

## **SALES PREDICTION AND ANALYSIS OF SUPERMARKETS USING RIDGE AND POLYNOMIAL REGRESSION TECHNIQUES**

**Arun B Prasad<sup>1</sup>**

<sup>1</sup>Associate Professor (Economics), Institute of Law Nirma University, Ahmedabad, India,  
Email: [1arunprasad16@gmail.com](mailto:1arunprasad16@gmail.com), Orcid ID: 0000-0002-6108-9219

**Ram Bhawan Singh<sup>2</sup>**

<sup>4</sup>Professor and Head, Department of Computer Science, Tula's Institute of Technology,  
Dehradun, India, India, Email: [2rambhawansingh@gmail.com](mailto:2rambhawansingh@gmail.com)

**Manish Garg<sup>3</sup>**

<sup>1</sup>Assistant Professor, School of Management and Commerce, Dev Bhoomi Uttarakhand  
University, Dehradun, India, Email: [3mailgargmanish@gmail.com](mailto:3mailgargmanish@gmail.com), Orcid ID:0000-0002-0505-9680

### **Abstract**

Due to the obvious rapid growth of global shops including e - shopping, daily competition amongst numerous shopping complexes as well as large super markets is growing fiercer & more aggressive. Every marketplace tries to attract the interest of consumers by providing tailored as well as restricted discounts, such that the quantity of revenues with each product can be projected for the company's managing inventory, shipping, including logistical services. Nowadays, supermarkets run own branches and franchises, known as Big Marts, keep records of every product's revenue information in order to forecast possible customers' needs & adjust inventory control. Monitoring the information warehouse's storage space is a common way to find abnormalities as well as general patterns. This generated data may be utilised by merchants like Big Mart to anticipate the future sales volume using different machine learning approaches. For estimating the sales of the company such as Sales -Mart, a predicting model is constructed utilizing XGBoost, Linear regression, Polynomial regression, as well as Ridge regression approaches, and that it reveals that the system provides the best model.

Keywords: Big Marts; Anomalies; sustainable; machine learning

### **1 Introduction**

Due to the obvious rapid growth of global shops including e - shopping, daily competition amongst numerous shopping complexes as well as large super markets is growing fiercer & more aggressive. Every marketplace tries to attract the interest of consumers by providing tailored as well as restricted discounts, such that the quantity of revenues with each product can be projected for the company's managing inventory, shipping, including logistical services [2]. The present machine learning model is quite powerful, and so it gives ways for estimating or forecasting sales in any type of company, which is incredibly useful in overcoming low-cost

prediction approaches. Better forecasting is always beneficial, both in designing & enhancing marketing plans for the business, that is especially beneficial [3-5].

Organisations from a wide range of industries have realised that they require more computer scientists. Institutions of higher learning are hurrying to put up data analyst training programmes.

Data science is being promoted as a trendy career option in magazines [6.] Nevertheless, when the notion fades into useless chatter, there is some disagreement regarding what precisely disenchantment is. In this paper, we suggest that there have been legitimate reasons why defining data science has indeed been difficult [7]. One explanation seems to be that data science is inextricably linked to other significant ideas such as big data as well as data-driven goal setting, that are both gaining traction & awareness [8]. Another issue is the obvious propensity, in the lack of intellectual programmes to educate differently, to equate what a practitioners performs admirably with the description of the practitioner's profession; this could lead to a neglect of the subject's basics [9].

## 2 Related Works

A significant amount of effort has been put in to date the region of deal forecasting. This section provides a quick overview of the major works in the topic of big-mart transactions. A variety of other Quantifiable techniques, such as regression, Auto-Regressive Incorporated Moving Average, as well as Auto-Regressive Exponentially Weighted moving, have been used to establish a few deal forecasting benchmarks [10]. Transactions forecasting, on the other side, is a complex subject that is impacted by both external and internal causes, because there are two important drawbacks to the quantifiable approach, as outlined in a combination of frequent stochastic relapsing with Auto-Regressive Integral Moving Average. The median technique to cope with day-to-day foodstuff transactions was suggested [11].

The individualized model's presentation was determined to be somewhat lower than the crossovers author's. The use of "Genetic Fuzzy Systems" as well as knowledge collecting to predict the printed circuit panel's transactions. K-means clumping supplied K groupings of all data and information in their document [12]. At around that time, all of the bundles were put into automated mode, complete with given dataset tweaking including rule-based extracting capabilities. Work in the industry of transactions measure as it is viewed. The administrators of a publishing industry used computer approaches to estimate the sales of freshly dispersed publications. Regarding earnings estimation, "artificial neural organisations" are also used [13]. The Radial "Base " Functional Neural Network is necessary to have an opportunity to take advantage for predicting agreements, and Fluff Convolutional Networks were designed with the goal of boosting prophetic efficacy.

However, data science entails far more than data-mining techniques. Data analysts that are effective would have to be able to see business challenges via the lens of information. Data analytic reasoning has a basic framework including basic concepts that should be grasped [14]. Several "traditional" disciplines of research are used in computer science. The basic concepts of statistical inference must be grasped. A significant percentage of what has historically been examined in the subject of statistical is crucial to data analytics. Data visualisation techniques & methodology are critical. In addition, there are some situations

when vision, imagination, good judgement, & understanding of a specific technology must always be applied [15]. A datascience viewpoint gives structure & standards to professionals, giving data scientists a foundation to methodically address difficulties of obtaining useful insights from data.

### 3. Proposed method

Figure 1 depicts the suggested model's architectural structure, which focuses on the various technique applications to the information. They calculate Reliability, MAE, MSE, RMSE, as well as finally arrive with the optimal yielding method. Here are several examples: Algorithms are employed.

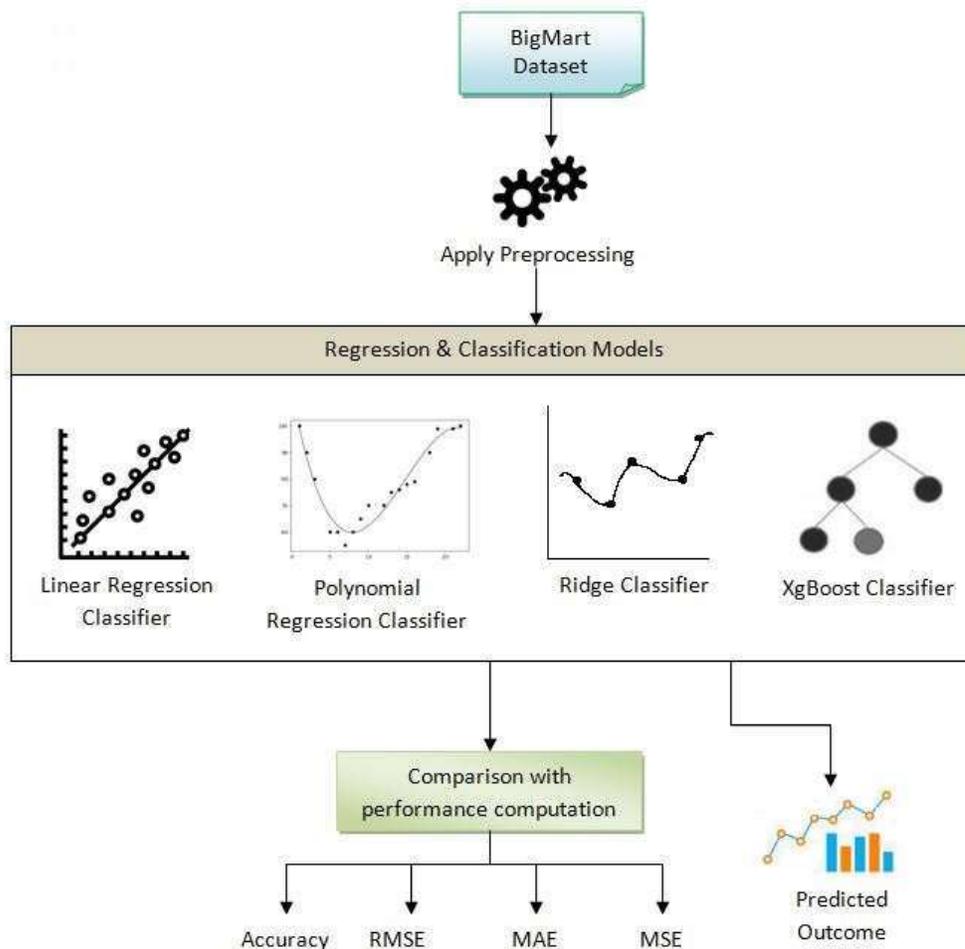


Figure 1: The proposed Architecture

#### 3.1 Approaches:

Machine learning is strongly tied to computational statistics as a topic, therefore possessing a background in statistically can help to clarify & apply machine learning methods. For individuals who haven't learned statistical, it's a good idea to start with a definition of

independent variables and dependent variable, which are two often used approaches for examining the connection between statistical method. Correlation coefficient is a statistical of the linear association that are neither reliant nor autonomous. At its most fundamental level, regression is used to investigate the connection among one reliant and one predictor variables. Recurrent statistics offer forecasting skills since they may be used to forecast the explanatory variables whenever the predictor variables is understood. Data mining techniques are indeed being created all the time.

For our objectives, we'll go through a handful of the most widely utilised techniques in machine learning at the time of posting.

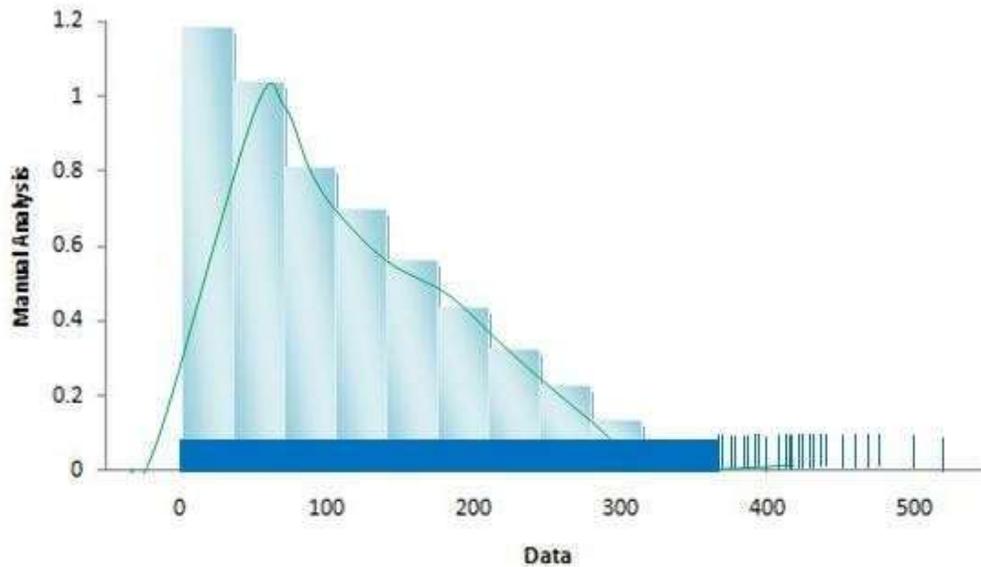
### **3.1.1 k-nearest neighbor:**

This k-nearest neighbour technique is an information processing technique that may be applied to both regression and classification problems. That k in k-nearest neighbour is a positive number, that is often low, and is sometimes shortened as k-NN. This inputs would be the k closest trained instances inside a region in either regression and classification. We'll concentrate on the k-NN classification method. The result of this function is class membership. It'll also place a new instance in the category that has the most members between its k closest neighbours.

The item is allocated to the category of the sole closest neighbor when  $k = 1$ .

### **3.1.2 Data set**

This database has been mined for important data elements at this stage. That is, it is attempting to discern knowledge from assumptions vs existing facts. This demonstrates that the parameters Endpoint width & Article mass have incomplete data, and also that the minimal value of Items Transparency is 0, that is not realistic. The year Outlets was founded ranges from 1985 to 2009. In this format, certain values might not be acceptable. As a result, we must translate them into the age of a certain output. The collection contains 1559 distinct goods and also 10 distinct retailers. There are 16 different values for the property Type of product. And yet there are 2 types of Product Saturated Fat, many of them are misspelt, such as normal rather than 'Normal' as well as reduced fat, LF rather than Minimal Fat. Figure 2 is an example. This dependent variable, Object Outlets Revenue, was found to have positive biased. So, on Product Outlets Revenues, a logarithmic activity was done to eliminate the distortion of the answer variables.



**Figure 2: xg boost regressor**

Both manual evaluation as well as automated information discovery software applications are used to graphically study and analyze correlation between two selected variables, the formation of the data source, the existence of exceptions, as well as the dissemination of datasets in order to expose shapes and items of interest, allowing dbas to gain a better understanding of the original data.

Database administrators must first comprehend and establish a complete representation of the study before collecting important data that can be analyzed, including such univariate, bivariate, multivariate, including principal components analyzation.

### 3.1.3 Linear Regression

Create a storyline that is disjointed. 1) a statistical trend that is regular or non-linear, & 2) a variability. If the labeling isn't straight, contemplate a modification. Foreigners could only recommend deleting them if there is indeed a non-statistical rationale if that's the situation. Connect the information to the linear regression lines and then use the remainder graphic to check the parameter estimates. If the expectations stated do not appears to be satisfied, a conversion may be required.

Convert the raw to the least squares if necessary, then draw a trendline using the converted information. Directly introduce procedure 1 if a modification has been performed. Keep stepping 5 if you haven't already. Construct the least-square slope of the line equation once a "good-fit" classical has been specified. Standard estimating, prediction, plus Rsquared mistakes were included.

The change in X that describes the entire variation in Y is defined. This could be simply calculated as:

^

$$= 1 - \frac{\sum (y_i - \hat{y}_i)^2}{\sum (y_i - \bar{y})^2}$$

### 3.1.4 Regression Algorithm

Polynomial Regression is a linear computation that analyses the connection between the variables of the study and use the most radical polynomials limiting. The following is the criteria for polynomials recurrence:  $y = b_0 + b_1x_1 + b_2x_1^2 + b_3x_1^3 + \dots + b_nx_1^n$ . It is sometimes referred to as an uncommon case of many straight relapses in ML. We changed it to polynomials relapsing correction to enhance the accuracy because they employ certain polynomials components to the various straightforward recurrence conditions. This material gathering used to prepare for polynomials recurrence is non-linear in character. It fits complicated & non-linear equations as well as information using a regression approach.

Ridge regression is a modeling tweaking approach that may be used to assess whatever dataset that has multi - collinearity. The L2 regularisation operation is performed using this way. Because least - square were impartial, andas well as the deviations are substantial if multi - collinearity difficulties develop, culminating in the predicted values being distant from the true numbers. "Advanced Gradient Boosting" is similar to gradient boosting but significantly more powerful. This has a branch method as well as a prediction method solution. This allows "xgboost" to run many times faster than the current gradient enhancing algorithms. It supports a number of goal abilities, such as recurrence, ordering, & evaluation. Because "xgboost" has a strong precognitive force but is often slow to organise, it is ideal for particular rivalry. It's also beneficial for bridge & identifying relevant components.

### 4. Results and discussions

Simple to complex machine learning techniques, including such Regression Analysis, Ridge Regression, Logistic Regression, Randomized Forests, as well as XGBoost, are used to anticipate BigMart's sales. With XGBoost techniques that have a reduced RMSE score, overall efficiency has been demonstrated. As a consequence, due to its speedy and relatively easy calculation, further Hyper - parameters Optimization was performed on XGBoost using the Bayesian Optimization approach, that resulted in the attainment of the least RMSE value as well as improved the model's fit to the fundamental findings. The outcome is a submissions document reporting Items Outlets Purchases for Item based on the Framework.

TABLE 1: Shows the linear regression result on the various parameter

Variable	Value
MSE	7.592
MAE	1.253
RMSE	2.724

TABLE 2: Shows the polynomial regression result on the various parameter

Variable	Value
MSE	6.214
MAE	2.965
RMSE	7.825

TABLE 3: Shows the Ridge regression result on the various parameter

Variable	Value
MSE	3.662
MAE	8.92
RMSE	1.915

TABLE 4: Shows the Xgboost regression result on the various parameter

Variable	Value
MSE	0.002
MAE	0.030
RMSE	0.033

TABLE 5: Shows the Xgboost regression frequency of item fat content

Variable	Value
Low fat	5124
Regular	2897
LF	305
Reg	115

TABLE 6: Comparison of MAE, MSE, RMSE with the Model

Model	MSE	MAE	RMSE
Linear Regression	7.532	1.168	2.733
Polynomial Regression	2.124	7.005	1.437
Ridge Regression	3.674	8.292	1.924
Xgboost Regression	0.001	0.030	0.0325

This randomized speculating method, as shown in the lift graph, delivers the targeted letter to all possible consumers. It's possible that half of the targeted customers will respond, that will provide as a benchmark for measuring the raise. This optimal model's maximum is about 48%, which means that with the precision of the mistake, you need only to send a letter to 48% of targeted buyers to receive 100% of the targeted consumer reaction. The real predictive

algorithm has a modest improvement around 60 and 75 percent when choosing the intended audience of 48 percent. Randomized guessing systems achieve poor response percentages over choice trees & clustering algorithms. The clustering algorithm has a number of respondents of 73.43 percent, whereas hierarchical clustering has a number of respondents of 61.66 percent. As a result, the decision tree algorithm improves the most, and the responsiveness is superior to the clustering modelling process.

## 5. Conclusion

Considering efficiency of many techniques on taking into account the amount as well as a study of, excellent performance, herein present a programme to use linear regression for projecting sales based on sales data from the past in this work. This strategy can improve the efficiency of regression analysis predictions, as well as polynomials analysis, Ridge regression, & Xgboost regression. As a result, they can infer that ridges & Xgboost analysis provide superior predictions in terms of effectiveness, MAE, and RMSE than linear and polynomials extrapolation. Important for decision making and developing a sales strategy in the future might aid in preventing unexpected working capital as well as better manage manufacturing, staffing, and finance requirements. Researchers could also explore using the ARIMA model, that displays a time - series data chart, in future employment. The results of machine learning techniques would assist in the selection of the most appropriate sales forecasting method, that BigMart will use to plan its product offerings.

## References

- [1] Yelne, A., & Theng, D. (2021, November). Stock Prediction and analysis Using Supervised Machine Learning Algorithms. In *2021 International Conference on Computational Intelligence and Computing Applications (ICCICA)* (pp. 1-6). IEEE.
- [2] Shen, S., Jiang, H., & Zhang, T. (2012). Stock market forecasting using machine learning algorithms. *Department of Electrical Engineering, Stanford University, Stanford, CA*, 1-5.
- [3] Shen, S., Jiang, H., & Zhang, T. (2012). Stock market forecasting using machine learning algorithms. *Department of Electrical Engineering, Stanford University, Stanford, CA*, 1-5.
- [4] Shen, S., Jiang, H., & Zhang, T. (2012). Stock market forecasting using machine learning algorithms. *Department of Electrical Engineering, Stanford University, Stanford, CA*, 1-5.
- [5] Balamurugan, K., Pavan, M. V., Ali, S. A., & Kalusuraman, G. (2021). Compression and flexural study on PLA-Cu composite filament using FDM. *Materials Today: Proceedings*, 44, 1687-1691.
- [6] Albuquerque, V., Sales Dias, M., & Bacao, F. (2021). Machine learning approaches to bikesharing systems: A systematic literature review. *ISPRS International Journal of GeoInformation*, 10(2), 62.

- [7] Behera, G., & Nain, N. (2019, September). A comparative study of big mart sales prediction. In *International Conference on Computer Vision and Image Processing* (pp. 421-432). Springer, Singapore.
- [8] Yelne, A., & Theng, D. (2021, November). Stock Prediction and analysis Using Supervised Machine Learning Algorithms. In *2021 International Conference on Computational Intelligence and Computing Applications (ICCICA)* (pp. 1-6). IEEE.
- [9] Garikapati, P., Balamurugan, K., Latchoumi, T. P., & Malkapuram, R. (2021). A Cluster Profile Comparative Study on Machining AlSi7/63% of SiC Hybrid Composite Using Agglomerative Hierarchical Clustering and K-Means. *Silicon*, 13(4), 961-972.
- [10] Latchoumi, T. P., & Parthiban, L. (2021). Quasi oppositional dragonfly algorithm for load balancing in cloud computing environment. *Wireless Personal Communications*, 1-18.
- [11] Karthikeyan, K. A., Balamurugan, K., & Rao, P. M. V. Studies on cryogenically treated WC-Co insert at different soaking conditions. *Materials and Manufacturing Processes*, 35(5), 345-355.
- [12] Albuquerque, V., Sales Dias, M., & Bacao, F. (2021). Machine learning approaches to bikesharing systems: A systematic literature review. *ISPRS International Journal of GeoInformation*, 10(2), 62.
- [13] Punam, K., Pamula, R., & Jain, P. K. (2018, September). A two-level statistical model for big mart sales prediction. In *2018 International Conference on Computing, Power and Communication Technologies (GUCON)* (pp. 617-620). IEEE.
- [14] Kaunchi, P., Jadhav, T., Dandawate, Y., & Marathe, P. (2021, October). Future Sales Prediction For Indian Products Using Convolutional Neural Network-Long Short Term Memory. In *2021 2nd Global Conference for Advancement in Technology (GCAT)* (pp. 1-5). IEEE.
- [15] Sharma, S., Deepika, D., & Singh, G. (2021, December). Intelligent Warehouse Stocking Using Machine Learning. In *2021 IEEE International Conference on Mobile Networks and Wireless Communications (ICMNBC)* (pp. 1-6). IEEE.