



# **iJRASET**

International Journal For Research in  
Applied Science and Engineering Technology



---

# **INTERNATIONAL JOURNAL FOR RESEARCH**

IN APPLIED SCIENCE & ENGINEERING TECHNOLOGY

---

**Volume:** 10    **Issue:** VIII    **Month of publication:** August 2022

**DOI:** <https://doi.org/10.22214/ijraset.2022.46068>

**[www.ijraset.com](http://www.ijraset.com)**

**Call:** ☎ 08813907089

**E-mail ID:** [ijraset@gmail.com](mailto:ijraset@gmail.com)

# Predictive Analysis for Big Mart Sales Using Machine Learning Algorithm

Pranav Kumar Sinha<sup>1</sup>, Neethi M V<sup>2</sup>

<sup>1</sup>PG Student, <sup>2</sup>Assistant Professor, The National Institute of Engineering, Mysuru, India

**Abstract:** Currently, Big Marts, the equivalent of supermarket run-canthers, keep track of each item's sales data in order to forecast implicit consumer demand and update force operation. In order to estimate the volume of bargains for each item for the association's stock control, transportation, and logistical services, each request aims to offer verified and limited time deals to attract numerous guests over time. By intentionally entangling the data store of the data storage, anomalies and broad trends are continuously uncovered. Retailers like Large Mart can use the performing data to predict future transaction volume utilising a variety of machine learning techniques, such as big bazaar. The present machine learning algorithm is very sophisticated and offers methods for predicting or reading deals with any kind of association, which is very beneficial to Always better prophecy is useful in creating and refining commercial marketing plans, which is particularly useful. The development of a prediction model utilising linear retrogression and Ridge retrogression methods for analysing the transactions of a company like Big- Mart, and it was found to perform better than models themselves. additional Measurable factors methods with regression, machine-accumulative (ARIMA), and Integrated Using Moving Average, (ARMA) machine-cumulative Moving normal, create many transactions that read morality.

**Keywords:** Linear Regression, Ridge Regression, Mean Absolute Error, Root Mean Square Error, Mean Square Error

## I. INTRODUCTION

Everyday competitiveness between colourful shopping centres and massive marts is getting advanced violent, and violent just because of the quick development of global promenades also online shopping. The growth of international malls and online shopping has led to an increase in the severity and acrimony of the competition between numerous shopping malls and massive supermarkets. Each request seeks to offer substantiated and limited time deals to attract numerous guests counting on a period of time, so that each item's volume of deals may be estimated for the association's stock control, transportation, and logistical services, in order to efficiently draw a big number of customers and determine the number of sales for each product, as well as for the business' logistics, distribution, and stock management requirements. The current machine learning is highly sophisticated and offers opportunities for forecasting or forecast demand for any type of organization in order to defeat low-cost prediction methods. For creating and enhancing market-specific marketing strategies, projections that are regularly updated are crucial. Always better vaticination is helpful, both in developing and perfecting marketing strategies for the business, which is also particularly helpful. But not all machine-learning techniques are equal, and not all of them are equally accurate. As a result, a machine-learning algorithm may be extraordinarily effective when applied to a particular problem but ineffective when applied to another. Due to this, Big Mart requires combining several machine-learning algorithms to produce a useful predictive model. projecting revenue with analytics. In order to find the most powerful predictive analytics We created a working prototype of a machine learning-based sales forecasting system for Big Mart. We must test the algorithm on Big Mart before launching this prototype. Genuine data from Mart. Consequently, we used Big Mart's sales data to test our prototype, and we used two variations to construct a machine-learning classifier model.

Proposed system is having Linear Regression is one of the easiest and most popular Machine Learning algorithms. It's a statistical system that's used for prophetic analysis. Linear retrogression makes prognostications for nonstop/ real or numeric variables similar as deals, payment, age, product price, etc. It Create a dispersed plot, There is a direct or complicated pattern (outliers) as well as friction in the data. If the marking is irregular, think of a metamorphosis. If there is a non-statistical base, it should only be advised to count non-natives in those circumstances. Using the residual plot (for the constant standard), connect the data to the least-squares line. the unity of friction, and they also support the model hypotheses (for the divagation thesis).

It may be essential to undergo a metamorphosis if the hypotheticals seem to be incorrect

Using the streamlined data and, if necessary, least places, create a retrogression line. So, it gives the linear values to predict.

The proposed system also allows Ridge regression in this while assessing the data that exhibits multicollinearity, crest retrogression is a model-tuning fashion employed. L2 regularisation is carried in this work. When least places are unprejudiced, multicollinearity problems do, and the dissonances are substantial, which causes a large gap between the anticipated and factual result.

## II. LITERATURE SURVEY

- 1) In this study, we examine to evaluate the forecasting performance of several linear and nonlinear models of total retail sales. Numerous conventional seasonal forecasting techniques, including the time series approach and the regression approach using seasonal dummy variables and trigonometric functions, are used because to the significant seasonal swings in retail sales. Neural networks, which are generalised nonlinear functional approximators, are used to implement the nonlinear versions of these methods. Deseasonalization and other seasonal time series modelling issues are also researched. We find that the nonlinear models outperform their linear counterparts in out-of-sample forecasting using repeated cross-validation samples, and that prior seasonal adjustment of the data can greatly enhance forecasting performance of the neural network model.
- 2) In this paper, we examine with the rising demand for such products over the past 10 years, we can observe that research on refurbished products has attracted more and more attention. We use a data-mining approach to conduct a thorough examination of the Indian e-commerce business in order to forecast the demand for reconditioned gadgets. Analysis is also done on how the variables and demand are affected by real-world conditions. Three arbitrary e-commerce websites' real-world datasets are taken into consideration for investigation. The collection, processing, and validation of data is done using effective algorithms. Based on the findings of this analysis, it is obvious that using the suggested approach, very accurate forecast can be achieved despite the effects of variable customer behaviour and market circumstances.
- 3) In this paper, we examine how A two-level strategy is used to estimate product sales from a certain outlet, and it outperforms any popular single model predictive learning algorithm in terms of predictive performance. The technique is applied on 2013 Big Mart Sales data. In order to anticipate outcomes accurately, data exploration, data transformation, and feature engineering are essential. The outcome showed that a two-level statistical method outperformed a single model approach because the former offered additional data that improved prediction.
- 4) In this paper we study about Support Vector Regression (SVM). Retrogression model construction grounded on sample data sets has been the main emphasis of previous ways in prognosticating review/ magazine deals. still, over-fitting can be a concern with these retrogression models. Support vector retrogression (SVR) was suggested as a unique approach to working the over-fitting issue in recent theoretical studies in statistics. SVR's thing is to attain the smallest structural threat rather than the smallest empirical threat, in discrepancy to classic retrogression models, which aim to minimize both. Support vector retrogression was therefore used in this work to break the soothsaying deals issue for journals and magazines. The results of the trial demonstrated that SVR is a better approach for this problem.

## III. DATA SETS

A group of data points that can be used by a computer for analysis and prediction as a single entity. collected data from the internet for the Kaggle.com website. The test data set in this study has 8542 rows and 12 classes, and it has been trained to produce the best prediction results.

Variable	Description	Relation to Hypothesis
Item_Identifier	Unique product ID	ID Variable
Item_Weight	Weight of product	Not considered in hypothesis
Item_Fat_Content	Whether the product is low fat or not	Linked to 'Utility' hypothesis. Low fat items are generally used more than others
Item_Visibility	The % of total display area of all products in a store allocated to the particular product	Linked to 'Display Area' hypothesis. More inferences about 'Utility' can be derived from this.
Item_Type	The category to which the product belongs	Not considered in hypothesis
Item_MRP	Maximum Retail Price (list price) of the product	ID Variable
Outlet_Identifier	Unique store ID	Not considered in hypothesis
Outlet_Establishment_Year	The year in which store was established	Not considered in hypothesis
Outlet_Size	The size of the store in terms of ground area covered	Linked to 'Store Capacity' hypothesis
Outlet_Location_Type	The type of city in which the store is located	Linked to 'City Type' hypothesis.
Outlet_Type	Whether the outlet is just a grocery store or some sort of supermarket	Linked to 'Store Capacity' hypothesis again.
Item_Outlet_Sales	Sales of the product in the particular store. This is the outcome variable to be predicted.	Outcome variable

Fig. 1 Attributes Information of Dataset



#### IV. PROPOSED WORK

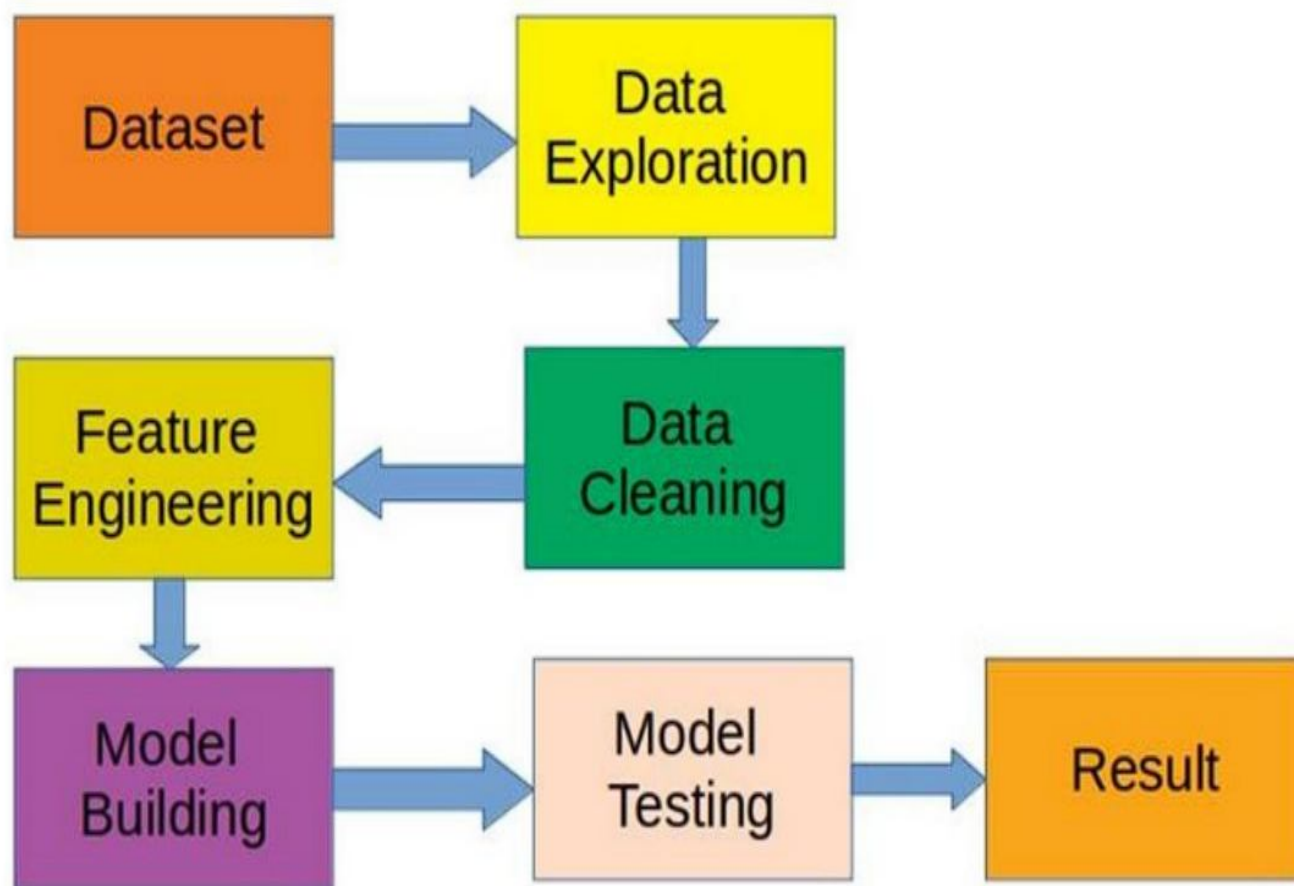
The proposed system gives most effective predictive analytics solution for sales forecasting realized the intended model's armature illustration, which focuses on the colourful algorithm operations to the dataset. We calculate the delicacy, MAE, MSE, and RMSE in this stage before choosing the stylish yield algorithm.

Furthermore, the system extends its functionalities by predicting the sales of outlet based on the trained datasets. Where the retailer uploads his sales chart and after that based on the best-chosen algorithms which gives optimal result with good accuracy the result is given. All the accuracy is shown in the form of graph and pie chart to better visualization. The system provides flexibility to the retailer and more effective and more adapted to handle massive data sets due to the inclusion of Ridge Regression and Linear Regression models. It also helps retailer to get how to improve his sales and fulfil the demands of customers.

#### V. METHODOLOGY

The proposed system utilizing the constructed system is referred to as "programme implementation". All procedures necessary to use the new programme are included in this. Confirming that the technology's processes are operating as anticipated is the organization's main objective after the planning phase. Prior to beginning the implementation process, a number of requirements must be satisfied. This system having any number of users can be supported by the system. An illustration of a non-functional need is this. The customer can watch the programme whenever it is convenient. The programme can be re-used, allowing the source code to be utilised to add additional capabilities with little to no changes.

Performance metrics will be provided by the programme we are creating.



Big Mart's data scientists gathered data from 10 businesses that were distributed across colourful locales, and each offered 1559 unique products. Using all the data, it's established what part particular item factors play and how they affect deals. The data collection comprises a variety of data types, similar as integer, pier, and object.

### A. Proposed Architecture Diagram

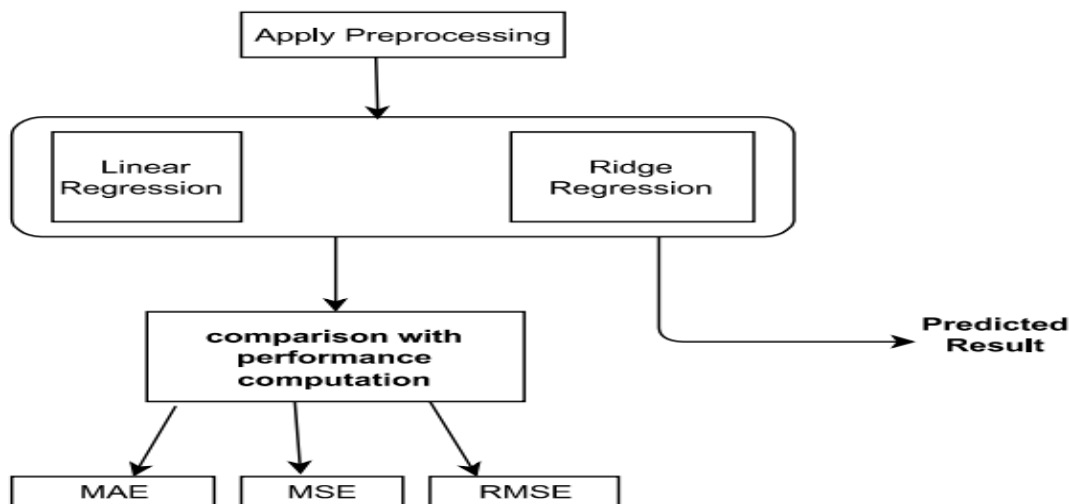


Fig. 2 Architecture Diagram

After pre-processing (cleaning and arranging) the data, the row data is prepared for constructing and ML model testing. The models concentrated on applying the two aforementioned algorithms to the datasets. The optimal yield algorithm is determined after computing the MAE, MSE, and RMSE.

Mean squared error	$MSE = \frac{1}{n} \sum_{t=1}^n e_t^2$
Root mean squared error	$RMSE = \sqrt{\frac{1}{n} \sum_{t=1}^n e_t^2}$
Mean absolute error	$MAE = \frac{1}{n} \sum_{t=1}^n  e_t $

Fig. 3 Mathematical formula performance computation

### B. Service Provider

Following the initial settings, the supplier tests and trains the datasets, compares accuracy using the MAE, MSE, and RMSE concepts, and prepares the machine to estimate the sales of large supermarkets.

### C. Remote User

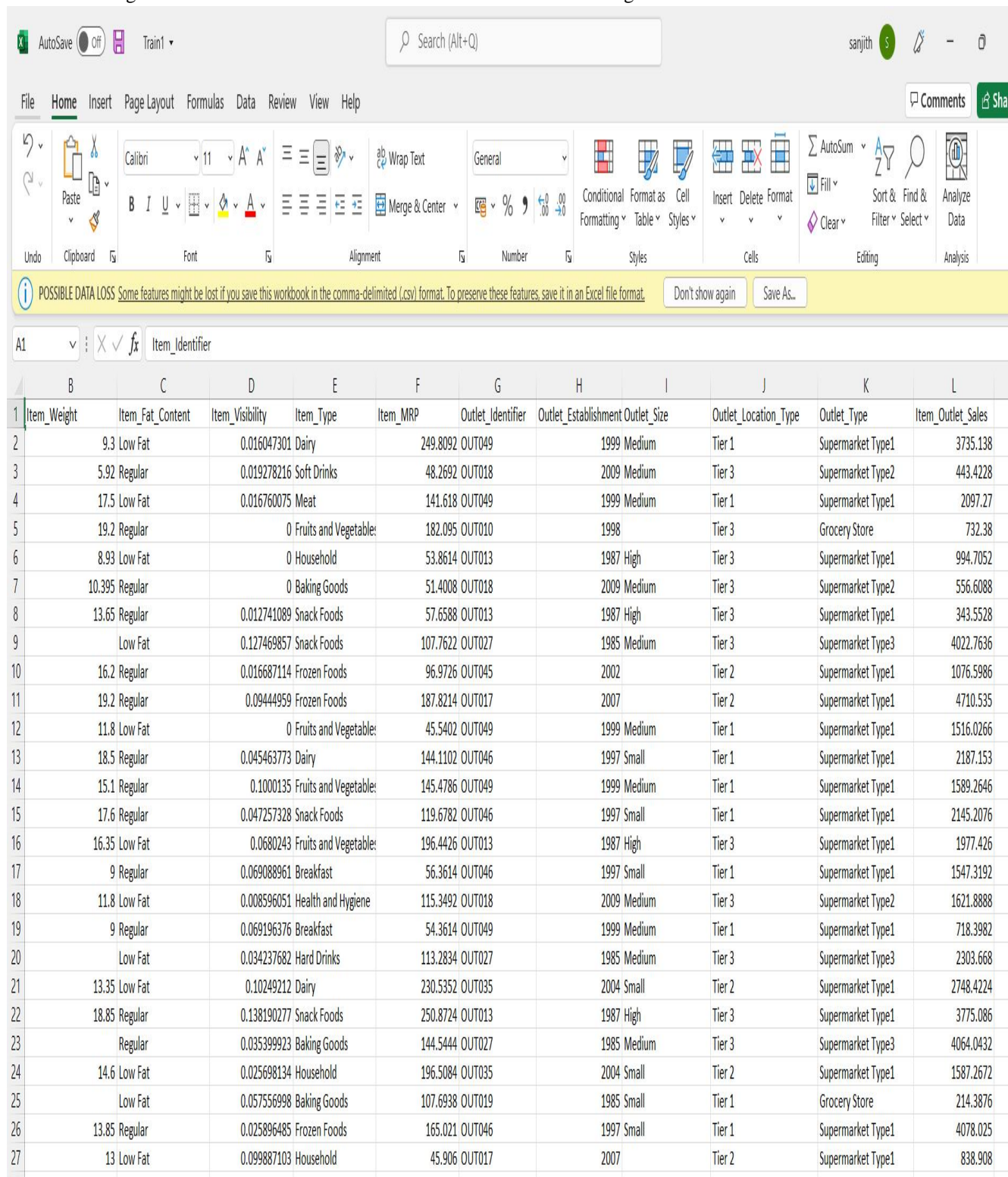
To get the most precise prediction result, the user must first register before they can connect into the site and input their sales forecast in xlxs format.

### D. View and Authorize Users

After the user uploads, the service provider will download the sales forecast after a short period of time, and after that, the business analysis team will meet in-depth with the store to discuss the profitability of sales and production.

## VI. RESULT ANALYSIS.

A subset of our real datasets called the "train dataset" is used by machine learning models to find and learn patterns. When a new input is provided based on data from a trained dataset, the trained dataset verifies the input and produces the most accurate and ideal results. The training datasets with all 12 columns and 8542 rows are shown in Fig. 3 below and are used to run the model.



	B	C	D	E	F	G	H	I	J	K	L	M
	Item_Weight	Item_Fat_Content	Item_Visibility	Item_Type	Item_MRP	Outlet_Identifier	Outlet_Establishment	Outlet_Size	Outlet_Location_Type	Outlet_Type	Item_Outlet_Sales	
2	9.3	Low Fat	0.016047301	Dairy	249.8092	OUT049	1999	Medium	Tier 1	Supermarket Type1	3735.138	
3	5.92	Regular	0.019278216	Soft Drinks	48.2692	OUT018	2009	Medium	Tier 3	Supermarket Type2	443.4228	
4	17.5	Low Fat	0.016760075	Meat	141.618	OUT049	1999	Medium	Tier 1	Supermarket Type1	2097.27	
5	19.2	Regular		0 Fruits and Vegetables	182.095	OUT010	1998		Tier 3	Grocery Store	732.38	
6	8.93	Low Fat		0 Household	53.8614	OUT013	1987	High	Tier 3	Supermarket Type1	994.7052	
7	10.395	Regular		0 Baking Goods	51.4008	OUT018	2009	Medium	Tier 3	Supermarket Type2	556.6088	
8	13.65	Regular	0.012741089	Snack Foods	57.6588	OUT013	1987	High	Tier 3	Supermarket Type1	343.5528	
9		Low Fat	0.127469857	Snack Foods	107.7622	OUT027	1985	Medium	Tier 3	Supermarket Type3	4022.7636	
10	16.2	Regular	0.016687114	Frozen Foods	96.9726	OUT045	2002		Tier 2	Supermarket Type1	1076.5986	
11	19.2	Regular	0.09444959	Frozen Foods	187.8214	OUT017	2007		Tier 2	Supermarket Type1	4710.535	
12	11.8	Low Fat		0 Fruits and Vegetables	45.5402	OUT049	1999	Medium	Tier 1	Supermarket Type1	1516.0266	
13	18.5	Regular	0.045463773	Dairy	144.1102	OUT046	1997	Small	Tier 1	Supermarket Type1	2187.153	
14	15.1	Regular	0.1000135	Fruits and Vegetables	145.4786	OUT049	1999	Medium	Tier 1	Supermarket Type1	1589.2646	
15	17.6	Regular	0.047257328	Snack Foods	119.6782	OUT046	1997	Small	Tier 1	Supermarket Type1	2145.2076	
16	16.35	Low Fat	0.0680243	Fruits and Vegetables	196.4426	OUT013	1987	High	Tier 3	Supermarket Type1	1977.426	
17	9	Regular	0.069088961	Breakfast	56.3614	OUT046	1997	Small	Tier 1	Supermarket Type1	1547.3192	
18	11.8	Low Fat	0.008596051	Health and Hygiene	115.3492	OUT018	2009	Medium	Tier 3	Supermarket Type2	1621.8888	
19	9	Regular	0.069196376	Breakfast	54.3614	OUT049	1999	Medium	Tier 1	Supermarket Type1	718.3982	
20		Low Fat	0.034237682	Hard Drinks	113.2834	OUT027	1985	Medium	Tier 3	Supermarket Type3	2303.668	
21	13.35	Low Fat	0.10249212	Dairy	230.5352	OUT035	2004	Small	Tier 2	Supermarket Type1	2748.4224	
22	18.85	Regular	0.138190277	Snack Foods	250.8724	OUT013	1987	High	Tier 3	Supermarket Type1	3775.086	
23		Regular	0.035399923	Baking Goods	144.5444	OUT027	1985	Medium	Tier 3	Supermarket Type3	4064.0432	
24	14.6	Low Fat	0.025698134	Household	196.5084	OUT035	2004	Small	Tier 2	Supermarket Type1	1587.2672	
25		Low Fat	0.057556998	Baking Goods	107.6938	OUT019	1985	Small	Tier 1	Grocery Store	214.3876	
26	13.85	Regular	0.025896485	Frozen Foods	165.021	OUT046	1997	Small	Tier 1	Supermarket Type1	4078.025	
27	13	Low Fat	0.099887103	Household	45.906	OUT017	2007		Tier 2	Supermarket Type1	838.908	

Fig. 4 Dataset with columns



After the initial setup has been completed the service provider can start the train and test dataset by that all 3 accuracy comparison computation as shown in the below figure Fig 4.

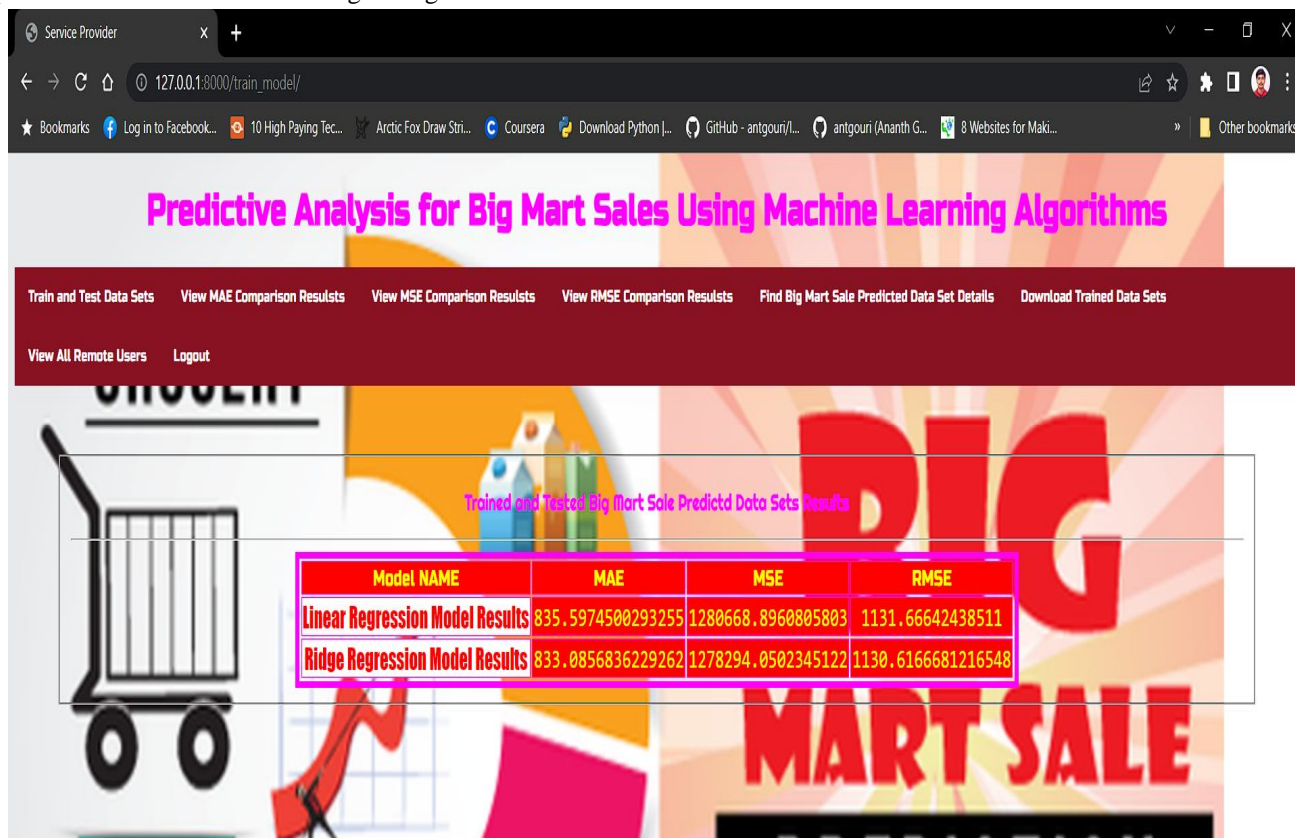


Fig. 5 Accuracy measurement

Without considering their direction, MAE calculates the average magnitude of the mistakes in a group of projections. The below figure Fig 5 shows the Mean Absolute Error bar graph result.

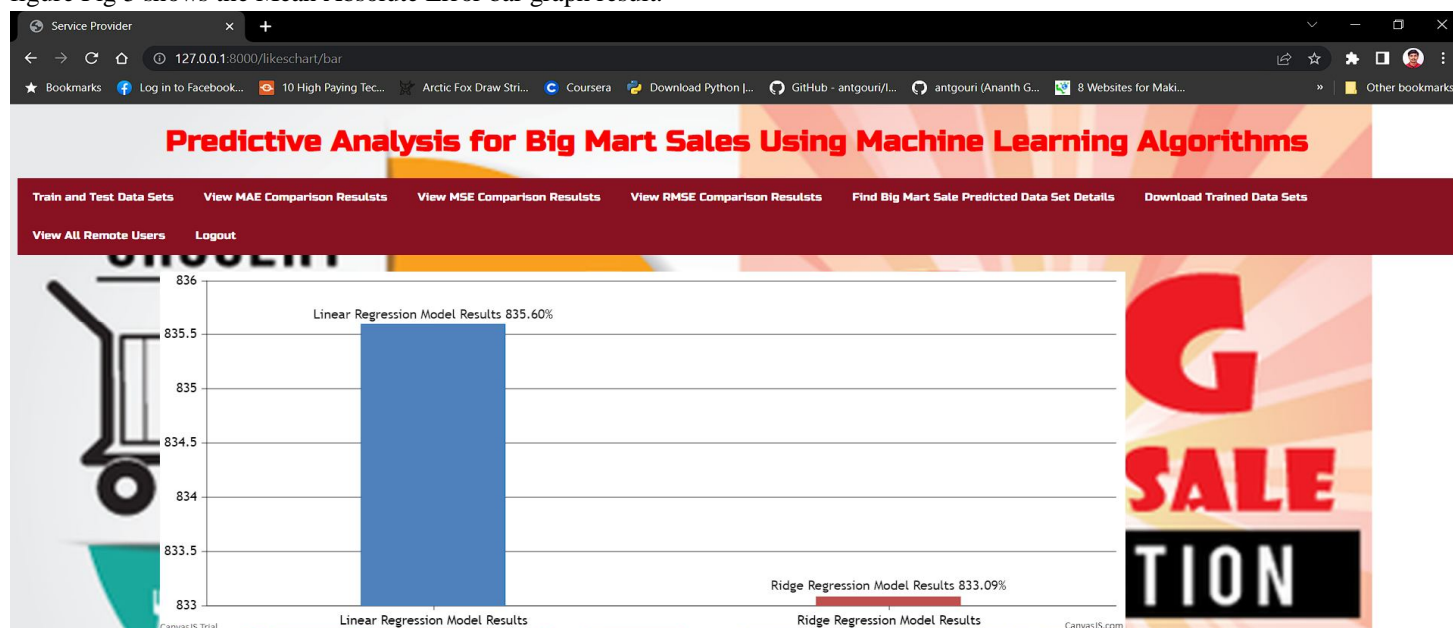


Fig. 6 Bar graph result of MAE

Perhaps the most basic and widely used loss function is the Mean Squared Error (MSE), which is frequently covered in beginner machine learning classes. The MSE is calculated by taking the difference between the predictions made by your model and the actual data, squaring it, and averaging it over the entire dataset. The below figure Fig 6 and 7 shows the pie chart and line graph measurement of it.

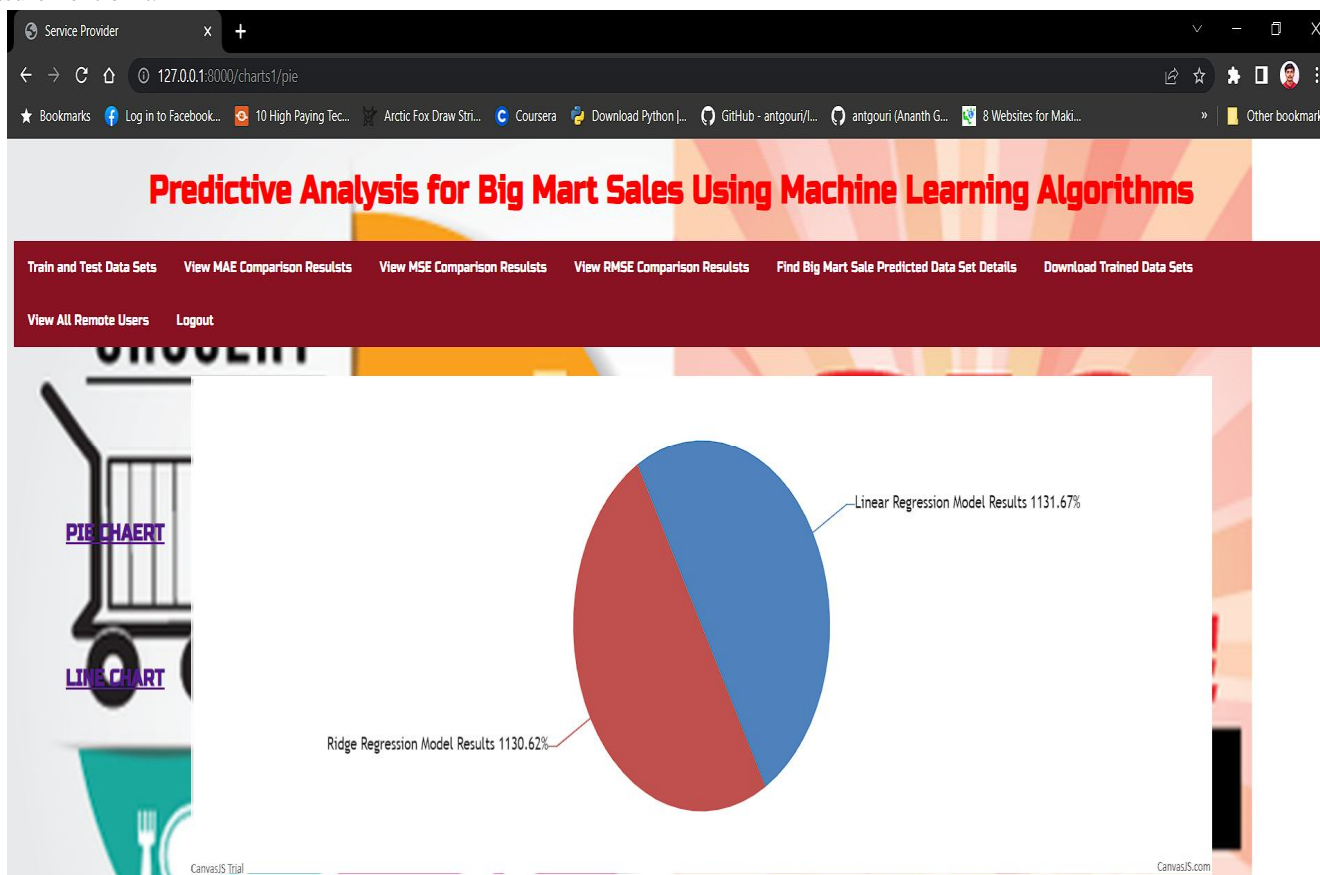


Fig. 7 Pie Chart of MSE

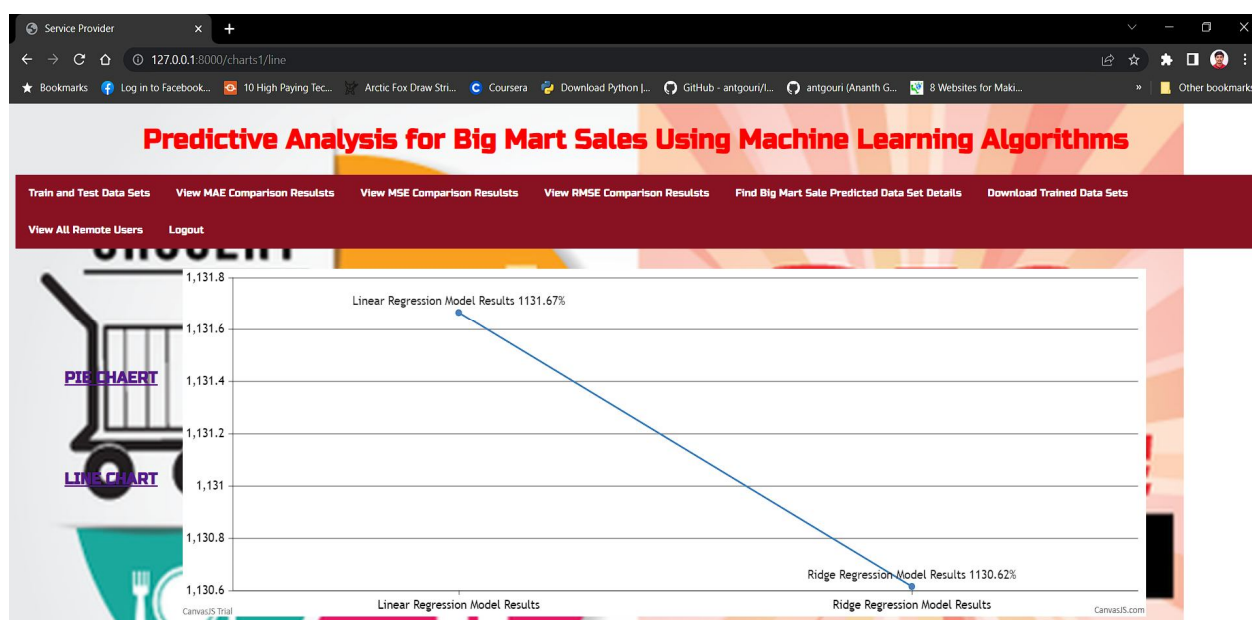


Fig. 8 Line chart of MSE



To reduce the root mean square error (RMSE), calculate the residual (difference between prediction and truth) for each data point, the norm of the residual, the mean of the residuals, and the square root of that mean. Since it requires and uses real measurements at each projected data point, RMSE is frequently utilised in supervised learning applications. The below figure Fig 8 and 9 shows the Root Mean Square error pie chart and line graph.

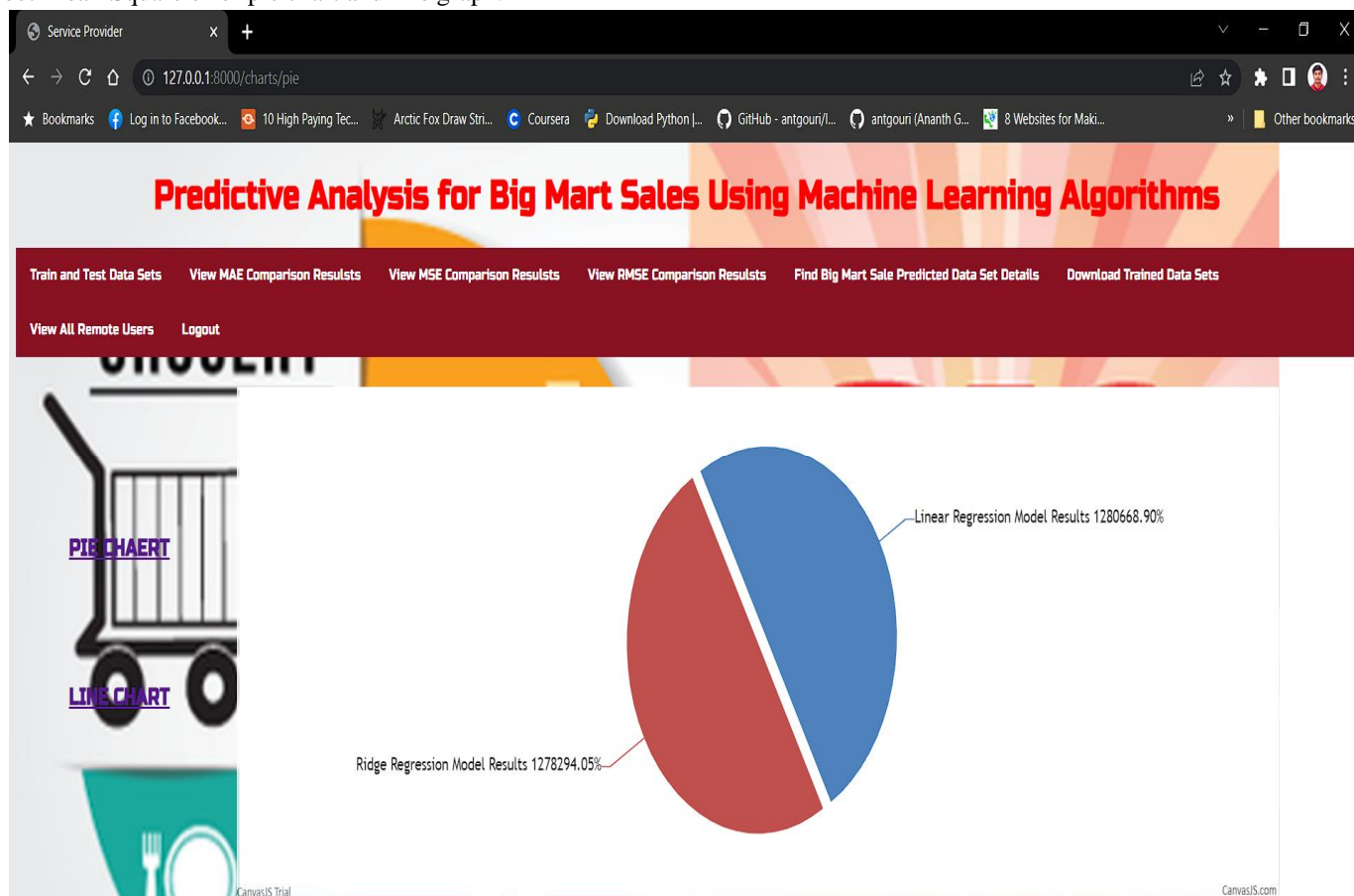


Fig. 9 Pie chart of RMSE

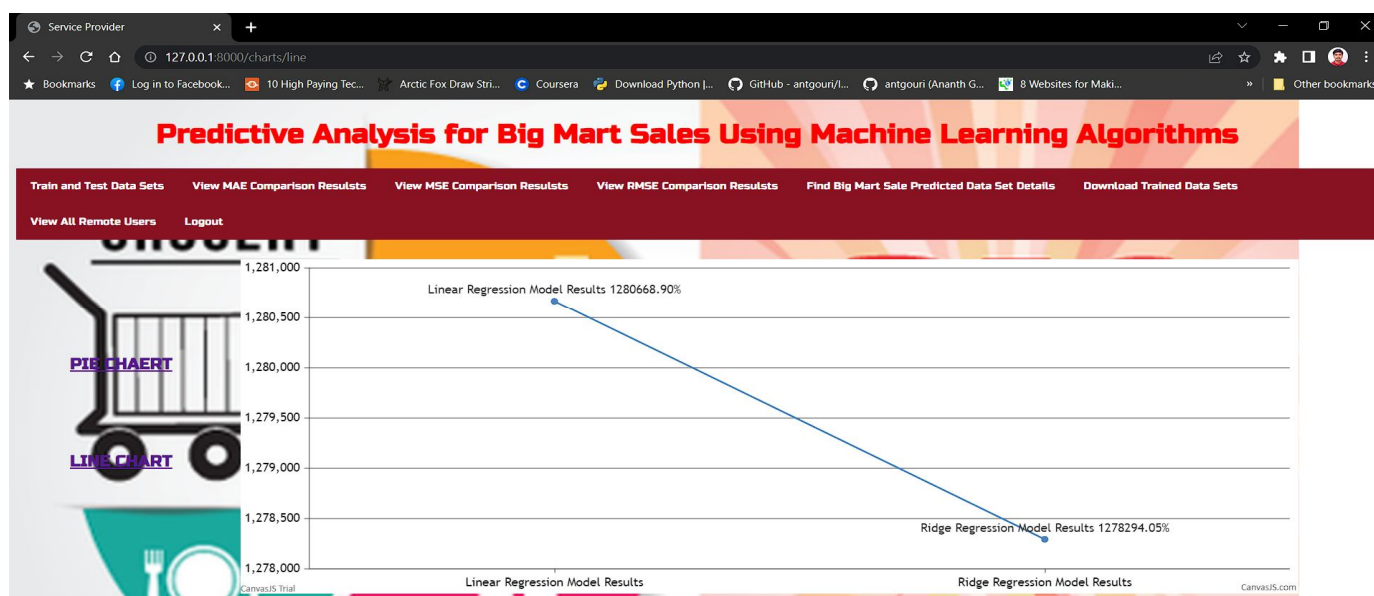


Fig. 10 Line chart of RMSE

## VII. CONCLUSION

The most efficient algorithm is one that, after examining the performance of colourful algorithms on profit data, employs a retrogression technique to forecast deals focusing on actual deal data. When using direct retrogression, prognostications may be more precise because using this technique. Ridge and linear retrogressions can also be found. Thus, we can conclude that the Ridge, MAE, RMSE, and MSE retrogression styles are the most effective. Regarding vaticination perfection, there are two retrogression styles: direct and linear. unborn child, Staffing, financial requirements, and transaction soothsaying will all make it easier to manage. making a business plan. The time series graph, which shows data through time, may also be used for future investigations the ARIMA simulation.

## REFERENCES

- [1] Ching Wu Chu and Guoqiang Peter Zhang, "A comparative study of linear and nonlinear models for aggregate retails sales forecasting", Int. Journal Production Economics, vol. 86, pp. 217- 231, 2003.
- [2] 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
- [3] Kumari Punam; Rajendra Pamula; Praphula Kumar Jain." A Two-Level Statistical Model for Big Mart Sales Prediction" IEEE 2018 International Conference on Computing, Power and Communication Technologies (GUCON).DOI: 10.1109/GUCON.2018.8675060.
- [4] Xiaodan Yua,b, Zhiquan Qib ,Yuanmeng Zhaoc." support Vector Regression for Newspaper/Magazine Sales Forecasting" Published by Elsevier B.V. 2013 International Conference on Information Technology and Quantitative Management Open access under CC BY-NC-N. DOI: 10.1016/j.procs.2013.05.134



10.22214/IJRASET



45.98



IMPACT FACTOR:  
7.129



IMPACT FACTOR:  
7.429



# INTERNATIONAL JOURNAL FOR RESEARCH

IN APPLIED SCIENCE & ENGINEERING TECHNOLOGY

Call : 08813907089  (24\*7 Support on Whatsapp)