

# Prediction of Bigmart Sales using Machine Learning

K. Kumara Swamy, G. Sathwika, A. Priyanka, A. Jyothsna

1. Assistant Professor, Department of Information Technology,  
Malla Reddy Engineering College For Women(UGC-Autonomous),  
Hyderabad, India  
Email : kankala.kumar24@gmail.com

2. Department of Information Technology,  
Malla Reddy Engineering College For Women(UGC-Autonomous),  
Hyderabad, India  
Email : gundasathwika14@gmail.com

3. Department of Information Technology,  
Malla Reddy Engineering College For Women(UGC-Autonomous),  
Hyderabad, India  
Email : priyankareddyannadi@gmail.com

4. Department of Information Technology,  
Malla Reddy Engineering College For Women(UGC-Autonomous),  
Hyderabad, India  
Email: adullajyothsna23@gmail.com

## Abstract:

presently, supermarket run- centres, massive Marts keep track of every individual item's sales knowledge so as to anticipate potential client demand and update inventory management. Anomalies and general trends square measure usually discovered by mining {the knowledge the info the information} warehouse's data store. For retailers like massive mercantile establishment, the ensuing knowledge will be accustomed forecast future sales volume victimization numerous machine learning techniques like massive mercantile establishment.

## I. INTRODUCTION

Everyday competitiveness between various shopping centres as and as huge marts is becoming higher intense, violent just because of the quick development of global malls also online shopping. Each market seeks to offer personalized and limited-time deals to attract many clients relying on period of time, so that each item's volume of sales may be estimated for the organization's stock control, transportation and logistical services. The current machine learning algorithm is very advanced and provides methods for predicting or forecasting sales any kind of organization, extremely beneficial to overcome low – priced used for prediction. Always better prediction is helpful, both in developing and

improving marketing strategies for the marketplace, which is also particularly helpful

## II. RELATED WORK

A great deal of work having been gotten really intended to date the territory of deals foreseeing. A concise audit of the important work in the field of big mart deals is depicted in this part. Numerous other Measurable methodologies, for example, with regression, (ARIMA) Auto-Regressive Integrated Moving Average, (ARMA) Auto-Regressive Moving Average, have been utilized to develop a few deals forecast standards. Be that as it may, deals anticipating is a refined issue and is influenced by both outer and

inside factors, and there are two significant detriments to the measurable technique as set out in A. S. Weigend et A mixture occasional quantum relapse approach and (ARIMA) Auto-Regressive Integrated Moving Average way to deal with every day food deals anticipating were recommend by N. S. Arun Raj and furthermore found that the exhibition of the individual model was moderately lower than that of the crossover model. E. Hadavandi utilized theincorporation of “Genetic Fuzzy Systems (GFS)” and information gathering to conjecture the deals of the printed circuit board. In their paper, K-means bunching delivered K groups of all information records. At that point, all bunches were taken care of into autonomous with a data set tuning and rule-based extraction ability. Perceived work in the field of deals gauging was done by P.A. Castillo, Sales estimating of new distributed books was done in a publication market the executives setting utilizing computational techniques. “Artificial neural organizations” are additionally utilized nearby income estimating. Fluffy Neural Networks have been created with the objective of improving prescient effectiveness, and the Radial “Base Function Neural Network (RBFN)” is required to have an incredible potential for anticipating deals. Dataset: collected the dataset form the internet for the website called kaggle.com .In this work all having test dataset and train dataset in the test data set having a 5000 dataset and in the train data having a 8000 data

### **III. SYSTEM REQUIREMENTS**

#### **A. Hardware Requirements**

- System: MINIMUM i3
- Hard Disk:40 GB.
- Ram: 4 GB.

#### **B. Software Requirements**

- Operating System: Windows 8.
- Coding Language: Python 3.7

### **IV. SYSTEM STUDY**

#### **1 FEASIBILITY STUDY**

The feasibility of the project is analyzed in this phase and business proposal is put forth with a very general plan for the project and some cost estimates.

During system analysis the feasibility study of the proposed system is to be carried out. This is to ensure that the proposed system is not a burden to the company. For feasibility analysis, some understanding of the major requirements for the system is essential.

This study is carried out to check the economic impact that the system will have on the organization. The amount of fund that the company can pour into the research and development of the system is limited. The expenditures must be justified. Thus the developed system as well within the budget and this was achieved because most of the technologies used are freely available.

This study is carried out to check the technical feasibility, that is, the technical requirements of the system. Any system developed must not have a high demand on the available technical resources. This will lead to high demands on the available technical resources. This will lead to high demands being placed on the client. The developed system must have a modest requirement, as only minimal or null changes are required for implementing this system.

### **V. EXISTING SYSTEM**

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 andshort-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 thecheap availability of computing and storage, it has become possible to use sophisticated machine learning algorithms for this purpose.

### **VI. PROPOSED SYSTEM**

The data scientists at Big Mart have collected 2013 sales data for 1559 products across 10 stores in different cities. Also, certain attributes of each product and store have been defined. The aim is to build a predictive model and find out the sales of each product at a particular store. Using this model, Big Mart will try to understand the properties of products and stores which play a key role in increasing sales.

## VII. ALGORITHM AND MODULES

**CNN Algorithm** : A Convolutional Neural Network (CNN) is a Deep Learning algorithm which can take in an input image, assign importance to various aspects/objects in the image and be able to differentiate one from the other. The pre-processing required in a CNN is much lower as compared to other classification algorithms. While in primitive methods filters are hand-engineered, with enough training, CNN have the ability to learn these filters/characteristics.

**OpenCV**: OpenCV is a cross-platform library using which we can develop real-time computer vision applications. It mainly focuses on image processing, video capture and analysis including features like face detection and object detection.

**TensorFlow**: TensorFlow is a free and open-source software library for machine learning and artificial intelligence. It can be used across a range of tasks but has a particular focus on training and inference of deep neural networks.

## VII. MODULES

PredictiveModeling:

In order to find a decent model to predict sales we performed an extensive search of various machine learning models available in R, in particular of those accessible through the caret wrapper. In the end, however, models from the h2o package yielded the best results for the task. In particular, deep learning neural networks h2o.deep learning and gradient boosting regression trees h2o.gbm performed particularly well. An ensemble of various such models, constructed in h2oEnsemble.R forms the basis of our submission. Here, we used only the 12 most important predictors to avoid over-fitting. To include some features we may have missed with this rather small sub set of predictors we supplemented the ensemble with a deep learning neural net using 23 predictors. One of the most essential and commonly used regression techniques is linear regression. It's one of the most basic regression techniques. The simplicity with which the results may be interpreted is one of its primary merits.

$$= \theta_0 + \theta_1 x_1 + \dots + \theta_r x_r + \epsilon$$
 Where Y - Variable to be Predicted X – Variables used for making a prediction  $\theta_0, 1 \dots r$  - Regression Coefficients

- Random Error Regardless of how well the model

is trained, tested, and validated, there will always be a variation between observed and predicted, which is irreducible

error, so we cannot rely entirely on the learning algorithm's predicted results. Data must meet several conditions for a successful linear regression model. One of them is the lack of multiple linear regression, which means that the independent variables should be correlated. The RMSE value obtained from this algorithm is 1200.37 Ridge Regression: When multiple regression impacts results, Ridge Regression is employed. In multiple regression, the least square estimates are objective, but their variances are large and vary from the real value. Ridge regression eliminates standard errors by introducing a degree of bias to regression computations. In Ridge Regression, the Linear Regression Loss function is extended to punish not just the number of square residuals but also the parameter estimations. The RMSE value obtained from this algorithm is 1200.37 Lasso Regression: The Least Absolute Shrinkage Selector Operator regression makes some Coefficients to be zero, given a chance and improves the model. Thus, Lasso regression enables feature selection. Even at small alpha's, our coefficients are reducing to absolute zero. Therefore, Lasso selects only some features while reduces the coefficients of others to zero. This property of Lasso regression is called feature selection. The RMSE value obtained from this algorithm is 1093.94

## IX .RESULT AND DISCUSSIONS

In this work, the effectiveness of various algorithms on the data on revenue and review of, best performance-algorithm, here propose a software to using regression approach for predicting the sales centered on sales data from the past the accuracy of linear regression prediction can be enhanced with this method, polynomial regression, Ridge regression, and Xgboost regression can be determined. So, we can conclude ridge and Xgboost regression gives the better prediction with respect to Accuracy, MAE and RMSE than the Linear and polynomial regression approaches. In future, the forecasting sales and building a

sales plan can help to avoid unforeseen cash flow and manage production, staff and financing needs more effectively. In future work we can also consider with the ARIMA model which shows the time series graph.

#### **X. ACKNOWLEDGMENT**

The authors would like to thank all our anonymous critics for their feedback on this paper, and our project guide, along with the faculty of Mallareddy Engineering college for women (Autonomous) for their unconditional support.

#### **XI .REFERENCES**

[1] Ching Wu Chu and Guoqiang Peter Zhang, "A comparative study of linear and nonlinear models for aggregate retail sales forecasting", *Int. Journal Production Economics*, vol. 86, pp. 217- 231, 2003.

[2] Wang, Haoxiang. "Sustainable development and management in consumer electronics using soft computation." *Journal of Soft Computing Paradigm (JSCP)* 1, no. 01 (2019): 56.- 2. Suma, V., and Shavige Malleshwara Hills. "Data Mining based Prediction of D

[3] Suma, V., and Shavige Malleshwara Hills. "Data Mining based Prediction of Demand in Indian Market for Refurbished Electronics." *Journal of Soft Computing Paradigm (JSCP)* 2, no. 02 (2020): 101- 110

[4] Giuseppe Nunnari, Valeria Nunnari, "Forecasting Monthly Sales Retail Time Series: A Case Study", *Proc. of IEEE Conf. on Business Informatics (CBI)*, July 2017.

[5] <https://halobi.com/blog/sales-forecasting-five-uses/>. [Accessed: Oct. 3, 2018]

[6] Zone-Ching Lin, Wen-Jang Wu, "Multiple Linear Regression Analysis of the Overlay Accuracy Model Zone", *IEEE Trans.*

on Semiconductor Manufacturing, vol. 12, no. 2, pp. 229 – 237, May 1999.

[7] O. Ajao Isaac, A. Abdullahi Adedeji, I. Raji Ismail, "Polynomial Regression Model of Making Cost Prediction In Mixed Cost Analysis", *Int. Journal on Mathematical Theory and Modeling*, vol. 2, no. 2, pp. 14 – 23, 2012.

[8] C. Saunders, A. Gammernan and V. Vovk, "Ridge Regression Learning Algorithm in Dual Variables", *Proc. of Int. Conf. on Machine Learning*, pp. 515 – 521, July 1998. *IEEE TRANSACTIONS ON INFORMATION THEORY*, VOL. 56, NO. 7, JULY 2010 3561.

[9] "Robust Regression and Lasso". Huan Xu, Constantine Caramanis, Member, IEEE, and Shie Mannor, Senior Member, IEEE. 2015 International Conference on Industrial Informatics-Computing Technology, Intelligent Technology, Industrial Information Integration." "An improved Adaboost algorithm based on uncertain functions". Shu Xinqing School of Automation Wuhan University of Technology. Wuhan, China Wang Pan School of the Automation Wuhan University of Technology Wuhan, China.

[10] Xinqing Shu, Pan Wang, "An Improved Adaboost Algorithm based on Uncertain Functions", *Proc. of Int. Conf. on Industrial Informatics – Computing Technology, Intelligent Technology, Industrial Information Integration*, Dec. 2015.

[11] A. S. Weigend and N. A. Gershenfeld, "Time series prediction: Forecasting the future and understanding the past", Addison-Wesley, 1994.

[12] N. S. Arunraj, D. Ahrens, A hybrid seasonal autoregressive integrated moving average and quantile regression for daily food sales forecasting, *Int. J. Production Economics* 170 (2015) 321-335

