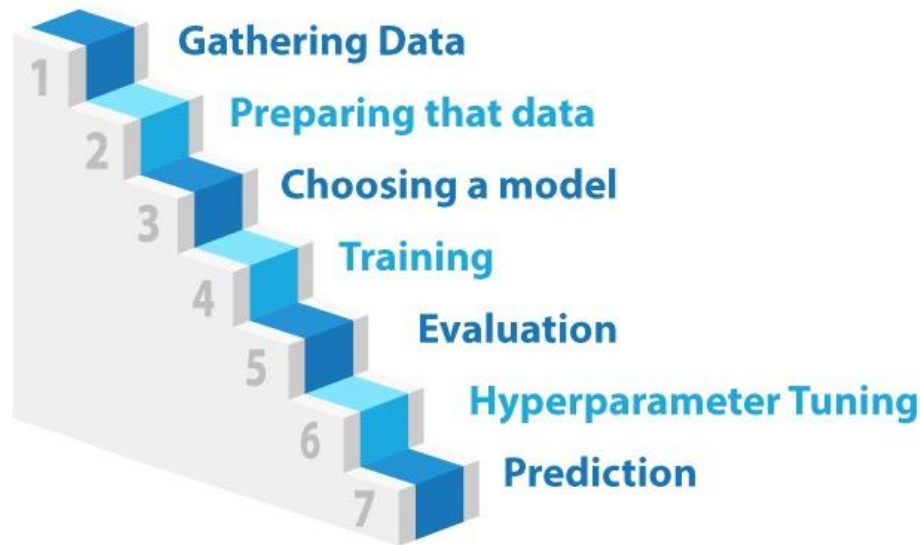


Fig 1: Machine Learning flow

STEP 1: Collecting Data:

Machines, as you may know, learn from the data you provide them with. It's critical to get trustworthy data so that your machine learning model can uncover the right trends. The accuracy of your model is determined by the quality of the data you provide the machine. If your data is faulty or obsolete, you'll get inaccurate results or forecasts that aren't useful. Make sure you utilize data from a reputable source, as it will have a direct impact on the model's conclusion. Good data is meaningful, has few missing and duplicated numbers, and accurately represents the many subcategories/classes[17].

STEP 2: Data Pre-Processing:

Putting all your info together and randomizing it. This ensures that data is dispersed uniformly and that the ordering has no effect on the learning process. Unwanted data, missing values, rows, and columns, duplicate values, data type conversion, and so on are all removed from the data [18]. It's possible that you'll need to rearrange the dataset and alter the rows and columns, as well as the indexes of rows and columns. Visualize the data to see how it's organized and to see the connections between the various variables and classes. Creating two sets of cleansed data: a training set and a testing set. The training set is the one from which your model learns. A testing set is used to assess your model's correctness after it has been trained.

STEP 3: Choosing a Model

After performing a machine learning algorithm on the obtained data, a machine learning model selects the output. It is critical to select a model that is appropriate for the work at hand. Over time, scientists and engineers have built a variety of models for diverse tasks such as speech recognition, picture recognition, prediction, and so on. Aside from that, you must determine if your model is best suited for numerical or categorical data and make the appropriate choice[19].

STEP 4: Training the Model

The most crucial phase in machine learning is training. To detect patterns and create predictions, you give the prepared data to your machine learning model during training[20]. Consequently, the model learns from the data and can complete the goal assigned. The model improves in predicting over time as it is trained.

STEP 5: Evaluating the Model

After you've trained your model, you'll want to see how it's doing. This is accomplished by putting the model to the test on previously unknown data. The testing set that you split our data into before is the unseen data utilized. If you test on the same data that was used for training, you won't receive an accurate result since the model is already familiar with the data and recognizes the same patterns it did before. This will provide you with a disproportionately high level of precision. When applied on testing data, you can receive a precise estimate of how your model will perform and how fast it will run[21].

STEP 6: Parameter Tuning

Examine whether your model's accuracy can be enhanced in any manner once you've constructed and tested it. This is accomplished by fine-tuning the parameters in your model. Parameters are the variables in the model that are set by the programmer. The accuracy will be at its highest for a certain value of your parameter[22]. Finding these settings is referred to as parameter tweaking.

STEP 7: Making Predictions

Finally, you'll be able to generate accurate predictions using your model on previously unknown data.

III. MATHEMATICS USED IN MACHINE LEARNING

Most of our real-world business problems are solved using four pillars of Machine Learning. These pillars are also used in many Machine Learning techniques. They really are.

- Statistics
- Probability
- Linear Algebra
- Calculus

FIGURE 2 shows the components of math used in ML [3].

Machine learning is all about dealing with data. We collect data from organizations or from any repositories like Kaggle, UCI etc., and perform various operations on the dataset like cleaning and processing the data, visualizing and predicting the output of the data. For all the operations we perform on data, there is one common foundation that helps us achieve all of this through computation-- and that is Math

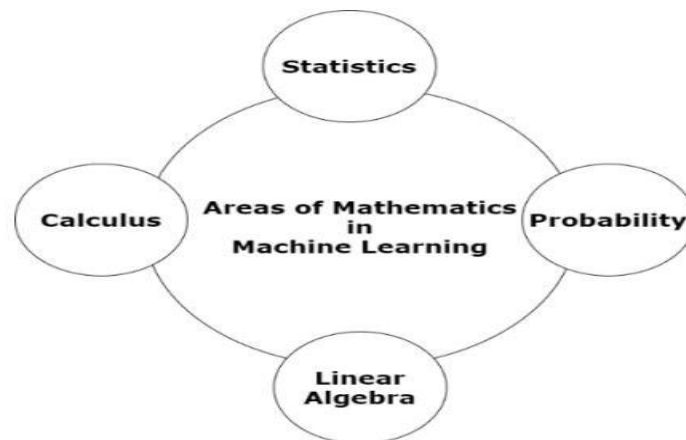


Fig 2: Fields of Maths in ML

Statistics:

Statistics is the core of everything [4]. It's utilized to draw conclusions based on evidence. It is concerned with statistical approaches for gathering, presenting, analyzing, and interpreting numerical data. Statistics are significant in the field of Machine Learning since they work with vast volumes of data and are a crucial aspect in an organization's growth and development. Data can be gathered via censuses, samples, primary and secondary data sources, and more. This stage assists us in identifying our objectives so that we may work on the next phases. The data obtained contains noise, incorrect data, null values, outliers, and other anomalies. We need to clean up the data and turn it into actionable findings. The information should be provided in an appropriate and succinct manner. It is one of the most important processes since it aids in the comprehension of the insights and serves as the foundation for additional data analysis. Condensation, Summarization, and Conclusion are examples of data analysis techniques that use central tendencies, dispersion, skewness, Kurtosis, co-relation, regression, and other techniques. Because the statistics do not speak for themselves, the interpretation process comprises making inferences from the data acquired. Factual learning is at the heart of any great machine learning. [8] Regression was used as an example of a unique factual strategy [9]. Based on the kind of analysis they conduct on the data, statistics utilized in Machine Learning may be split into two groups. Descriptive statistics and inferential statistics are the two types of statistics.

a) Descriptive Statistics:

- Interested in characterizing and summarizing the target demographic
- It only works with a tiny set of data.
- Pictorial representations are used to demonstrate the outcome.
- Mean, Median, and Mode are Central and Range measurements, whereas Standard Deviation, variance, and other measures of Variability are employed in Descriptive Statistics.

b) Inferential Statistics:

- Methods for making population-level choices or predictions based on sample data.
- It can handle a large amount of data.
- Compares, tests, and forecasts future results.
- The likelihood scores represent the results.
- Inferential statistics are unique in that it draws inferences about a population based on information other than the data provided.
- Inferential statistics use procedures such as hypothesis testing, sampling distributions, and analysis of variance (ANOVA).
- Machine Learning Algorithms rely heavily on statistics. A Data Analyst's job in the industry is to derive conclusions from data, and he or she needs and is reliant on statistics to do so.

Probability:

The term probability refers to the possibility of a specific event occurring based on previous experiences. It is used in the field of Machine Learning to forecast the likelihood of future events.

An event's probability is calculated as:

- $P(\text{Event}) = \text{Favorable Outcomes} / \text{Total Number of Possible Outcomes}$

An event is a group of results from an experiment in the realm of probability. The probability of an event occurring is represented by $P(E)$, and E is referred to as an Event. Any occurrence has a probability ranging from 0 to 1. A Trial is a condition in which the event E may or may not occur.

Some of the basic concepts required in probability are as follows

- Joint Probability: $P(A \cap B) = P(A) \cdot P(B)$, Only when events A and B are independent of one another is this form of probability conceivable.
- Conditional Probability: When it is known that another event B has already occurred, it is the chance of event A occurring and is denoted by $P(A|B)$ i.e., $P(A|B) = P(A \cap B) / P(B)$
- Bayes theorem: It is defined as the application of probability theory results, which entails estimating unknown probabilities and making judgments based on fresh sample data. In the existence of additional data, it is beneficial for addressing business challenges. The popularity of this theorem stems from its use in modifying an existing set of probabilities (Prior Probability) with new knowledge and generating a new set of probabilities (Posterior Probability).

The link between the Conditional Probabilities of occurrences is explained by Bayes theorem." This theorem is useful in determining the 'Specificity' and 'Sensitivity' of data, and it works best with uncertainty samples of data. This theorem is crucial in constructing the CONFUSION MATRIX.

The performance of Machine Learning Models or Algorithms that we design is measured using a confusion matrix, which is a table-like structure. This is useful for calculating True Positive Rates, True Negative Rates, False Positive Rates, False Negative Rates, Precision, Recall, F1-score, Accuracy, and Specificity when constructing the ROC Curve from provided data.

We need to pay more attention to probability distributions, which are divided into Discrete and Continuous types, as well as Likelihood Estimation Functions. The Naive Bayes Algorithm is a probabilistic machine learning algorithm that assumes input characteristics are independent [6]. Different learning methods, such as Nave Bayes [6] and Bayesian Networks [7], are based on probability.

Calculus:

This is a field of mathematics that aids in the study of quantity change rates. It is concerned with improving the performance of machine learning algorithms or models. It is impossible to compute probabilities on data without comprehending calculus, and we cannot extrapolate plausible outcomes from the data we collect.

Integrals, limits, derivatives, and functions are the basic topics of calculus. Differential statistics and inferential statistics are the two forms of statistics. It is used to train deep Neural Networks using back propagation methods. Differential Calculus breaks down the provided data into little chunks in order to determine how it changes. To figure out how much there is, Inferential Calculus combines (joins) the little bits. Calculus is mostly used to improve the performance of Machine Learning and Deep Learning algorithms. It is utilized to provide quick and effective solutions. Calculus is utilized in algorithms such as Gradient Descent and Stochastic Gradient Descent (SGD), as well as optimizers like Adam, Rms Drop, and Adadelta. Calculus is primarily used by Data Scientists in the development of numerous Deep Learning and Machine Learning models. They assist in the optimization of data and the production of superior data outputs by extracting intelligent insights from it. Calculus may be utilised to implement pattern learning. For analysis, the ML model employs a variety of state and control combinations [11], [12].

Linear Algebra:

The emphasis of Linear Algebra is on computation. It is utilized for Deep Learning and plays an important part in comprehending the background theory of Machine Learning. It offers us a better understanding of how algorithms function in real life and allows us to make better judgments. It mostly focuses on vectors and matrices.

- A single integer is referred to as a scalar.
- A vector is a numerical array that is expressed in a row or column and has just one index to access it (i.e., either Rows or Columns)
- A matrix is a two-dimensional array of integers that may be accessed using both indices and keys (i.e., by both rows and columns)
- A tensor is a set of integers that are arranged in a grid in a certain order and have a variable number of axes.

All these numerical operations on the data are computed using the Python library's Numpy module. The Numpy library performs fundamental operations on vectors and matrices, such as addition, subtraction, multiplication, and division, and returns a meaningful value. The N-d array is how Numpy is expressed. Without Linear Algebra, machine learning models cannot be constructed, complicated data structures cannot be managed, and matrices operations cannot be done. Linear Algebra is used as a platform to present all the model outcomes. Linear Algebra is used in the development of several Machine Learning algorithms such as Linear, Logistic Regression, SVM, and Decision Trees. We may also create our own ML algorithms using Linear Algebra. When dealing with data, Data Scientists and Machine Learning Engineers use Linear Algebra to create their own algorithms[5].

IV. CONCLUSION

Mathematics is a critical area to focus on for machine learning enthusiasts and aspirants, and it is necessary to have a good foundation in Math. Every notion you learn in Machine Learning, every tiny algorithm you build or use to solve a problem has a direct or indirect mathematical connection. The arithmetic ideas used in machine learning are based on the fundamental math that we study in 11th and 12th grades. We obtain theoretical knowledge at that point, but in the field of Machine Learning, we get to encounter the practical applications of arithmetic that we studied before. Taking a Machine Learning Algorithm, finding a use case, then solving and understanding the math behind it is the greatest approach to become familiar with the ideas of Mathematics. To come up with machine learning solutions to real-world issues, we need to have a solid knowledge of math. A solid understanding of arithmetic principles also aids in the development of problem-solving abilities.

V. REFERENCES

- [1] Srinivas Pyda, Srinivas Kareenhalli, "Mathematics and Machine Learning, International Conference on Mathematics and Computing", 2018, pp. 135-153
- [2] J. Clerk Maxwell, A Treatise on Electricity and Magnetism, 3rd ed., vol. 2. Oxford: Clarendon, 1892, pp.68-73.
- [3] Sanjeev Great learning team (Jan 19,2022) What is machine learning ? How machine learning works and future of it ? <https://www.mygreatlearning.com/blog/what-is-machine-learning/>
- [4] Md. Kosher, "A combination of Mathematics, Statistics, and Machine Learning to Detect Fraud, National Mathematics Conference", Bangladesh, 2020.
- [5] Gennady Grabarnik, Luiza Kim-Tyan, Serge Yaskolko, "Addressing Prerequisite for STEM Classes

- Using an Example of Linear Algebra for a Course in Machine Learning”, The Twelfth International Conference on Mobile, Hybrid, and On-line Learning, 2020, pp. 21_26
- [6] Somya Goel, Sanjana Rosshan, Rishabh Tyagi, Sakshi Agarwal, “Augur Justice: A Supervised Machine Learning to predict Outcomes of Divorce court cases”, Fifth International Conference on Image Information Processing”, 2019, pp. 280-285.
- [7] Kevin Fong-Rey Liu, Jia-Shen Chen,” Prediction and assessment of student learning outcomes in calculus A decision support of integrating data mining and Bayesian belief networks”, 3rd International Conference on Computer Research and Development,2011, IEEE.
- [8] Aaron N. Richter,b,*, Taghi M. Khoshgoftaar, “A review of statistical and machine learning methods for modeling cancer risk using structured clinical data”, Artificial Intelligence In Medicine, vol 90, 2018, pp. 1_14
- [9] Zhenru Wang, Tijie Shi, “ Prediction of the Admission Lines of College Entrance Examination based on machine learning”, 2nd IEEE International Conference on Computer and Communications”, 2016, pp. 332-335
- [10] N. Milosevic, A. Dehghantanha, and K-K. R. Choo, “Machine learning aided Android malware classification”, Comput. Elect. Eng. , vol. 61, 2017, pp. 266_274
- [11] Mohan S Acharya, Asfia Armaan, Aneeta S Antony, “A Comparison of Regression Models for Prediction of Graduate Admissions”, Second International Conference on Computational Intelligence in Data Science, 2019, IEEE.
- [12] Le, DN., Parvathy, V.S., Gupta, D. et al. IoT enabled depthwise separable convolution neural network with deep support vector machine for COVID-19 diagnosis and classification. Int. J. Mach. Learn. & Cyber. (2021). <https://doi.org/10.1007/s13042-020-01248-7>
- [13] Anupama, C.S.S., Sivaram, M., Lydia, E.L. et al. Synergic deep learning model-based automated detection and classification of brain intracranial hemorrhage images in wearable networks. Pers Ubiquit Comput (2020). <https://doi.org/10.1007/s00779-020-01492-2>
- [14] Shankar, K., Sait, A. R. W., Gupta, D., Lakshmanaprabu, S. K., Khanna, A., & Pandey, H. M. (2020). Automated detection and classification of fundus diabetic retinopathy images using synergic deep learning model. Pattern Recognition Letters, 133, 210-216
- [15] Ghosh, S., Rana, A. and Kansal, V., “A Statistical Comparison for Evaluating the Effectiveness of Linear and Nonlinear Manifold Detection Techniques for Software Defect Prediction” International Journal of Advanced Intelligence Paradigms (IJAIP), 12(3/4), pp 370- 391,ISSN 1755-0394, Inderscience, DOI 10.1504/IJAIP.2019.098578, 2019
- [16] Apache 2.0 (2009) <https://spark.apache.org/docs/latest/index.html>
- [17] Yuji Roh, Geon Heo, Steven Euijong Whang,”A Survey on Data Collection for Machine Learning ” , 2019
- [18] Vivek Agarwal ,”Research on Data Preprocessing and Categorization Technique for Smartphone Review Analysis”,2015.
- [19] Sebastian Raschka,”Model Evaluation, Model Selection, and Algorithm Selection in Machine Learning“, 2020.
- [20] Dr. S. Veena,T. Shankari,S. Sowmiya,M. Varsha,”A SURVEY ON TOOLS USED FOR MACHINE LEARNING ”,2020.
- [21] Jianlong Zhou , Amir H. Gandomi, Fang Chen and Andreas Holzinger ,”Evaluating the Quality of Machine Learning Explanations:A Survey on Methods and Metrics”,2021.
- [22] Tong Yu, Hong Zhu,”Hyper-Parameter Optimization: A Review of Algorithms and Applications”,2020.