# UTILIZING MACHINE LEARNING ALGORITHMS FOR PREDICTIVE ANALYSIS OF BIG MART SALES

**[1] SHAIK.SANA, [2] MRS. CH. DEEPTI**

[1] PG Scholar in the department of MCA at QIS College of Engineering & Technology (AUTONOMOUS), Vengamukkapalem, Ongole- 523272, Prakasam Dt., AP., India.

[2] Professor in the department of CSE/MCA at QIS College of Engineering & Technology (AUTONOMOUS), Vengamukkapalem, Ongole- 523272, Prakasam Dt., AP., India.

**ABSTRACT :**
The utilization of machine learning algorithms for predictive analysis has become increasingly prevalent, particularly in the domain of retail sales forecasting. In this study, we explore the application of machine learning techniques to predict sales in a large retail chain, specifically Big Mart. Leveraging a comprehensive dataset containing information on product attributes, store demographics, and historical sales records, we employ various machine learning algorithms, including linear regression, decision trees, random forest regressor,hyper parameter tuning and XG boost regressor, to build predictive models. Through rigorous experimentation and evaluation, we assess the performance of these models in accurately forecasting sales for different products and stores. Our findings demonstrate the efficacy of machine learning algorithms in capturing complex patterns and relationships within the data, thereby enabling more accurate and reliable sales predictions for Big Mart and other retail chains.
**INDEX:**Machine Learning Algorithms,Big Mart Sales,Linear Regression, Random Forest,XG Boost,Predictive Analysis

## INTRODUCTION

With the rapid development of global malls and stores chains and the increase in the number of electronic payment customers, the competition among the rival organizations is becoming more serious day by day. Each organization is trying to attract more customers using personalized and short-time offers which makes the prediction of future volume 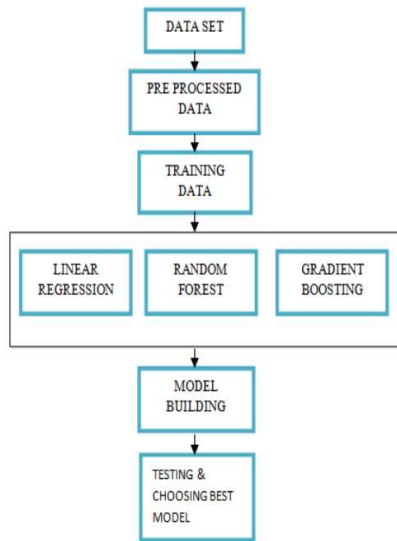of sales of every item an important asset in the planning and inventory management of every organization, transport service, etc. Due to the cheap availability of computing and storage, it has become possible to use sophisticated machine learning algorithms for this purpose. In this paper, we are providing forecast for the sales data of big mart in a number of big mart stores across various location types which is based on the historical data of sales volume. According to the characteristics of the data, we can use the method of multiple linear regression analysis and random forest to forecast the sales volume.

## PROBLEM STATEMENT :

The predictive analysis of sales in retail settings, particularly for large retail chains like Big Mart, poses significant challenges due to the complexity and variability of factors influencing sales outcomes. Traditional statistical methods often struggle to capture the intricate patterns and relationships present in large and heterogeneous datasets comprising product attributes, store demographics, and historical sales records. Moreover, manual forecasting processes are labor- intensive and time-consuming, limiting their scalability and efficiency in dynamic retail environments. Thus, there is a pressing need for advanced analytical techniques that can leverage the wealth of data available to retailers and provide accurate and timely sales predictions. Utilizing machine learning algorithms for predictive analysis of Big Mart sales presents an opportunity to address these challenges and unlock valuable insights to optimize inventory management, marketing

strategies, and overall business performance. However, deploying machine learning models in real-world retail settings requires addressing varioustechnical and practical considerations, including data preprocessing, feature selection, model interpretability, and scalability, to ensure the effectiveness and usability of predictive analytics solutions

## SYSTEMARCHITECTURE:



## METHODOLOGY :

The steps followed in this task, beginning from the dataset preparation to obtaining results are represented in Fig.1.



**Fig1:** Steps followed for obtaining results
Pre-processing of Dataset:
 Big Mart''s data scientists have collected sales data of their 10 stores established at different locations with each store having 1559 different products as per 2013 data collection. Using all the observations it is deduced what role certain properties of an item play and how they cloud affect their sales. The dataset is displayed in Fig.2 on using head() function on the dataset variable.



**Fig2:** Screenshot of Dataset
The data set consists of various data types such as integer, float and, object as shown in Fig.3.



Fig3: Various datatypes used in the Dataset In the raw data, there could be various types of underlying patterns which also gives deeper knowledge about subject of interests and provides useful insights about the problem. But caution should be observed while dealing with the data as it may contain null values, or redundant values, or ambiguity values, which also demands for pre-processing of data. Therefore data exploration becomes mandatory. Various factors that are important by statistical means like mean, standard deviation, median, count of values and maximum value etc. are shown in Fig.4 based on the numerical variables of our dataset.

**Fig4:** Numerical variables of the Dataset

Pre-processing of this dataset involves analysis on the independent variables like checking for null values ineach column and then replacing or feeding supported appropriate data types, so that analysis and model fitting is carried out its way to accuracy. Shown above are some of the representations that are obtained using Pandas tools which gives information about variable count for numerical columns and model values for categorical columns. Maximum and minimum values in numerical columns, along with their percentile values for median, plays an important factor in deciding which value will be chosen at priority for future data exploration tasks and analysis. Data types of different columns are further used in label processing and one-shot encoding scheme during model building.

Models and Metrics:

Three determining models were chosen to fit the data: Linear Regression, Random Forest, and Gradient Boosting. All models came from scikit-learn library in Python and were executed in Python 3.7. The metrics used to determine the performance of models were Root Mean Square Error (RMSE), Mean Absolute Error (MAE), and the coefficient of determination ($R^2$).

$$\text{RMSE} = \sqrt{\frac{1}{n}\sum_{i=1}^{n}(x_i - \bar{x})^2}$$

$$\text{MAE} = \frac{1}{n}\sum_{i=1}^{n}|x_i - \bar{x}|$$

$$R^2 = 1 - \frac{\sum_{i=1}^{n}(y_i - \hat{y}_1)^2}{\sum_{i=1}^{n}(y_i - \bar{y})^2}$$

Linear regression model with ordinary least squares (OLS) method was used. This model assumes that the relationship between the dependent variable y and the independent variable x is linear, and generates the set of coefficients $\beta i$ 's which minimizes the residual sum of squares between the observed values in dataset and the targets predicted by the model. Here, the intercept term was included to achieve better result. The formular is as following:

$$y = \beta 0 + \beta 1x1 + \beta 2x2 + \beta 3x3 + \cdots$$

Random Forest is an ensemble learning method that operates by building a number of decision trees at training time and uses averaging to improve prediction accuracy and reduce overfitting. For regression tasks, the average predicted values of individual trees are returned.

Gradient Boosting is a method of combining several simple models, which are typically decision trees, into a composite model. It has a forward stage-wise fashion since simple models are added one by one and each new model takes a step in the direction minimizing the prediction error. As more simple models are combined, the final completed model gets stronger. This algorithm uses gradient descent to minimize losses, which is where the term "gradient" comes from. When training the model, learning rate and the number of boosting stages were set to 0.1 and 100, respectively.

**RESULTANALYSIS :**

## CONCLUSION :

In conclusion, the utilization of machine learning algorithms for predictive analysis of Big Mart sales presents a transformative opportunity for the retail industry. Through the application of advanced analytical techniques and leveraging vast amounts of data, retailers like Big Mart can gain valuable insights into sales patterns, customer behavior, and market trends. By accurately forecasting sales across product categories and store locations, Big Mart can optimize inventory management, pricing strategies, and promotional activities, leading to improved operational efficiency and profitability. The objective of this framework is to predict the future sales from given data of the previous year's using machine Learning techniques. how different machine learning models are built using different algorithms like Linear regression, Random forest regressor, and XG booster algorithms. These algorithms have been applied to predict the final result of sales. We have addressed in detail about how the noisy data is been removed and the algorithms used to predict the result. Based on the accuracy predicted by different models we conclude that the random forest approach is the best models. Our predictions help big marts to refine their methodologies and strategies which in turn helps them to increase their profit.

## FUTURE ENHANCEMENT :

The future of utilizing machine learning algorithms for predictive analysis of Big Mart salesholds several exciting prospects for further research and development. To begin, in order to make sales projections that are more accurate and resilient, we need to investigate more sophisticated machine learning approaches like reinforcement learning algorithms and deep learning frameworks.

Sales trend forecasting in Big Mart's ever-changing retail environment might benefit from deep learning models, such as convolutional neural networks (CNNs) and transformer-based architectures, which have shown potential in capturing intricate patterns and connections in large-scale data. Big Mart and other comparable retail chains may benefit from predictive analytics solutions that are more agile and responsive by developing adaptive forecasting models that learn and adapt to changing market circumstances. This can be achieved by employing reinforcement learning methods.Moreover, future research efforts should focus on integrating external data sources and incorporating advanced analytics techniques, such as sentiment analysis and geospatial analytics, into sales prediction models. By leveraging data from social media, weather forecasts, and economic indicators, retailers like Big Mart can gain deeper insights into consumer behavior and external factors influencing sales patterns. Furthermore, the integration of advanced analytics techniques, such as sentiment analysis, can help retailers understand customer preferences and sentiments, enabling personalized marketing strategies and product recommendations. Additionally, geospatial analytics can provide valuable insights into regional variations in demand and help optimize store locations and distribution networks for Big Mart,ultimately leading to improved sales performanceandcustomer satisfaction.

## REFERENCE :

1.    Srikanth veldandi, et al. "Design and Implementation of Robotic Arm for Pick and Place by using Bluetooth Technology." Journal of Energy Engineering and Thermodynamics, no. 34, June 2023, pp. 16–21. https://doi.org/10.55529/jeet.34.16.21.
2.    Srikanth veldandi., et al. "Grid Synchronization Failure Detection on Sensing the Frequency and Voltage beyond the Ranges." Journal of Energy Engineering and

Thermodynamics, no. 35, Aug. 2023, pp. 1–7. https://doi.org/10.55529/jeet.35.1.7.

3. Srikanth veldandi, et al. "Intelligents Traffic Light Controller for Ambulance." Journal of Image Processing and Intelligent Remote Sensing, no. 34, July 2023, pp. 19–26. https://doi.org/10.55529/jipirs.34.19.26.

4. Srikanth veldandi, et al. "Smart Helmet with Alcohol Sensing and Bike Authentication for Riders." Journal of Energy Engineering and Thermodynamics, no. 23, Apr. 2022, pp. 1–7. https://doi.org/10.55529/jeet.23.1.7.

5. Srikanth veldandi, et al. "An Implementation of Iot Based Electrical Device Surveillance and Control using Sensor System." Journal of Energy Engineering and Thermodynamics, no. 25, Sept. 2022, pp. 33–41. https://doi.org/10.55529/jeet.25.33.41.

6. Srikanth veldandi, et al "Design and Implementation of Robotic Arm for Pick and Place by using Bluetooth Technology." Journal of Energy Engineering and Thermodynamics, no. 34, June 2023, pp. 16–21. https://doi.org/10.55529/jeet.34.16.21.

7. Srikanth, V. "Secret Sharing Algorithm Implementation on Single to Multi Cloud." Srikanth | International Journal of Research, 23 Feb. 2018, journals.pen2print.org/index.php/ijr/article/view/11641/11021.

8. V. Srikanth. "Managing Mass-Mailing System in Distributed Environment" v srikanth | International Journal & Magazine of Engineering, Technology, Management and Research, 23 August. 2015. http://www.ijmetmr.com/olaugust2015/VSrikanth-119.pdf

9. V. Srikanth. "SECURITY, CONTROL AND ACCESS ON IOT AND ITS THINGS" v srikanth | INTERNATIONAL JOURNAL OF MERGING TECHNOLOGY AND ADVANCED RESEARCH IN COMPUTING, 15 JUNE. 2017. http://ijmtarc.in/Papers/Current%20Papers/IJMTARC-170605.pdf

10. V. Srikanth. "ANALYZING THE TWEETS AND DETECT TRAFFIC FROM TWITTER ANALYSIS" v srikanth | INTERNATIONAL JOURNAL OF MERGING TECHNOLOGY AND ADVANCED RESEARCH IN COMPUTING, 20 MARCH. 2017. http://ijmtarc.in/Papers/Current%20Papers/IJMTARC-170309.pdf

11. V. Srikanth. "A NOVEL METHOD FOR BUG DETECTION TECHNIQUES USING INSTANCE SELECTION AND FEATURE SELECTION" v srikanth | INTERNATIONAL JOURNAL OF INNOVATIVE ENGINEERING AND MANAGEMENT RESEARCH, 08 DECEMBER. 2017. https://www.ijiemr.org/public/uploads/paper/976_approvedpaper.pdf

12. V. Srikanth. "SECURED RANKED KEYWORD SEARCH OVER ENCRYPTED DATA ON CLOUD" v srikanth | INTERNATIONAL JOURNAL OF INNOVATIVE ENGINEERING AND MANAGEMENT RESEARCH, 08 Febraury. 2018. http://www.ijiemr.org/downloads.php?vol=Volume-7&issue=ISSUE-02

13. V. Srikanth. "WIRELESS SECURITY PROTOCOLS (WEP,WPA,WPA2 & WPA3)" v srikanth | Journal of Emerging Technologies and Innovative Research (JETIR), 08 mAY. 2019. https://www.jetir.org/papers/JETIRDA06001.pdf

14. V. Srikanth, et al. "Detection of Fake Currency Using Machine Learning Models." Deleted Journal, no. 41, Dec. 2023, pp. 31–38. https://doi.org/10.55529/ijrise.41.31.38.

15. V. Srikanth, et al. "A REVIEW ON MODELING AND PREDICTING OF CYBER HACKING BREACHES." 25 Mar. 2023, pp. 300–305. http://ijte.uk/archive/2023/A-REVIEW-ON-MODELING-AND-PREDICTING-OF-CYBER-HACKING-BREACHES.pdf.

16. V. Srikanth, "DETECTION OF PLAGIARISM USING ARTIFICIAL NEURAL NETWORKS." 25 Mar. 2023, pp. 201–209. http://ijte.uk/archive/2023/DETECTION-OF-

PLAGIARISM-USING-ARTIFICIAL-NEURAL-NETWORKS.pdf.

17. V. Srikanth, "CHRONIC KIDNEY DISEASE PREDICTION USING MACHINELEARNINGALGORITHMS." 25 January. 2023, pp. 106–122. http://ijte.uk/archive/2023/CHRONIC-KIDNEY-DISEASE-PREDICTION-USING-MACHINE-LEARNING-ALGORITHMS.pdf.

18. Srikanth veldandi, et al. "View of Classification of SARS Cov-2 and Non-SARS Cov-2 Pneumonia Using CNN". journal.hmjournals.com/index.php/JPDMHD/article/view/3406/2798.

19. Srikanth veldandi, et al. "Improving Product Marketing by Predicting Early Reviewers on E-Commerce Websites." Deleted Journal, no. 43, Apr. 2024, pp. 17–25. https://doi.org/10.55529/ijrise.43.17.25.

20. Srikanth veldandi, et al."Intelligents Traffic Light Controller for Ambulance." Journal of Image Processing and Intelligent Remote Sensing, no. 34, July 2023, pp. 19–26. https://doi.org/10.55529/jipirs.34.19.26.

21. Veldandi Srikanth, et al. "Identification of Plant Leaf Disease Using CNN and Image Processing." Journal of Image Processing and Intelligent Remote Sensing, June 2024, https://doi.org/10.55529/jipirs.44.1.10.

22. Veldandi Srikanth, et al. "Human-AI Interaction Using 3D AI Assistant" International Conference on Emerging Advances and Applications in Green Energy (ICEAAGE-2024), 15, Feb 2024, https://www.researchgate.net/publication/380971799_ICEAAGE-2024_Conference_Proceedings_Final

23. Veldandi Srikanth, et al. "Voice Based Assistance for Traffic Sign Recognition System Using Convolutional Neural Network" International journal of advance and applied research (IJAAR) ISSN – 2347-7075, 4th, April 2024, https://ijaar.co.in/wp-content/uploads/2021/02/Volume-5-Issue-4.pdf

24. Veldandi Srikanth, et al. "Convolutional Neural Network Based Heart Stroke Detection" International journal of advance and applied research (IJAAR) ISSN – 2347-7075, 4th, April 2024, https://ijaar.co.in/wp-content/uploads/2021/02/Volume-5-Issue-4.pdf

25. Srikanth veldandi, et al. "Data Analytics Using R Programming Lab | IOK STORE." IOK STORE, ww.iokstore.inkofknowledge.com/product-page/data-analytics-using-r-programming-lab.

26. Srikanth veldandi, et al. "Data Structures Laboratory Manual | IOK STORE." IOK STORE, www.iokstore.inkofknowledge.com/product-page/data-structures-laboratory-manual.

27. Srikanth veldandi, et al. "Cyberspace and The Law: Cyber Security | IOK STORE." IOK STORE, iokstore.inkofknowledge.com/product-page/cyberspace-and-the-law.

**AUTHOR PROFILE:**

Mrs. Chepuri. Deepti, currently working as an Assistant Professor in the Department of Computer Science and Engineering, QIS College of Engineering and Technology, Ongole, Andhra Pradesh. She did her BTech from Uttar Pradesh Technical University, Lucknow, M.Tech from JNTUK, Kakinada. Her area of interest is Machine Learning, Artificial intelligence, Cloud Computing and Programming Languages.

Ms.Shaik.Sana, currently pursuing Master of Computer Applications at QIS College of Engineering and Technology (Autonomous), Ongole, Andhra Pradesh. She Completed BCA from Sri Nagarjuna Degree College, Ongole, Andhra Pradesh. Her areas ofinterest are Machine learning and Cloud Computing.