

Hybrid Machine Learning Method for Product Sales Forecasting in E-Commerce

Ravi Kumar

Centre for Interdisciplinary Research in Business and Technology,
Chitkara University Institute of Engineering and Technology,
Chitkara University, Punjab, India
ravi.kumar.orp@chitkara.edu.in

Abstract - The term "e-commerce" refers to the practice of buying and selling products using the internet. E-commerce was created so that consumers wouldn't have to leave the comfort of their homes to purchase something. So, they can buy it online and have it delivered to their house in a few days. Online retailing has recently developed in a favorable environment because of the robust backing of national legislation. This year's pandemic has made the e-commerce sector's contribution to the growth of the country's economy more noticeable. The quantity and effectiveness of e-commerce networks and e-commerce businesses are growing under such circumstances. A network must've been able to better match user demands and perform admirably in all areas of collaboration and administration if it hopes to keep its competitive edge. Accurately predicting the sales volume of e-commerce networks at this time is crucial. This article presents a novel hybrid machine learning (ML) algorithm to accurately predict revenue growth for an online store. A hybrid cat-XGBoost (HC-XG) algorithm is suggested to predict the sale. The main goal of the proposed work is to forecast e-commerce sales, thus i the final selling value is included in the solution since it has a significant impact on sales in extremely competitive and price-sensitive contexts like e-commerce. We next continue with the development of the suggested solution and evaluate it for accuracy, error, and overall efficiency, and compare it to some competitor solutions. Thus, the proposed approach is compared with other existing models for sales forecasting and demonstrate that it provides the best overall performance in product sales forecasting.

Keywords- E-commerce, machine learning (ML), product sales forecasting, Hybrid cat-XGBoost (HC-XG)

I. INTRODUCTION

Online shopping has been more widespread over the last 2 decades as consumers have realized the convenience of not having to leave their residences to do their shopping. Recent years have seen an explosion in e-commerce, with many retailers making the switch to online marketplaces. E-commerce platforms make it harder to price things due to the large number of items offered. Pricing itself is often subject to change due to the several variables that affect it [1]. Managing a company's resources effectively is difficult in the ever-changing and complicated world of e-commerce. In response to these difficulties, several intelligent systems have been created, such as sales forecasting. The ability to accurately predict future sales is useful for optimizing manufacturing supply chains, allocating resources, and coordinating activities. Sales forecasts are only as good as the data they're based on. It is detrimental to E-commerce selection efficiency if stockouts or overstocks

occur due to inaccurate estimates. Common methods for predicting future sales have historically relied on a statistical analysis of time series, which uses only past sales data. Products having a consistent or seasonal sales pattern are easy for these strategies to manage. Nevertheless, the sales patterns of goods sold via e-commerce are substantially less predictable, making the accuracy of forecasts made using these approaches often inadequate [2]. A variety of data from several disciplines, including historical patterns, price, customer details, promotions, marketing channels, and changes in design, must be taken into account to develop sales projections. In addition, one must properly predict market trends, keep an eye on rivals, and take other company strategies into account. Sales of the product go through three phases over the long term: development, stabilization, and decrease. In the near run, they are impacted by factors including pricing, promotions, the period, and internet ranking. Sales variations are abrupt, brutal, and difficult to forecast, especially in e-commerce contexts where not all relevant data is accessible. Because of different possible assumptions, sales may exhibit a linear trend of rise or decline within a certain time, but some periods may exhibit the features of nonlinear volatility [3]. The process of predicting may make use of a wide range of methods, including methodological techniques, analysis of time series and extrapolation, and causality models. The first kind may or may not think about just the past, and relies on qualitative material like expert opinion and knowledge about extraordinary occurrences. In contrast, the second method is heavily reliant on past data since it is concerned only with tendencies and the alterations to those patterns. The third is sophisticated enough to explicitly account for exceptional occurrences and makes use of highly precise and particular knowledge regarding interactions among system parts. The past is crucial to forecasting models, just as it is to time series evaluation and contemplations [4]. The goal of every business, online or otherwise, is to serve the public by providing some kind of service or good. Hence, demand predictions will play a significant role in its manufacturing and subsequent decisions. In addition to boosting client satisfaction and profitability, arranging production and reducing stock on hand promptly may help businesses manage surplus inventory and set acceptable prices and discounts for products and services. Predicting sales will have an impact on transportation planning. E-commerce businesses are more able to adapt to changes in the economy and consumer preferences than their brick-and-mortar

counterparts. Thus, online stores must anticipate future sales volume [5]. In this work, a Hybrid Cat-XGBoost algorithm is suggested to predict sales.

The remaining sections of this article are as follows. The relevant research on sales prediction is covered in Section II, and the prediction model is presented in Section III. The outcomes of the review procedure are shown in Section IV. The study is concluded in Section V, which also highlights the work that will be done in the future.

II. RELATED WORKS

Reference [6] offered DSF, a revolutionary deep convolutional neural architecture for e-commerce sales prediction. In DSF, the process of estimating future sales is framed as an iterative sequential learning issue. To represent the effect of competitive relations when an increase in brand awareness is conducted for a target object or some exchangeable equivalents, they construct marketing communication over a network on top of the decoding. Substantial tests are carried out on two feature sets from the Taobao E-Commerce network, each representing a distinct area. Reference [7] includes the cross-series data into a single model since the objective and scope hierarchical on an e-commerce platform comprises a lot of related goods where the sales problem occurred might be associated. That was accomplished by worldwide training an LSTM that takes use of the non-linear demand connections present in E-commerce information in the attached hierarchy. In addition to the predicting methodology, they provide a methodical pre-processing foundation to address the problems in the e-commerce industry. Moreover, they provide a variety of product grouping techniques to support LSTM learning methods in cases when a product stock's sales trends are dissimilar. Reference [8] has proven that advertising is a useful tool for forecasting product demand in business-to-consumer (B2C) online marketplaces. In predicting the likelihood of a transaction, it stresses the significance of the product details, visuals, and promotional context. Artificial Neural Networks have been shown to perform better than linear regression techniques for load forecasting, as detailed in much previous research. The authors of that research look at how well linear programming, Classification Tree compositions, and deep learning techniques perform in comparison to one another. Reference [9] addressed current challenges in the manufacturing and global scale forecast of farming, a yield estimation model based on the Back Propagation Neural network algorithm is provided, as well as a set of techniques to optimize the BPNN. Due to the model's lengthy training period, it is straightforward to get into the locally optimal issue when using the typical BPNN. The density of the artificial neuron is decreased using the basic Johnson procedure, and the hidden state multilayer perceptron is constructed using that approach. Meanwhile, the data-gathering approach is utilized to sort out the information. The parameters are determined using the particle swarm optimization method. Reference [10] assists consumers in keeping tabs on the price of the product, that

research suggests a web tool called PriceCop. Because of its useful forecasting capabilities, users of PriceCop will be able to research products' prices for the following day and make more informed decisions. Linear Regression is used to create the price forecasting model. It is standard practice to utilize LR as a forecast and a means of determining the results of an experiment. To gauge how well the LR method performs, it is compared to Least Squares Support Vector Machine - Artificial Bee Colony. Traditionally, LSSVM-ABC was suggested for use in predicting stock prices. Reference [11] analyzed the K-Means network model that has been employed in that study to categorize clients who have efficiently purchased clothes products. PCA was utilized to reduce the dimensionality of a variety of consumer and product variables. The main aim of that research is to identify the many trademark, item, and pricing connections that customers may have based on their purchasing behaviors. The outcome demonstrates that the clusters created by the methodology based on PCA and K-Means are comparable, and the outcomes are appropriate based on comments from current consumers. They also satisfy the needs of the customers according to how much (price range) they wish to spend when they purchase online. Reference [12] assists fashion shops improve prediction performance, this research is based on load forecasting from an information approach, using both ML methods and finding relevant predictor factors. The effectiveness of ML methods was shown by comparing the prediction results that were achieved. A top fashion retail corporation used the suggested method to predict the sales for recently introduced seasonal goods in the absence of previous data [13], [14].

III. PROPOSED METHODOLOGY

The frameworks of Hybrid cat-XGBoost (HC-XG) ensemble learning approaches are suggested for forecasting the sale in this part. These collective learnings provide a methodical way to combine the prediction skills of several methods. Fig. 1 depicts the suggested systems' flowchart.

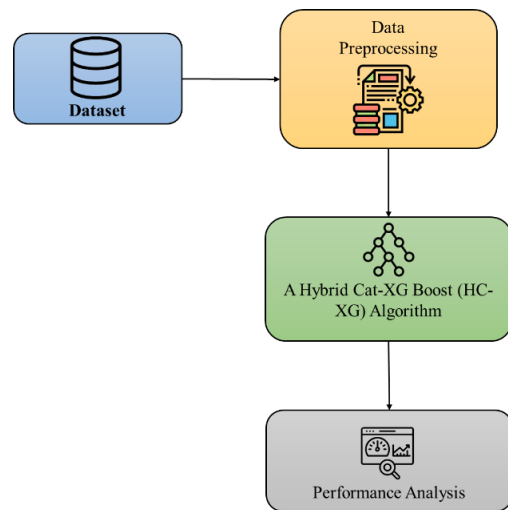


Fig.1 Framework of the Proposed Method

A. Data selection

The dataset for the research will come from commercialized data gathering. This is because we will be making a sales prediction, which requires the use of all previously collected data information. A small number of online shops will contribute to the source code, a publicly accessible operational dataset. The given commercial order historical data number is roughly 100,000. These statistics, which range from 2016 to 2018, are all supplied.

B. Preprocessing

The term "data preparation" refers to the steps used before actual analysis to eliminate any anomalies or abnormalities from the chosen data. This necessitates removing superfluous information from the data set, which has no practical use. In this database, for instance, product reviews are there, despite the fact that they aren't necessary for predicting sales. These reviews are, thus, noise that must be eliminated. Not only that, but we need to deal with missing transactions and price information in the right way by using the approximate value or the median interpolation to maintain data consistency if the dataset provides them. At this point, we may additionally take into consideration the data's chronological order and any previously established shifts.

C. A Hybrid cat-XGBoost (HC-XG) algorithm

A popular machine learning approach is XGBoost (eXtreme Gradient Boosting). It may be used for supervised learning including sorting, classifying, and forecasting. Designed to "push the limit of the processing boundaries of machines to produce a scalable, portable, and accurate library," it is constructed around the concepts of support vector architecture. For regression and classification problems, xgboost is an ML approach that creates a prediction model in the form of a group of weak forecasts, often decision trees.

The XGBoost replaces the search strategy by using the first and second generations of the error function. To boost the computation efficiency, it employs pre-ordering and base station techniques. The signal in the form is used at each iteration to exclude the weak classifier (decision trees) from the model. The XGBoost's objective function o is denoted as:

$$Q^{(d)} = \sum_{j=1}^m \left[f(x_i, \hat{x}_j^{(d-1)}) + \text{bld}(y_j) + \frac{1}{2} \text{oj}l_d^2(y_j) \right] + \Omega(l_d) \quad (1)$$

where j seems to be the j th specimen, d denotes the repetition, x_j is the actual value of the j th specimen, and $\hat{x}_j^{(d-1)}$ denotes the anticipated result of $d - 1$ th recursion; bi and oj denote both the 1st and 2nd derivatives; and $\Omega(l_d)$ denotes the term of normalization.

Further dealing with ML issues, XGboost, a gradient tree boosting method, is often used. Summarizing numerous tree classifications is the core notion underlying the gradients tree enhancing method.

The continuity formula describes the number of created characteristics, f , for a set k of training data:

$$M = \{(j_n, i_n)\} \quad (2)$$

Where j_n is written as j_n Gf, i_n as i_n R, and $|M|$ equals k . In addition, the standard HC-XGBoost model estimates the target value with the help of L exponential functions. Then,

$$\hat{i}_n = \theta(j_n) = \sum_{l=1}^L k_f(j_n), k_f \in K, \quad (3)$$

Each regression tree has a unique pattern, designated by bi , that may be utilized to convert training data into the corresponding position of a tree branch. U represents the whole collection of leaves on the tree. Each k_f has its own distinct bi and oj leaf values in a regression model.

The categorization of leaves for a particular training sample is achieved by following decision processes, and the final estimated output for the corresponding leaves is generated by summing the scores that are acquired from weights. Then, a regularised function is used to provide the group of functions used in this classification tree using the calculation below.

$$H(\theta) = \sum_n H(\hat{i}_n, i_n) + \sum_f \lambda(k_f) \quad (4)$$

This is the logistic regression design, which has a complication defined as,

$$\lambda(k) = \zeta U + (2)^{-1} \Gamma \|o\|^2 \quad (5)$$

Where H is a differentiable convex error term used to assess how far the estimated \hat{i}_n deviates from the original i_n . Using a complication factor to rule out the fitting problem helps ensure that the regularised value on the final predicted weights is as smooth as possible. In this case, the regularised algorithm chooses a basic approximated functional predictive tree model. In contrast to other classification trees, the effectiveness of this one is enhanced by the fact that its structure lends itself well to parallelization.

It is challenging to optimize the parameters of the logistic regression model described in equation (4) using conventional optimization algorithms. That's why we use an adaptive training scheme for multivariate regression. Assuming that, at the u^{th} iteration, the predicted output is $\hat{i}_n^{(u-1)}$ for the m th case, we know that the optimization of the regularised function necessitates the use of the variable k_f ,

$$H^u = \sum_{n=1}^r h(im, \hat{i}_n^{(u-1)} + k_u(j_n)) + \lambda(k_f) \quad (6)$$

Here, the make good decisions of the regression method are optimized with the help of the k_f parameter. The regularised function is optimized quickly using a 2nd approximation, as shown in the equation below.

$$H^u \approx \sum_{n=1}^r h([(im, \hat{i}_n^{(u-1)}) + e_n k_u(j_n)) + 2^{-1} p_n k_u^2(j_n)] + \lambda(k_f) \quad (7)$$

Equation (8) is simplified by removing the constant components to yield,

$$\tilde{H}^u = \sum_{n=1}^r [(e_n k_u(j_n)) + 2^{-1} p_n k_u^2(j_n)] + \lambda(k_u) \quad (8)$$

Then, calculate N_t as, for the system decided d of leaves.

$$N_t = \{m|b(j_n = t)\} \quad (9)$$

Next, we find by reducing equation (8) to its lowest terms that

$$\tilde{H}^u = \sum_{n=1}^r [(e_n k_u(j_n)) + 2^{-1} p_n k_u^2(j_n)] + \zeta U + (2)^{-1} \Gamma \sum_{t=1}^U o_t^2 \quad (10)$$

$$\tilde{H}^u = \sum_{t=1}^U [(\sum_{n \in N_t} e_n) o_t + 2^{-1} (\sum_{n \in N_t} p_n + \Gamma) o_t^2] + \zeta U \quad (11)$$

Then, using the corresponding optimization equation, get their final, best values:

$$\tilde{H}^u(b) = -(2^{-1}) \cdot \sum_{t=1}^U \frac{(\sum_{n \in N_t} e_n)^2}{\sum_{n \in N_t} p_n + \Gamma} + \zeta U \quad (12)$$

Specifically, where the scoring method shown above may be used to evaluate the merit of tree designs in equation (10). This score is compared to a wide variety of regularised variables and is employed in classifying tree models. In this case, the suggested tree model classifies the very first leaf of the classification trees before expanding it with further leaves. Just after a split, it is reasonable to assume that the moved node's case set is N_W and the correct node's case set is N_C . If we suppose that $N = N_W \cup N_C$, then the continuity formula describes the loss reduction term H_g after the split:

$$H_g = (2^{-1}) \cdot \left[\frac{(\sum_{n \in N_W} e_n)^2}{\sum_{n \in N_W} p_n + \Gamma} + \frac{(\sum_{n \in N_C} e_n)^2}{\sum_{n \in N_C} p_n + \Gamma} - \frac{(\sum_{n \in N} e_n)^2}{\sum_{n \in N} p_n + \Gamma} \right] - \zeta \quad (13)$$

In this case, the split alternatives in a boosted model may be evaluated using equation (13). By integrating the classifications of binary trees, the suggested boost model may also be used for a multi-class classification stage.

$$\frac{\partial \hat{f}}{\partial z} = \frac{\partial \delta(z)}{\partial z} = \delta(z)(1 - \delta(z)) = \hat{f}(1 - \hat{f}) \quad (14)$$

Equation (14) displays a sigmoid activation function characteristic and is utilized to further derive the loss function.

Algorithm 1: HC-XG Boost method

- Begin with the base node
- Evaluating each valid split that results in loss minimization while traversing all servers and variables on that specific point using the complexity search:

$$S = \text{loss}(L) - (\text{loss}(LR) + \text{loss}(RR))$$

where, S - attains a group, L - is an individual's root, LR - is the left root, RR - is the right root

- If the minimal split benefit (Hyperparameter) is not favorable and the split occurred on the strongest branch, then tree pruning is not appropriate.
 - Due to the nature of the issue, the hyperparameters object has been changed to reg: linear.
 - There is a 0.05 threshold for the learning algorithm.
-

1) Category Boosting

The analytics of randomization and the statistics of focusing on a certain target are the primary areas of interest in the Category Boosting method. It's useful for a broad variety of business problems and can handle data from several different sources. Converting categories to integers does not require any special preparation of the data. To quantify classifications, we employ an assortment of quantitative data in conjunction with our internal categorization characteristics. It's able to process huge datasets with little memory use. The performance of the classifier is reduced, and hence, more flexible algorithms are produced. A method of category classification known as Target Based with Condition precedent (TBS) is used to organize data. To mimic the structure of courses online, we linearly get training sets over the duration. Each simplicity's goal is derived entirely from empirical data. To capture higher-order correlations, we employ a mixture of classification features.

Here's how CatBoost functions:

- Constructing the data' subgroups
- Putting numerical values to the words used for categorization
- Quantifying the category characteristics

By using extremely randomized trees, as implemented by CatBoost, breaking criteria are consistently applied throughout all levels of the tree. Assessment prediction times may be greatly reduced with such a tree since it is well-balanced and resistant to over-matching.

IV. RESULTS AND DISCUSSION

This section includes HC-XG tests and a comparative analysis using various approaches. The approaches utilized for assessments include LSTM (Long Short -Term Memory), BPNN (Back Propagation Neural Networks), LSSVM-ABC (Least Squares Support Vector Machine with Artificial Bee Colony), and the suggested HC-XG. The simulation is carried out using a sales forecast. Comparing differences is done using the following metrics: Accuracy, R^2 , Mean Absolute Error, and Root Mean Square Error.

A) Accuracy

The standard value serves this function and is sometimes referred to as the corrected arithmetic mean and the reversed correctness of the advanced level. The accuracy with which the encrypted data is sorted may be used to gauge the reliability of a prediction system. Fig.2 displays the precision of the proposed and existing approaches. Table 1 displays the outcomes of the accuracy comparison.

The accuracy is expressed by,

$$Accuracy = TP + TN \div TP + TN + FP + FN \quad (15)$$

TABLE I. Comparison of Accuracy

Methods	Accuracy (%)
LSTM	75
BPNN	62
LSSVM-ABC	80
Proposed	97

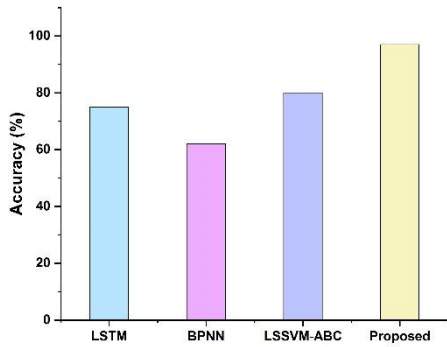


Fig.2 Accuracy Comparison

B) R^2

The estimation coefficient, often known as R^2 shows the fraction of the variation found in the dependent data that the linear regression framework explains. Since R is not based on a scale, its square root will always be smaller than one whether the parameters are small or large.

$$R^2 = 1 - \frac{\sum (x_j - \hat{x})^2}{\sum (x_j - \bar{x})^2} \quad (16)$$

TABLE II. Comparison of R^2

Methods	R^2 (%)
LSTM	80
BPNN	62
LSSVM-ABC	75
Proposed	95

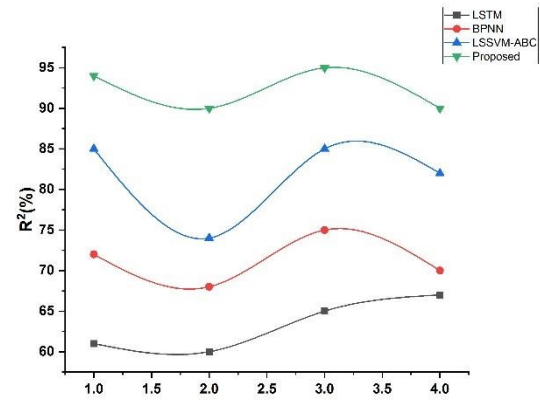


Fig.3 R^2 Comparison

The suggested and existing methods of r^2 comparison are shown in Fig.3. When compared with other existing methods our recommended method achieves a high r^2 value. Table 2 displayed the comparison value of r^2 .

C) Mean Absolute Error

The MAE is a statistical measure of how far off the mark the dataset's predictions were from the actual amounts. It is a statistical median of the dataset's error terms. MAE was expressed by,

$$MAE = \frac{1}{M} \sum_{j=1}^M |x_j - \hat{x}| \quad (17)$$

TABLE III. Comparison of MAE

Methods	MAE (%)
LSTM	92
BPNN	70
LSSVM-ABC	80
Proposed	62

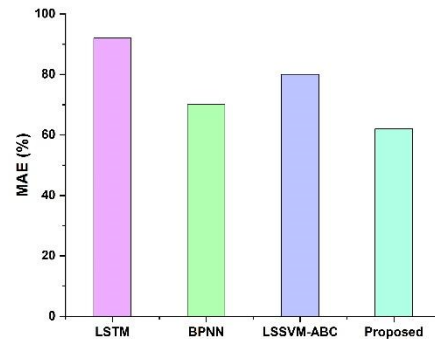


Fig.4 MAE Comparison

It was observed that the HC-XG Mean Absolute Error (MAE) was lower than that of the existing methods, indicating that the proposed method has a lower residual (prediction error). RMSE comparison graph was displayed in Fig.4 and Table 3 depicts the comparison analysis.

D) Root Mean Square Error

Frequently used metrics for gauging the accuracy of a forecast are the root mean square error or root mean square deviation. Using the Euclidean distance, it displays how much the predicted values deviate from the actual values. The RMSE is calculated by first finding the remaining (deviation between forecasting and fact) for each data point, then finding the norm of the error terms, then finding the average of the residuals, and finally finding the square root of the mean. Due to its reliance on measured data at every anticipated value, RMSE finds most of its usefulness in supervised learning settings. RMSE equation is expressed by,

$$RMSE = \sqrt{\frac{\sum_{j=1}^M \|x(j) - \hat{x}(j)\|^2}{M}} \quad (18)$$

TABLE IV. Comparison of RMSE

Methods	RMSE (%)
LSTM	90
BPNN	62
LSSVM-ABC	75
Proposed	55

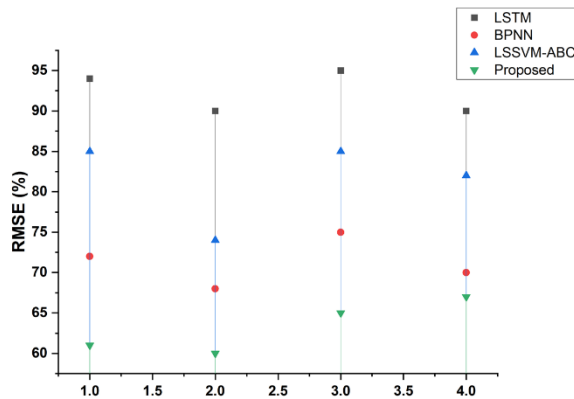


Fig.5 RMSE Comparison

It was demonstrated that our proposed method achieves less error when compared with other existing methods. Fig.5 displayed the comparison graph of RMSE and Table 4 shows the comparison values of the RMSE.

V. CONCLUSION

A lot of tiny firms now distribute their items on e-commerce platforms instead of opening physical stores because of the platform's explosive growth in the past ten years. Providing a tool that is simple to use and offers these tiny company owners a basic understanding of how to sell their items properly would be extremely helpful given their limited financial resources. Long-term economic growth would be aided by strengthening such enterprises, which would also encourage the establishment and growth of regional industries. The results show that the HC-XG Boost methodology is capable of delivering optimal price solutions with the least amount of error. However, the application of ensemble approaches produced output that was dependable, effective, and potentially even had a reduced error rate, giving consumers a successful outcome and a higher level of pleasure.

REFERENCES

- [1] A. Namburu, P. Selvaraj and M. Varsha, "Product pricing solutions using hybrid machine learning algorithm". Innovations in Systems and Software Engineering, pp.1-12, 2022.
- [2] K. Zhao and C. Wang, "Sales forecast in E-commerce using convolutional neural network". arXiv preprint arXiv:1708.07946, 2017.
- [3] K.N Vavliakis, A. Siallis, and A.L Symeonidis, "Optimizing Sales Forecasting in e-Commerce with ARIMA and LSTM Models". In WEBIST (pp. 299-306), 2021.
- [4] S. Steinker, K. Hoberg and U.W Thonemann, "The value of weather information for e-commerce operations". Production and Operations Management, 26(10), pp.1854-1874, 2017.
- [5] Z. Huo, "Sales prediction based on machine learning". In 2021 2nd International Conference on E-Commerce and Internet Technology (ECIT) (pp. 410-415). IEEE, 2021.
- [6] Y. Qi, C. Li, H. Deng, M. Cai, Y. Qi, and Y. Deng, "A deep neural framework for sales forecasting in e-commerce". In Proceedings of the 28th ACM international conference on information and knowledge management (pp. 299-308), 2019.
- [7] SivaramKrishnan, M., Rajini, A. R., Logu, K., Kumarasamy, M., Jayaprakash, S., Gandhi, R. R., & Ramkumar, M. S. (2022, October). Leaf disease identification using machine learning models. In AIP Conference Proceedings (Vol. 2519, No. 1, p. 050068). AIP Publishing LLC..
- [8] S. Rai, A. Gupta, A. Anand, A. Trivedi, and S. Bhaduria, "Demand prediction for e-commerce advertisements: A comparative study using state-of-the-art machine learning methods". In 2019 10th international conference on computing, communication and networking technologies (ICCCNT) (pp. 1-6). IEEE, 2019.
- [9] Y. Zhang, "Application of improved BP neural network based on e-commerce supply chain network data in the forecast of aquatic product export volume". Cognitive Systems Research, 57, pp.228-235, 2019.
- [10] M.Z Shahrel, S. Mutalib, and S. Abdul-Rahman, "PriceCop-Price Monitor and Prediction Using Linear Regression and LSVM-ABC Methods for E-commerce Platform". International Journal of Information Engineering & Electronic Business, 13(1), 2021.
- [11] S. Bandyopadhyay, S.S Thakur, and J.K Mandal, "Product recommendation for e-commerce business by applying principal component analysis (PCA) and K-means clustering: benefit for the society". Innovations in Systems and Software Engineering, 17(1), pp.45-52, 2021.
- [12] M. Kharfan, V.W.K Chan and T. Firdolas Efendigil, "A data-driven forecasting approach for newly launched seasonal products by leveraging machine-learning approaches. Annals of Operations Research, 303(1-2), pp.159-174, 2021.
- [13] Dogra, Aditya, Akshina Soni, and Yogeeshia Pai. "An experimental study on the mechanical properties Of basalt and banana fiber reinforced Hybrid polymer composites." International Journal of Mechanical and Production Engineering Research and Development 9.1 (2019): 263-270.

- [14] Khairandish, Mohammad Omid, R. Gurta, and Meenakshi Sharma. "A hybrid model of faster R-CNN and SVM for tumor detection and classification of MRI brain images." *Int. J. Mech. Prod. Eng. Res. Dev* 10.3 (2020): 6863-6876.