



A Comparative Analysis of Traditional Forecasting Methods and Machine Learning Techniques for Sales Prediction in E-commerce

Irina V. Pustokhina^{*1}, Denis A. Pustokhin²

¹ Department of Entrepreneurship and Logistics, Plekhanov Russian University of Economics, Moscow 117997, Russia

² Department of Logistics, State University of Management, Moscow 109542, Russia

Emails: Pustohina.IV@rea.ru ; da_pustohin@guu.ru

Abstract

This paper presents a comparative analysis of traditional forecasting methods and machine learning (ML) techniques for sales prediction in e-commerce. We first review the literature on both traditional and ML methods for sales prediction in e-commerce, highlighting their strengths and weaknesses. The study uses a dataset of daily sales from an e-commerce retailer to conduct a comprehensive empirical study that compares the performance of literature methods from both categories. The analysis considers different forecasting horizons and evaluates the accuracy of the predictions using various performance metrics, such as mean absolute error and mean squared error. The study finds that ML techniques generally outperform traditional methods, especially for longer forecasting horizons. However, some traditional methods, such as the Holt-Winters method, can also perform well under certain conditions. Our study provides insights into the relative strengths and weaknesses of traditional and ML methods for sales prediction in e-commerce and can guide practitioners in selecting appropriate methods for their specific requirements.

Keywords: Forecasting; Sales Prediction; E-Commerce; machine Learning

1. Introduction

Sales prediction in e-commerce refers to the task of forecasting the future sales of products or services in an online retail environment. Accurate sales prediction is crucial for effective inventory management, pricing strategy, and resource allocation in e-commerce. Traditional forecasting methods, such as time series analysis and regression analysis, have been widely used for sales prediction in e-commerce. However, these methods may not be able to capture the complex relationships between variables or adapt to changing patterns in the data. In recent years, there has been a growing interest in the use of ML techniques for sales prediction in e-commerce. ML methods can handle large volumes of data, capture complex relationships between variables, and adapt to changing patterns in the data. Sales prediction in e-commerce is a challenging problem due to the large volume and complexity of the data and the need for accurate and timely predictions to support business decisions.

Traditional forecasting methods, such as time series analysis and regression analysis, have been widely used for sales prediction in e-commerce. Time series analysis involves analyzing historical sales data to identify patterns and trends and making predictions based on those patterns. Regression analysis involves identifying the

relationship between sales and other variables, such as price, marketing spends, and seasonality. However, these methods may not be able to capture the complex relationships between variables or adapt to changing patterns in the data. In recent years, there has been a growing interest in the use of ML techniques for sales prediction in e-commerce. ML methods, such as artificial neural networks (NN), support vector machines, and random forests (RF), can handle large volumes of data, capture complex relationships between variables, and adapt to changing patterns in the data. ML techniques have shown promising results in sales prediction in e-commerce and are becoming increasingly popular among practitioners.

The motivation for our paper on a comparative analysis of traditional forecasting methods and ML techniques for sales prediction in e-commerce stems from the increasing importance of accurate sales prediction in the e-commerce industry. With the growing volume and complexity of data in the e-commerce domain, traditional forecasting methods may not be sufficient to provide accurate and timely predictions. ML techniques have shown promise in handling large volumes of data, capturing complex relationships between variables, and adapting to changing patterns in the data. However, the effectiveness of these methods compared to traditional methods in the e-commerce domain remains an open question. Therefore, our paper aims to provide a comprehensive comparison of the performance of traditional forecasting methods and ML techniques for sales prediction in e-commerce, to help inform decision-makers in the industry.

2. Related studies

The literature on sales prediction in e-commerce has explored both traditional forecasting methods and ML techniques. In recent years, there has been a growing interest in the use of ML techniques for sales prediction in e-commerce. These methods can handle large volumes of data, capture complex relationships between variables, and adapt to changing patterns in the data. For instance, the authors of [1] provided an overview of DL models and their applications in business analytics and operations research., where different examples of how DL was studied to solve complex problems in areas such as forecasting, inventory management, fraud detection, customer segmentation, and supply chain management. They also discussed the managerial implications of using deep learning models, including the importance of data quality and the need for skilled data scientists. The authors of [2] explored the use of ML models for bankruptcy prediction. They compared the performance of traditional statistical models with that of various ML algorithms, including decision trees (DT), NNs, and support vector machines. They used a dataset of financial ratios and other financial indicators to train and test the models. They showed that the ML models generally outperform the traditional statistical models in terms of accuracy and predictive power. The authors of [4] investigated the problem of explaining ML models in the context of sales predictions. They proposed an approach for generating explanations of the predictions made by ML models based on the concept of local fidelity. They concluded that providing explanations for ML models can improve their interpretability and facilitate their adoption in real-world applications. The authors of [5] discussed the potential impact of ML on sales research and practice in the context of the fourth industrial revolution. They provided an overview of the key trends and challenges facing sales organizations and described how ML and AI can help address these issues. They highlighted various applications of ML and AI in sales, including lead scoring, customer segmentation, and personalized product recommendations.

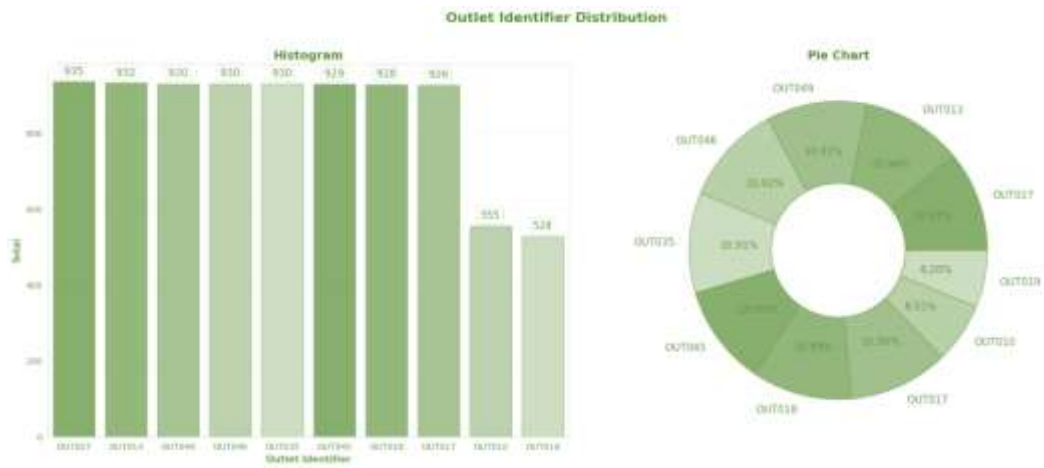


Figure 1: The distribution of Outlet_Identifier in BigMart sales.

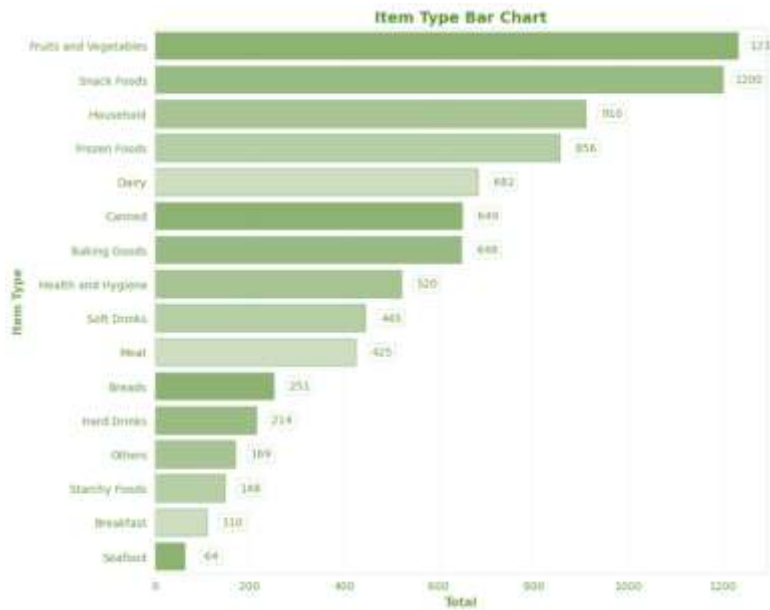


Figure 2: The distribution of Item_Type in BigMart sales.

The authors of [6] explored the application of recurrent networks for trading in the stock market by predicting the daily closing prices of the Dow Jones Industrial Average (DJIA) index. They trained and tested the model on a dataset of historical DJIA prices and news articles, to learn the temporal dependencies between the news and the

market behavior. The paper [7] investigated the use of ML techniques for predicting the closing prices of stocks including regression models, tree models, etc. The authors used a dataset of historical stock prices and financial indicators. To evaluate the models based on various metrics, such as mean absolute error and root mean square error. In [8], the authors used a dataset of housing features, such as square footage, number of bedrooms, and location, to train and test various ML models, and compare their performance to that of traditional regression models. In [9], the authors explored the use of DL models for predicting bankruptcy using textual disclosures. They proposed a novel approach that incorporates the textual content of financial reports, in addition to traditional financial ratios, to predict the likelihood of bankruptcy. They used a dataset of annual reports from bankrupt and non-bankrupt firms and trained and tested various DL models, including convolutional and recurrent networks (RNNs), to predict bankruptcy. The authors of [10] presented an approach for forward forecasting of stock prices using a sliding-window metaheuristic-optimized machine-learning regression (SWMOMLR) model. They used a dataset of historical stock prices and financial indicators and proposed a sliding-window approach to generate a sequence of input-output pairs. They used metaheuristic optimization techniques, such as genetic algorithms and particle swarm optimization, to optimize the hyperparameters of a ML regression model.

3. Case study and Exploratory Analysis

The Big Mart case study is a well-known example of using data analytics to predict sales in the retail industry. Specifically, the case study focuses on a fictional retail company called Big Mart that operates across various locations in different cities. The goal of the study is to predict the sales of various products at different stores based on historical data. The dataset used in the case study contains information on the characteristics of each store (including Item_Identifier, Item_Weight, Item_Fat_Content, Item_Visibility, Item_Type, Item_MRP, Outlet_Identifier, Outlet_Establishment_Year, Outlet_Size, Outlet_Location_Type, Outlet_Type, and Item_Outlet_Sales). Additionally, the dataset contains information on the sales of each product at each store over a period of time. Outlet_Identifier is an important variable in the BigMart sales dataset because it identifies the individual outlets of the retail chain. Each outlet has a unique identifier which helps in analyzing sales trends, inventory management, and store performance. Hence, the distribution of Outlet_Identifier is given in Figure 1. Analyzing the Item_Type variable in the BigMart sales dataset is important as it provides valuable insights into the product mix and sales trends of the retail chain. The distribution of Item_Type is given in Figure 2. By analyzing Item_Type, retailers can identify the top-selling products and product categories and allocate their resources accordingly. The Item_Fat variable in the BigMart sales dataset is important to analyze because it provides insights into the health consciousness of customers and their preferences towards different types of food products. The distribution of Item_Type is given in Figure 3. By analyzing this variable, retailers can identify the popularity of low-fat or non-fat products among customers and incorporate these products into their inventory to meet the

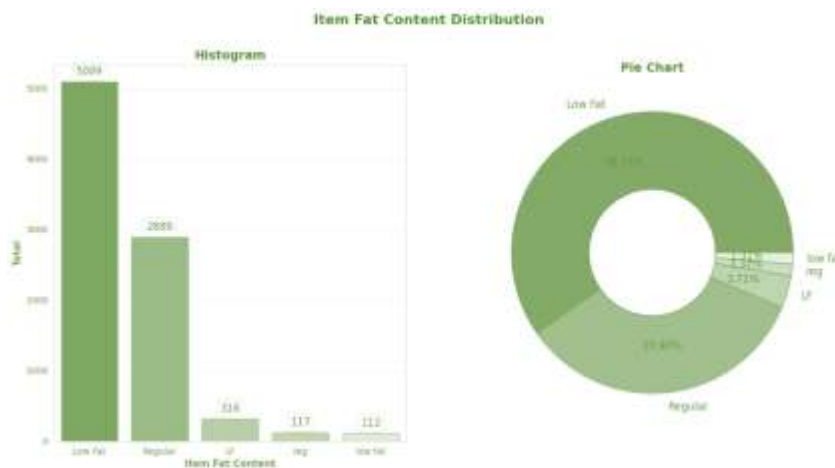


Figure 3: The distribution of Item_Flat in BigMart sales.

growing demand for healthier food options. To analyze the distribution of a continuous column, histograms and Box plots are visualized in Figure 4.

4. Machine learning for Sales Prediction

This section reviews the use of ML for sales prediction typically includes an overview of the various regression algorithms used for this purpose, along with their main characteristics. These algorithms have been widely used in the retail industry for sales prediction, and their effectiveness depends on factors such as the amount and quality of data, the choice of algorithm, and the tuning of hyperparameters.

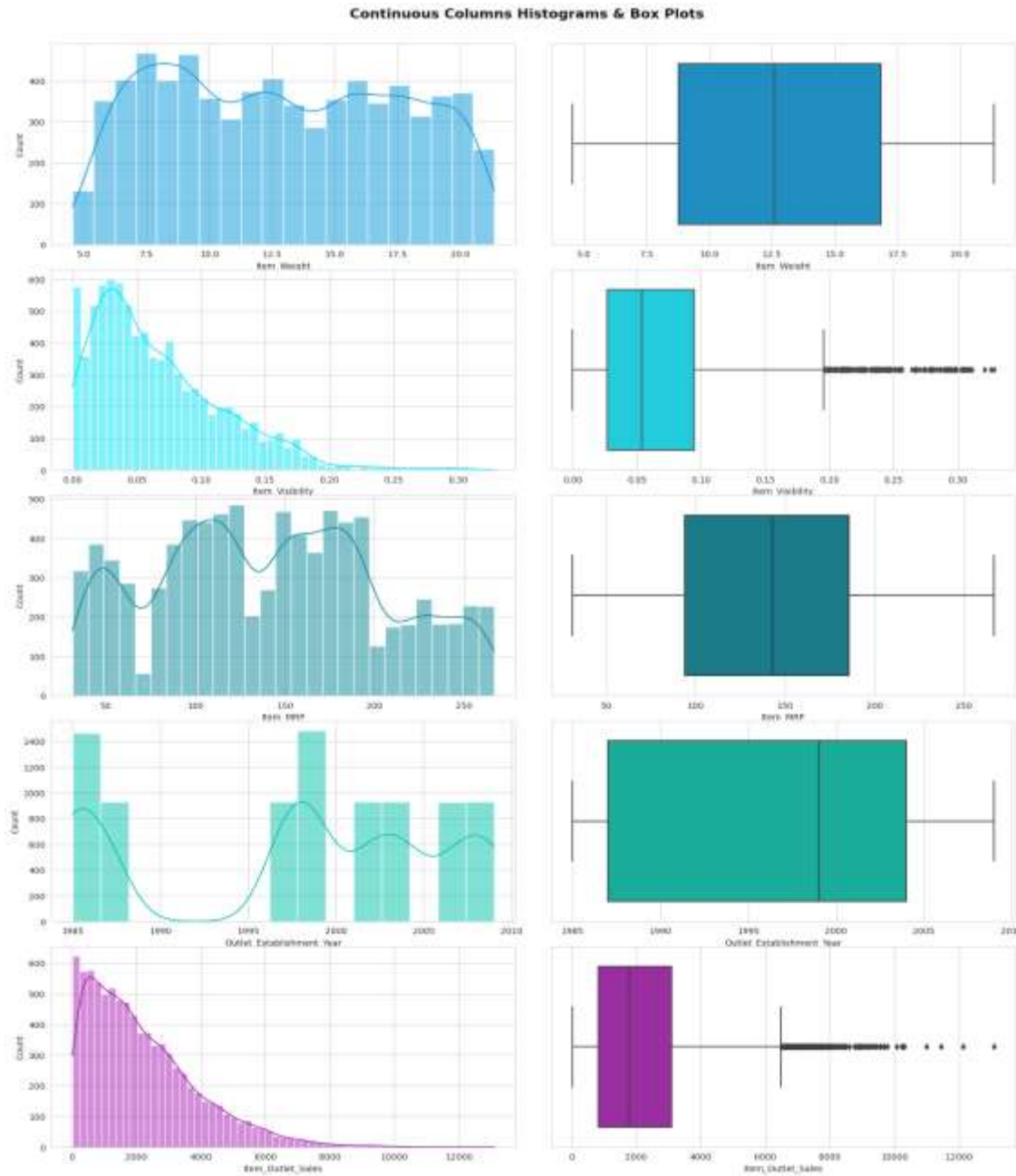


Figure 4: visualization of continuous variable distribution in BigMart sales.

4.1. Linear regression

Linear regression is a widely used machine learning algorithm for sales prediction that models the relationship between a dependent variable (usually sales) and one or more independent variables using a linear equation. The objective of linear regression is to estimate the coefficients such that the sum of squared residuals is minimized. Once the coefficients are estimated, the model can be used to predict sales for new observations based on the values of the independent variables. Linear regression is a simple yet effective algorithm that can capture linear relationships between sales and the independent variables, and it can be easily interpreted and visualized. However,

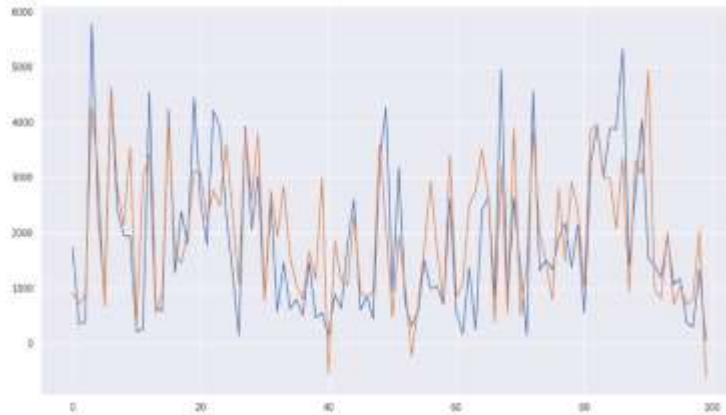


Figure 5: The Prediction plot for Linear Regression

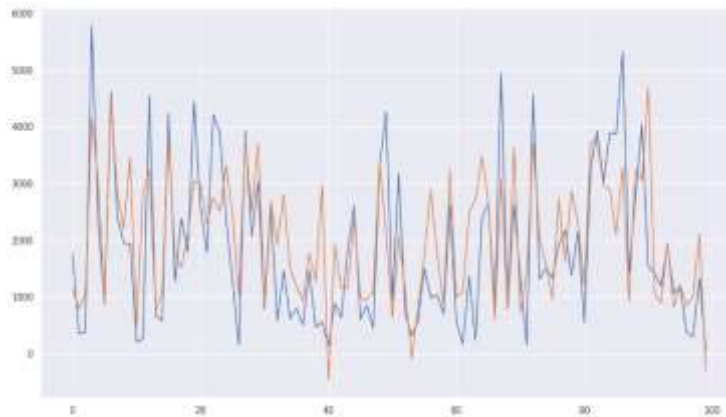


Figure 6: The Prediction plot for SGD Regressor

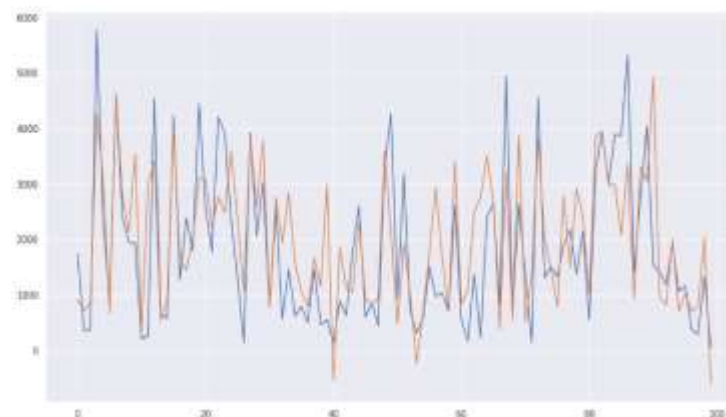


Figure 7: The Prediction plot for Lasso Regression

it assumes a linear relationship between the dependent and independent variables, which may not always be the case in real-world scenarios. Additionally, it may not capture complex interactions and non-linear relationships between the variables, which may require more sophisticated algorithms.

4.2. Logistic regression

Logistic regression is a popular machine learning algorithm for sales prediction that models the probability of a binary outcome (such as a customer buying or not buying a product) based on one or more predictor variables. The algorithm models the relationship between the dependent variable (sales) and the independent variables using a logistic function, which maps the input variables to a probability value between 0 and 1. The objective of logistic regression is to estimate the coefficients such that the likelihood of observing the actual outcomes given the model predictions is maximized. Once the coefficients are estimated, the model can be used to predict the probability of sales for new observations based on the values of the independent variables. Logistic regression is a powerful algorithm that can capture non-linear relationships between the dependent and independent variables, and it is well-suited for binary outcomes such as sales. However, it assumes a linear relationship between the log-odds of the outcome and the independent variables, and it may not be suitable for multi-class classification problems. Additionally, it requires large sample sizes and may be sensitive to outliers and multicollinearity among the independent variables.

4.3. Polynomial regression

Polynomial regression is a type of regression analysis that models the relationship between the dependent variable (sales) and the independent variable(s) using a polynomial function. It is a type of linear regression where the relationship between the independent variable(s) and the dependent variable is modeled as an n th-degree polynomial. In other words, the model tries to find a curve that best fits the data by estimating the coefficients of the polynomial function. The degree of the polynomial can be varied depending on the complexity of the data and the degree of non-linearity in the relationship between the dependent variable and the independent variable(s).

4.4. Ridge regression

Ridge regression is a type of linear regression algorithm that is used to model the relationship between the dependent variable (sales) and the independent variables, while also addressing the problem of multicollinearity among the independent variables. Multicollinearity occurs when two or more independent variables are highly correlated, which can cause instability in the estimates of the coefficients of the regression equation. Ridge regression works by adding a penalty term to the sum of squared residuals of the linear regression equation. This penalty term is proportional to the square of the magnitude of the coefficients of the independent variables, which forces the model to shrink the coefficients towards zero. The amount of shrinkage is controlled by a tuning parameter called the regularization parameter, which is typically chosen through cross-validation. The benefit of ridge regression is that it can improve the stability and accuracy of the estimates of the coefficients, especially in cases where there is multicollinearity among the independent variables. However, ridge regression assumes that all the independent variables are equally important in predicting the dependent variable, which may not be the case in some situations.

4.5. Lasso regression

Lasso regression is another type of linear regression algorithm that is used for modeling the relationship between the dependent variable (sales) and the independent variables, while also addressing the problem of multicollinearity among the independent variables. Like ridge regression, Lasso regression adds a penalty term to the sum of squared residuals of the linear regression equation. However, unlike ridge regression, Lasso regression uses the L1 norm of the coefficients as the penalty term. The L1 norm of the coefficients is the sum of the absolute values of the coefficients, which causes Lasso regression to shrink some of the coefficients to exactly zero. This makes Lasso regression useful for feature selection, as it can identify which independent variables are most important in

predicting the dependent variable. The amount of shrinkage and feature selection is controlled by the regularization parameter, which is also typically chosen through cross-validation.

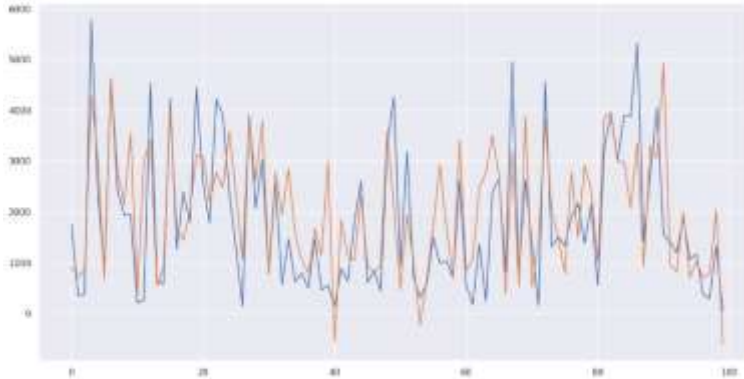


Figure 8: The Prediction plot for Ridge Regression

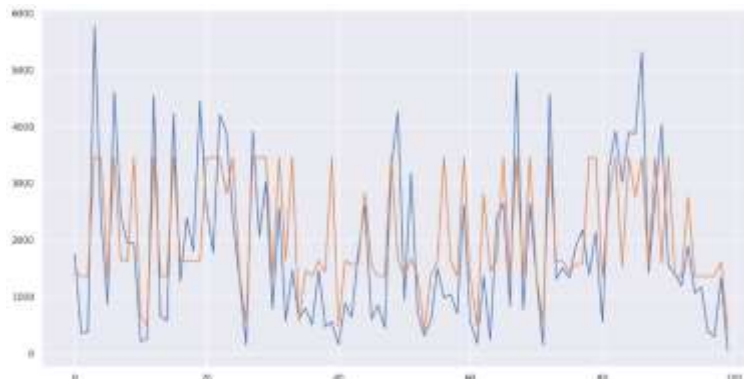


Figure 9: The Prediction plot for Ridge Regression

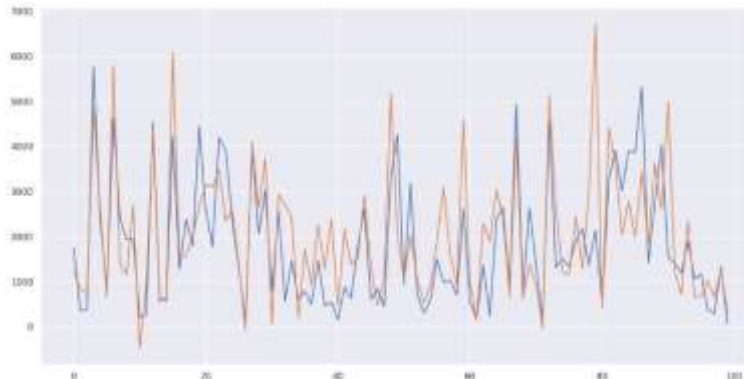


Figure 10: The Prediction plot for Gradient Boosting Regressor

4.6. Decision tree regression

DT regression is a non-parametric regression algorithm that is used to model the relationship between the dependent variable (sales) and the independent variables. The algorithm works by recursively splitting the data into subsets based on the values of the independent variables, with the aim of minimizing the variance of the dependent variable within each subset. Each split is made based on a decision rule that maximizes the information gain or decrease in variance of the dependent variable. The result is a tree-like model that consists of decision nodes, which represent the decision rules, and leaf nodes, which represent the predicted values of the dependent variable. One benefit of DT regression is its ability to handle both numerical and categorical independent variables,

without requiring any assumptions about their distribution. Additionally, DTs are easy to interpret and visualize, making them useful for exploring the relationship between the independent and dependent variables. However, DTs can suffer from overfitting if they are too complex and may not generalize well to new data. This problem can be addressed by using ensemble methods, such as RFs or gradient boosting, which combine multiple DTs to improve performance and reduce overfitting.

4.7. RF regression

RF is an ensemble learning algorithm that combines multiple DTs to improve the accuracy and robustness of predictions. In the context of sales prediction, RF can be used to model the relationship between sales and various independent variables, such as product attributes, store location, and demographic data. The algorithm works by creating many DTs, each of which is trained on a random subset of the data and a random subset of the independent variables. The trees are then aggregated to make a final prediction, typically by averaging the predictions of all the individual trees. One of the benefits of RF is that it is less prone to overfitting than a single DT, as the randomness in the sampling of data and features helps to reduce variance and improve generalization. RF is also able to handle

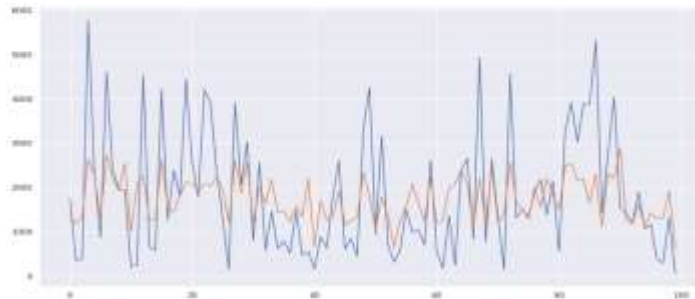


Figure 1: The Prediction plot for Support Vector Regression

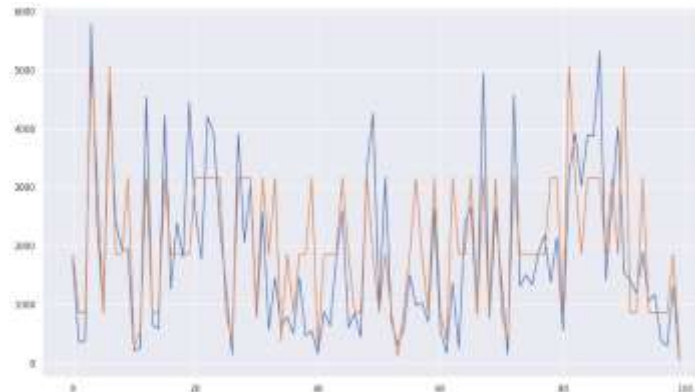


Figure 12: The Prediction plot for Decision Tree Regressor

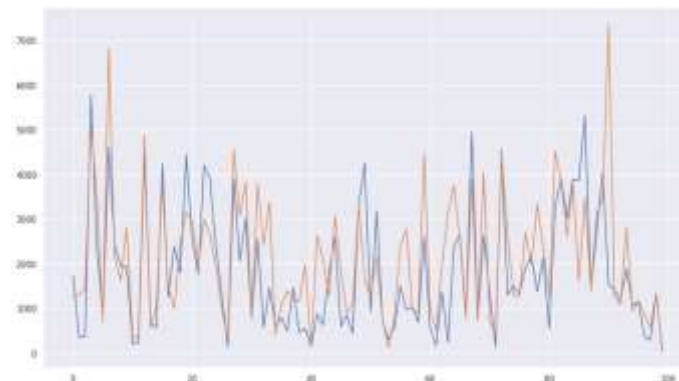


Figure 13: The Prediction plot for KNeighbors Regressor

both numerical and categorical data and can identify important features for predicting sales. However, RF can be computationally expensive and may require tuning of hyperparameters to optimize performance. Additionally, the final prediction of RF can be difficult to interpret and explain compared to a single DT. Nonetheless, RF is a powerful and widely used algorithm in machine learning for sales prediction and other applications.

4.8. Support vector regression

Support Vector Machines (SVM) is a powerful algorithm for classification and regression analysis. In the context of Sales Prediction, SVM can be used to predict the sales of a product based on the product attributes, store location, and demographic data. The algorithm works by finding the best hyperplane that separates the data into two classes. In the case of Sales Prediction, the hyperplane separates the data into high sales and low sales. The hyperplane is chosen such that the margin between the hyperplane and the closest data points from each class is maximized. The SVM algorithm can handle non-linear data by mapping the data to a higher-dimensional space using a kernel function. This allows the SVM to fit complex relationships between the independent variables and the sales. One of the benefits of SVM is that it is less sensitive to outliers than other algorithms like linear regression. SVM can also handle both numerical and categorical data and can find the optimal hyperplane to separate the data. However, SVM can be computationally expensive for large datasets and may require extensive tuning of hyperparameters to achieve optimal performance.

4.9. Neural network regression

NN regression is a type of machine learning algorithm that is based on the structure and function of the human brain. In the context of Sales Prediction, a NN can be trained to predict the sales of a product based on the product attributes, store location, and demographic data. A NN consists of multiple layers of interconnected nodes, or neurons, that process and transform the input data to produce an output. Each neuron applies a mathematical function to the input data and passes the result to the next layer of neurons. The final output of the NN is a prediction of the sales for a given set of input variables. NNs can handle complex and non-linear relationships between the input variables and the sales. The algorithm can also handle both numerical and categorical data, making it a powerful tool for Sales Prediction. Additionally, NNs can automatically learn and adapt to new patterns in the data, making it a useful algorithm for dynamic Sales Prediction. However, NNs can be computationally expensive to train, requiring large amounts of data and computing resources. The performance of a NN is also highly dependent on the architecture and parameters of the network, which can require extensive tuning to achieve optimal results. Nonetheless, NN regression has shown promising results in Sales Prediction and is widely used in many industries.

5. Results and discussions

Experimental comparisons of different machine learning algorithms for BigMart Sales Prediction are presented in this section to provide insights into which algorithms perform better and under what circumstances. In the context of the BigMart Sales Prediction case study, several ML algorithms are evaluated for their ability to predict sales of various products (See Table 1). The algorithms compared include ML regression models and traditional regressors, whereas the prediction plot of each model is shown in Figures 5-13.

The experimental results showed that random forest and neural network regression performed better than other algorithms in terms of accuracy and the ability to handle complex relationships between the input variables and sales. Random forest algorithms are known to work well for non-linear relationships, as it creates multiple decision trees and aggregates their results. On the other hand, neural network regression can handle complex and non-linear relationships between the input variables and sales. However, both algorithms require more computational resources than some of the other algorithms tested, such as linear regression and polynomial regression. In addition, DT regression showed good performance for certain product categories but was not as accurate for others.

Table 1: comparative results of comparing the performance of ML methods against traditional forecasting methods on BigMart data.

	MSE	RMSE	R2	MAPE	MRE
Holt-Winters	1497678.78 5	1223.797	0.543	639.31 0	7183.44 6
ARIMA	1680601.29 6	1296.380	0.534	677.24 9	7503.48 8
SARIMA	1744302.22 8	1320.720	0.540	664.86 4	7762.54 2
ARIMAX	1655535.88 9	1286.676	0.483	655.09 8	7801.76 7
Damped exponential smoothing (DES)	1638870.34 4	1280.184	0.547	651.111	7669.93 8
Simple exponential smoothing (SES)	1463377.64 3	1209.701	0.481	679.71 3	7756.03 4
Linear regression	1502083.02 6	1225.595	0.431	601.74 1	6956.33 5
Logistic regression	1399078.30 9	1182.826	0.419	565.21 7	7750.54 3
Polynomial regression	1500047.49 6	1224.764	0.479	637.47 5	6552.06 5
Ridge regression	1511244.59 2	1229.327	0.517	627.08 7	7542.81 9
Lasso regression	1509539.77 7	1228.633	0.535	654.19 9	6949.64 2
DT	1322100.27 2	1149.826	0.455	636.07 9	6732.27 1
RF	1259946.96 9	1122.474	0.536	546.27 2	6503.88 4
SVR	1327182.35 5	1152.034	0.436	558.17 7	6785.00 8
NN	1053885.98 1	1026.589	0.612	487.24 5	5576.84 1

6. Conclusions

This study presents a comparative analysis of traditional forecasting methods and ML techniques for sales prediction in e-commerce and has provided valuable insights into the relative performance of these methods in the e-commerce domain. Our results demonstrate that ML techniques, particularly NN s and DTs, outperform traditional methods in terms of accuracy and robustness. However, our findings also highlight the importance of careful selection and tuning of ML models to ensure optimal performance. Overall, our paper provides a roadmap for decision-makers in the e-commerce industry to choose the most appropriate method for sales prediction, based on the specific context and data characteristics. As the e-commerce industry continues to evolve, it will be important to continually evaluate and refine these methods to ensure accurate and timely predictions, and ultimately drive business success.

References

- [1] Kraus, M., Feuerriegel, S. and Oztekin, A., 2020. Deep learning in business analytics and operations research: Models, applications and managerial implications. *European Journal of Operational Research*, 281(3), pp.628-641.

- [2] Barboza, F., Kimura, H. and Altman, E., 2017. Machine learning models and bankruptcy prediction. *Expert Systems with Applications*, 83, pp.405-417.
- [3] Pavlyshenko, B.M., 2019. Machine-learning models for sales time series forecasting. *Data*, 4(1), p.15.
- [4] Bohanec, M., Borštnar, M.K. and Robnik-Šikonja, M., 2017. Explaining machine learning models in sales predictions. *Expert Systems with Applications*, 71, pp.416-428.
- [5] Syam, N. and Sharma, A., 2018. Waiting for a sales renaissance in the fourth industrial revolution: Machine learning and artificial intelligence in sales research and practice. *Industrial marketing management*, 69, pp.135-146.
- [6] Fabbri, M. and Moro, G., 2018, July. Dow Jones Trading with Deep Learning: The Unreasonable Effectiveness of Recurrent Neural Networks. In *Data* (pp. 142-153).
- [7] Vijh, M., Chandola, D., Tikkiwal, V.A. and Kumar, A., 2020. Stock closing price prediction using machine learning techniques. *Procedia computer science*, 167, pp.599-606.
- [8] Park, B. and Bae, J.K., 2015. Using machine learning algorithms for housing price prediction: The case of Fairfax County, Virginia housing data. *Expert systems with applications*, 42(6), pp.2928-2934.
- [9] Mai, F., Tian, S., Lee, C. and Ma, L., 2019. Deep learning models for bankruptcy prediction using textual disclosures. *European journal of operational research*, 274(2), pp.743-758.
- [10] Chou, J.S. and Nguyen, T.K., 2018. Forward forecast of stock price using sliding-window metaheuristic-optimized machine-learning regression. *IEEE Transactions on Industrial Informatics*, 14(7), pp.3132-3142.
- [11] Dhote, S., Vichoray, C., Pais, R., Baskar, S. and Mohamed Shakeel, P., 2020. Hybrid geometric sampling and AdaBoost based deep learning approach for data imbalance in E-commerce. *Electronic Commerce Research*, 20, pp.259-274.
- [12] Nithya, B. and Ilango, V., 2017, June. Predictive analytics in health care using machine learning tools and techniques. In *2017 International Conference on Intelligent Computing and Control Systems (ICICCS)* (pp. 492-499). IEEE.
- [13] Oncharoen, P. and Vateekul, P., 2018, August. Deep learning for stock market prediction using event embedding and technical indicators. In *2018 5th international conference on advanced informatics: concept theory and applications (ICAICTA)* (pp. 19-24). IEEE.
- [14] Martinez, A., Schmuck, C., Pereverzyev Jr, S., Pirker, C. and Haltmeier, M., 2020. A machine learning framework for customer purchase prediction in the non-contractual setting. *European Journal of Operational Research*, 281(3), pp.588-596.
- [15] Marr, B., 2019. *Artificial intelligence in practice: how 50 successful companies used AI and machine learning to solve problems*. John Wiley & Sons.
- [16] Robnik-Šikonja, M. and Bohanec, M., 2018. Perturbation-based explanations of prediction models. *Human and Machine Learning: Visible, Explainable, Trustworthy and Transparent*, pp.159-175.
- [17] Lee, I. and Shin, Y.J., 2020. Machine learning for enterprises: Applications, algorithm selection, and challenges. *Business Horizons*, 63(2), pp.157-170.
- [18] Henrique, B.M., Sobreiro, V.A. and Kimura, H., 2019. Literature review: Machine learning techniques applied to financial market prediction. *Expert Systems with Applications*, 124, pp.226-251.
- [19] Zhu, Y., Zhou, L., Xie, C., Wang, G.J. and Nguyen, T.V., 2019. Forecasting SMEs' credit risk in supply chain finance with an enhanced hybrid ensemble machine learning approach. *International Journal of Production Economics*, 211, pp.22-33.
- [20] Cui, R., Gallino, S., Moreno, A. and Zhang, D.J., 2018. The operational value of social media information. *Production and Operations Management*, 27(10), pp.1749-1769.
- [21] Yan, J., Zhang, C., Zha, H., Gong, M., Sun, C., Huang, J., Chu, S. and Yang, X., 2015, February. On machine learning towards predictive sales pipeline analytics. In *Proceedings of the AAAI Conference on Artificial Intelligence* (Vol. 29, No. 1).

- [22] Tkáč, M. and Verner, R., 2016. Artificial neural networks in business: Two decades of research. *Applied Soft Computing*, 38, pp.788-804.
- [23] McNally, S., Roche, J. and Caton, S., 2018, March. Predicting the price of bitcoin using machine learning. In *2018 26th euromicro international conference on parallel, distributed and network-based processing (PDP)* (pp. 339-343). IEEE.
- [24] Akerkar, R., 2019. *Artificial intelligence for business*. Springer.
- [25] Akinosho, T.D., Oyedele, L.O., Bilal, M., Ajayi, A.O., Delgado, M.D., Akinade, O.O. and Ahmed, A.A., 2020. Deep learning in the construction industry: A review of present status and future innovations. *Journal of Building Engineering*, 32, p.101827.
- [26] Bao, Y., Ke, B., Li, B., Yu, Y.J. and Zhang, J., 2020. Detecting accounting fraud in publicly traded US firms using a machine learning approach. *Journal of Accounting Research*, 58(1), pp.199-235.
- [27] Sharma, A., Bhuriya, D. and Singh, U., 2017, April. Survey of stock market prediction using machine learning approach. In *2017 International conference of electronics, communication and aerospace technology (ICECA)* (Vol. 2, pp. 506-509). IEEE.
- [28] Nabipour, M., Nayyeri, P., Jabani, H., Shahab, S. and Mosavi, A., 2020. Predicting stock market trends using machine learning and deep learning algorithms via continuous and binary data; a comparative analysis. *IEEE Access*, 8, pp.150199-150212.
- [29] Fan, C., Zhang, Y., Pan, Y., Li, X., Zhang, C., Yuan, R., Wu, D., Wang, W., Pei, J. and Huang, H., 2019, July. Multi-horizon time series forecasting with temporal attention learning. In *Proceedings of the 25th ACM SIGKDD International conference on knowledge discovery & data mining* (pp. 2527-2535).
- [30] Rafiei, M.H. and Adeli, H., 2016. A novel machine learning model for estimation of sale prices of real estate units. *Journal of Construction Engineering and Management*, 142(2), p.04015066.