

Automobile Spare Parts Demand Forecasting Using Machine Learning

Oguzhan Akan, David Shin
California State University Northridge Computer Science

Index Terms—Machine Learning, Gradient Boosting, KNN, K Nearest Neighbors, Random Forest, Demand Forecasting, Spare Parts

I. INTRODUCTION

Automobile companies have various spare part products in their warehouses to cater towards customers' needs. The range of spare parts allow the industries to have modularity in their vehicles parts for repair and maintenance, however, it is difficult to predict when and what parts customers need. In order to address this issue, demand forecasting using machine learning helps industries recognize patterns to predict future sales. This allows their warehouses to be filled with parts customers will most likely need, reduce the shelf life of unnecessary parts and increase production in factories of parts that are highly requested. In return, production costs are lower and sales are higher and faster.

A. Related Work

As mentioned previously, in order to prevent shortage of products and to reduce wasted resources, existing works have used a wide range of tools to improve the performance of demand forecasting models. Also, common challenges for demand forecasting spare parts are the intermittent data and stochastic characteristics [?], [1], [2].

1) *Logistic Regression and Support Vector Machine Hybrid (LRSVM)*: The authors designed the LRSVM to "synthetically evaluate autocorrelation of demand time series and relationship of explanatory variables with demand of spare part." By implementing a hybrid model, LRSVM was able to perform better than the Croston's method, Markov bootstrapping method, IFM method and the SVM method. The hybrid method allowed for the SVM forecast results of nonzero occurrences to be replaced with explanatory variables to formulate a more accurate prediction [3].

2) *Grey's Forecasting Model (GFM) and BP Neural Network (BPNN) Hybrid*: In this research, a hybrid model using the Grey's Forecasting model and BP Neural Network. The GFM was used to produce new data series from the original data with reduce noise to be used for forecasting. The BPNN then uses the new data series to generate prediction results. By combining these two models, the authors were able to produce higher forecasting precision than the models executed separately [4].

3) *Fractional Order Discrete Grey Model (FOGM)*: Qiu et al proposed a derivation of the Grey's model, FOGM, which "... can effectively deal with the problem with "small sample" and "poor information", and can predict the data sequence with very few data samples." The genetic algorithm is used to estimate fractional order and generation coefficient for the FOGM to minimize the MAE of the predicted value. With this modification, the researchers are able to produce a model that is closer to the actual value than a single gray model [5].

4) *Support Vector Machine (SVM)*: In this research, the SVM was used to make predictions of spare parts for the Canadian Armed Forces. In comparison to the ARIMA, Naive, Croston and SES, the researchers concluded that the SVM had a lower average A-MAPE than the other models that were compared. They also found that the Asymmetric Error, which is designed to "... highlight the costs incurred when over-forecasting, and the loss of operational availability when under-forecasting (orders go unfilled)," was lower for the SVM for intermittent demand series than any other models [6].

5) *Gradient Boosting Model (GB) and Deep Neural Network (DNN)*: Authors of this experiment compared the GB and DNN models in retail for forecasting univariates sale time series. They proposed either of the two models can be relied on when there is a lack of causal factors or historical data. In this research, the GB performed better than the DNN and mentioned that the GB may have the tendency to overfit the data while the DNN were vulnerable to vanishing and exploding gradients [7].

B. Limitations

The limitations of our project include newly added products, missing parameters, outliers and the time restriction [8]. As a project is demand forecasting, newly added products may not have enough historical data to make an accurate prediction. Some of the missing parameters are price and size. This is a limitation to our project because these parameters may have strong correlations to sale orders. The price of an item may be correlated to sale orders because if the price of a product is expensive, the sale orders may be less. The size may affect sale orders because if the object is too small or too big, it may be ordered more or less respectively. And finally, the time restriction is a limitation to our project because we may not have as much time as we would like to train and test our models for improved performance.

C. Hypothesis

In our research, our objective is to conclude which of the three models: Gradient Boosting (GB), K-Nearest Neighbors (KNN) and Random Forest (RF) will perform better for demand forecasting automobile spare parts. Our research compares three popular models and analyze them comprehensively for future researchers to have a direct comparison of these models in demand forecasting automotive spare parts.

We have chosen the Gradient Boosting and Random Forest methods for this experiment because they are ensemble models. In our research, we found that hybrid models produced more effective results, which led us to believe that ensemble models would produce more positive results than using a single model [3]–[5]. To further verify this, we have selected the K-Nearest Neighbors method as a comparison method. With this knowledge, we predict the Gradient Boosting method will outperform the other models.

II. DATASET

The dataset used in this research is obtained from Dogus Automotive, a private automobile distributor based in Turkey. The company distributes cars of many brands including Volkswagen, Seat, Audi, Bentley, and Lamborghini and their spare parts. Dogus made its spare part catalog and order data public for their data competition in 2020, so our dataset is publicly available.

A. Dataset Components

The dataset is made up of two tabular-format tables which can be joined together using **part ID** column. Here are the details of the tables:

- **catalog:** The catalog table consists of information related to each automobile part. There are 5,000 parts in this table and it has 5 descriptive features for each part: *part_definition_id*, *part_product_class_id*, *common_part_catalog_id*, *preferred_supplier_id*, *part_family_id*. Each descriptive feature refers to some ID number sourced in other tables that are not available. Therefore, we decide to use the ID numbers scratch in the feature selection phase. In order to decide the approach we want to go by with each feature, we first look how many unique values exist in each of them:

- **orders:** The orders table consists of information about orders having 206,306 rows of data. A sample order information can be as simple as the related *part_id*, *order_quantity*, *date*, and *firm_id*. Therefore, we do not have any descriptive features in this table yet. Still, we can create features from scratch using Time Series Analysis as the history of each parts' orders matters.

III. METHODOLOGY

For our research we selected three algorithms with one being a base estimator and other two being different kinds of ensemble Machine Learning algorithms which we found there were not many researches using these algorithms in this field. The three algorithms we will be experimenting in this project

are K-Nearest Neighbors, Random Forest and Gradient Boosting. To implement the research, we used Jupyter Notebooks in Python 3 with the data manipulation and modeling tools such as Sklearn, Pandas, and Numpy. For data visualization we made use of Seaborn and Matplotlib.

A. Algorithms

1) *K-Nearest Neighbors:* K-Nearest Neighbors (KNN) is a classification technique that uses a given set of "neighbors" to classify new data. The K represents the number of nodes the algorithm will use to classify a new data and it determines the new data classification by finding the distance of the new node's neighbors [9].

2) *Random Forest:* Random Forest is an ensemble learning method that utilizes decision trees to make predictions by taking either the mode of the decision trees or the mean prediction of each tree [10].

3) *Gradient Boosting:* Gradient Boosting is a machine learning algorithm for regression and classification problems. It produces predictions, usually by using decision trees, and minimizes error by modifying the weights of the trees that produce the best outcomes [11].

B. Validation Strategy

We decided to move forward with Time Series Cross Validation (TSCV) method in this research as it fit the ML goal's better in terms of seasonality. In TSCV, we select K folds just as we do in the K-fold Cross Validation, but we separate the folds according to a Time Series split instead of a random and/or shuffled split [12]. The intuition behind this choice is that we see a gradually decreasing trend in orders table (Figure 1).

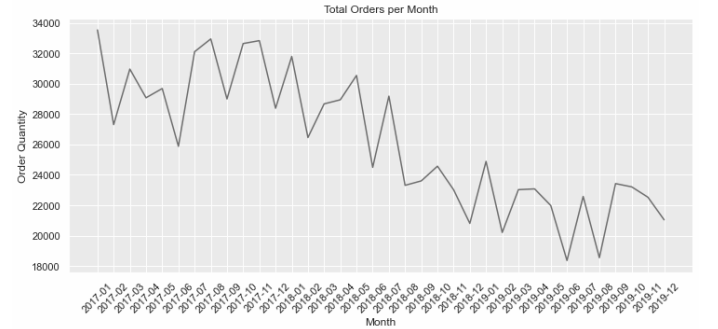


Figure 1. Decreasing trend in total order quantities with respect to each month.

Now that we have defined our validation strategy, we can select the time ranges we want build the models on. We selected 3 Folds with each fold's validation period being 3 months (Figure 2), based on the fact that 3 months of order data consisting of 15,000 rows is an adequate number to objectively evaluate performance metrics.

Fold\Period	2017-01 2019-03	2019-04 2019-06	2019-07 2019-09	2019-10 2019-12
1	TRAINING	VALIDATION		
2	TRAINING		VALIDATION	
3	TRAINING			VALIDATION

Figure 2. Time Series Cross Validation (TSCV) with 3 folds used in this research.

C. Preprocessing

1) *Handling Outliers*: Before applying bins to the data, a clamp transformation has been applied to it as the data had outliers. We cut the data into bins while setting the maximum quantity to 72. That is because after observing the frequencies, values greater than 72 can be considered as outliers. Since the lowest value of an order quantity can be minimum zero, no transformation is applied to the lower bound. After applying clamp transformation to the target feature, *order_quantity*, it is seen that the target is much skewed to the left (positively skewed) with a skewness value of 4.51 (Figure 3). We did not apply any normalization to this feature as it is the target feature we want to predict, thus cannot be transformed to any other scale.

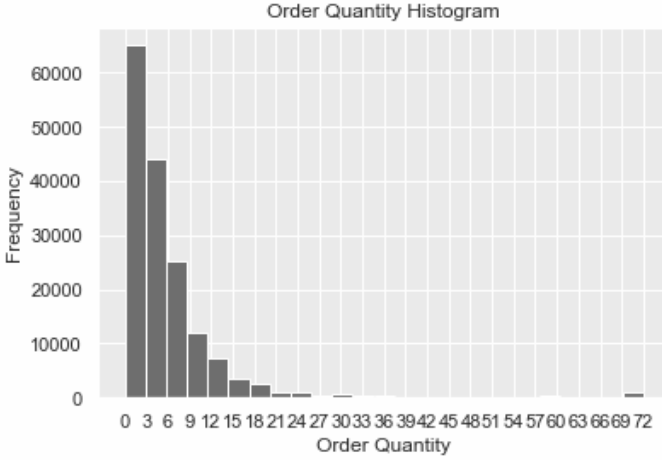


Figure 3. *order_quantity_bin* histogram graphic that shows the skewness in the data.

2) *Data Imputation*: Since the *orders* table has the rows at which some order exists, it does not give an information about the days where there is no order, which may be an insightful detail when building the model. Therefore, we inserted zero orders to the days where there is no order for a specific *part_id*. After this imputation, the dataset's initial size of 206,306 rows increased to 237,805.

3) *Handling Duplicate Rows*: Since the aim of the project is to predict the monthly orders for each *part_id*, we need to have exactly 1 row for each part in the training data. Therefore, the raw data, where there may be more than one order for one part in some month, is grouped by the *part_id* and returned their sum or *order_quantity*. Ultimately, our training data consisted of 167,059 rows.

4) *Time Series Analysis*: "Time series analysis is a method to establish a stochastic model for time series data based on its property, and utilizes the stochastic model to predict the long term trend" [13]. In order to predict the order quantities for a given month, it is a good practice to generate time series features in addition to the descriptive features of *part_ids*. We generate lag-N features where N represents some number of months. For example, *order_quantity_lag3* represents the order quantity 3 months before. Then we take the sum, mean, and standard deviations of lag-N features for $1 \leq N \leq 6$ to test which lag best explains the target feature.

5) *Clustering Categorical Features*: The descriptive features of *part_ids* are in ID format, meaning that they cannot be

used for modeling purposes as a lot of Sklearn's algorithms cannot handle categorical features. Thus, we need a way to convert these categorical features into continuous features. We use Binary Encoding, which encodes information to a bit, for the features that have a few number of categories [14]. For the rest, we will apply a K-means clustering with $K=5$. K-means clustering is an unsupervised learning algorithm that groups data points into K clusters by calculating their Euclidean distance [15]. The algorithm first randomly creates K center points which are called centroids, then in every iteration, it assigns the data points to the centroid which has minimum distance to it.

The reason why we applied K-means clustering to some of our categorical features is avoiding the possible overfitting during modeling. As it can be seen in Table I, there are too many distinct values exist for *part_definition_id*, *part_product_class_id*, and *part_family_id* that can cause the model to overfit.

feature	# of Unique Values
<i>part_definition_id</i>	1918
<i>part_product_class_id</i>	446
<i>common_part_catalog_id</i>	29
<i>preferred_supplier_id</i>	5
<i>part_family_id</i>	225

Table I: Number of unique values exist in the data for each descriptive feature of a *part*.

Therefore, we used normalized order quantity averages and standard deviations of each unique value in each feature to assign them clusters. Figure 4 shows before and after clustering of the three features.

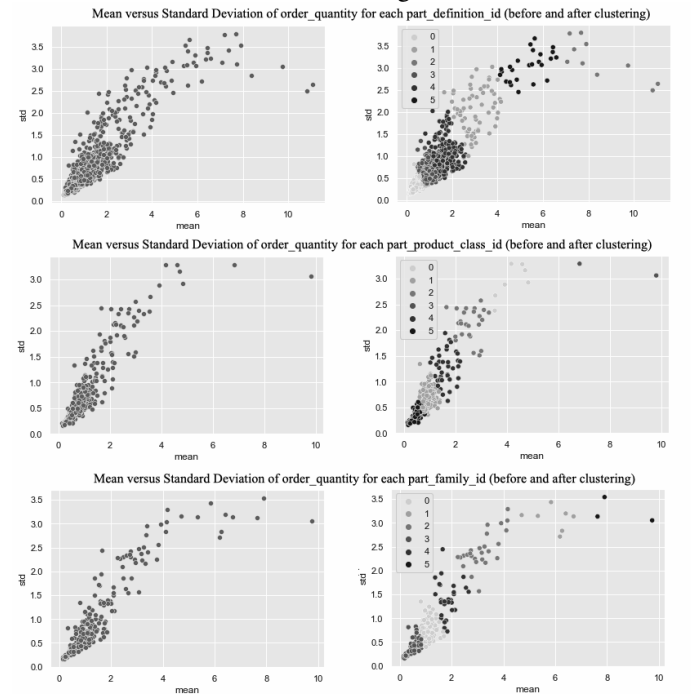


Figure 4. Scatterplot of three descriptive features, mean versus standard deviation of *order_quantity*, before and after K-means clustering.

6) *Binary Encoding Categorical Features*: As we clustered three of five descriptive features of *parts*, we need to transform the remaining two features, *common_part_catalog_id* and *preferred_supplier_id*, into either a continuous variable or a binary variable in order to be able to use them for modeling. We chose to convert *common_part_catalog_id* with a default transformation as it still have a lot of unique values. *part_family_id*, however, can be transformed to a binary variable using binary encoding. Binary encoding is a categorical feature transformation method that creates N new binary features where N is the number of unique values for that categorical feature [14]. Thus, in our case, we create 5 new features for *preferred_supplier_id* and specify them with a suffix *_binenc[N]*.

7) *Handling Null Values*: When we create new features, particularly in Time Series Analysis, some values have null value due to non-existent previous data. Therefore, we need to fill those null values in order to build the models. As a common approach, we first tried to fill the nulls with their category means. After model trials, however, we decided to continue with filling with category maximums because we see a decreasing trend in order quantities. The intuition behind is that the lag- N features bring the order quantities N months ago, and they can be even higher than the category maximum for this specific case.

8) *Final Training Data*: The training data consists of 167,059 rows with the target feature, *order_quantity_bin*, and descriptive features:

- *time_diff_mean*
- *time_diff_std*
- *order_quantity_lag[N]* where $1 \leq N \leq 6$
- *order_quantity_lag[N]_mean* where N in (3, 6)
- *order_quantity_lag[N]_std* where N in (3, 6)
- *pdid_cluster*
- *pfid_cluster*
- *pcid_cluster*
- *psid_binenc[N]* where $1 \leq N \leq 5$
- *psid*
- *cpid*

D. Feature Selection

Before building the models we want to eliminate some of the features we have if they are highly correlated with each other. But how do we decide which one to drop if we have two correlated features? In order to build better models, we need to know which one is describing the outcome better than the other. Therefore, we first use a SVM algorithm from Sklearn's Linear Support Vector Regression module in order to get coefficients of each feature. After applying the SVM, we see the most and least correlated features with the target feature (Figure 5).

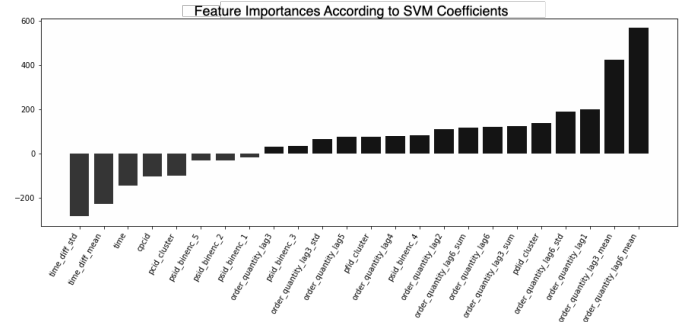


Figure 5. Bar chart of feature importance according to Support Vector Machine Regressor coefficients.

Now that at least we know which feature performs better than the other, we can build a correlation matrix and determine a maximum threshold that can two features be correlated. We selected 0.90 as the threshold and starting from the most important feature, we scanned the rest of the features and excluded them when necessary. Figure 6 shows the correlation matrix where darker colors represent higher correlation.

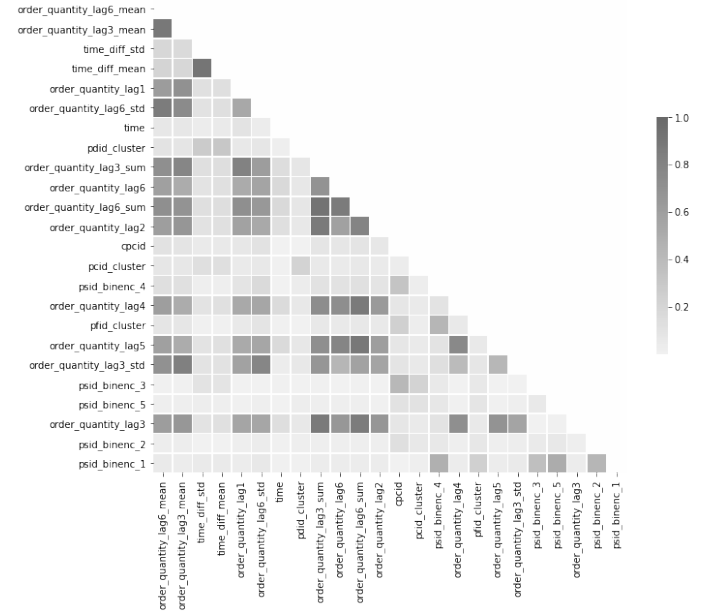


Figure 6. Correlation matrix of features ordered from most important to least important.

After running the scan, we drop the following features due to high correlation with other features: *time_diff_mean*, *order_quantity_lag6_sum*

IV. PRELIMINARY RESULTS

After inserting missing rows, correcting the data, clustering categorical features, filling null values, dropping correlated features, deciding on the modeling algorithms and validation strategy, we finally build the KNN, RF, and GB models. We use Sklearn's KNeighborsRegressor, RandomForestRegressor, and GradientBoostingRegressor modules respectively to build our models. There are many parameters used in each model, but in this research we will focus only on some of them for the sake of simplicity.

For evaluating performances of each type of model, we decided on considering Mean Absolute Error (MAE) and Mean Squared Error (MSE) as these metrics are commonly

used, and MAE ranges is in the target's scale, thus making it easier to interpret the results.

A. KNN Model

In the KNN model we set $n_neighbors$ parameter to 5 in order to limit computing power. The remaining parameters are in their default settings. For example, the $weights$ parameter is set to 'uniform', meaning that all the points in each neighborhood are weighted equally.

After running the KNN model, we obtained a MAE of 3.8670 from our validation data. Table II represents the metrics details against each fold.

Dataset_Metric	Fold 1	Fold 2	Fold 3	Average
mse_train	37.58	38.04	38.69	38.10
mae_train	3.48	3.52	3.56	3.52
mse_valid	49.83	48.11	47.51	48.48
mae_valid	3.92	3.81	3.85	3.86

Table II: Training and validation results of KNN model.

B. Random Forest Model

RF model can have many parameters regarding tree creation or ensembling, and here are our settings different from the Sklearn's default parameters:

- $n_estimators=100$
- $criterion='mse'$
- $max_depth=7$
- $min_samples_split=100$
- $min_samples_leaf=40$
- $max_leaf_nodes=45$

After running the RF model, we obtained a MAE of 3.6239 from our validation data. Table III represents the metrics details against each fold.

Dataset_Metric	Fold 1	Fold 2	Fold 3	Average
MSE (train)	45.69	46.32	47.11	46.37
MAE (train)	3.95	3.99	4.04	4.00
MSE (validation)	43.01	40.08	41.00	41.36
MAE (validation)	3.66	3.57	3.62	3.62

Table III: Training and validation results of RF model.

C. Gradient Boosting Model

In GB model we set the parameters as the same as they are being used in RF model. Since the two algorithms are both based on decision trees, they have similar parameters in creating the trees. They differ in the way they ensemble, and we let them be in their default for now.

After running the GB model, we obtained a MAE of 3.8039 from our validation data. Table IV represents the metrics details against each fold.

Dataset_Metric	Fold 1	Fold 2	Fold 3	Average
MSE (train)	50.33	51.31	52.33	51.32
MAE (train)	4.14	4.20	4.27	4.21
MSE (validation)	44.89	42.86	43.37	43.71
MAE (validation)	3.79	3.75	3.86	3.80

Table IV: Training and validation results of GB model.

V. DISCUSSION

A. Model Comparison

1) *Performance*: Contrary to KNN, RF and GB models performed better in validation sets as opposed to training sets. This occurred because of the decreasing trend in the data. Although we divided the training and validation sets using time split, we needed to include more features that are related to seasonality. Nonetheless, both RF and GB performed better than KNN in validation sets.

As we draw line plots of each model's predictions and the actual outcome, we see a very similar trend in RF and GB, while KNN's trend is closer to the actual outcome (Figure 7). Although KNN's trend is similar to actual outcomes, we must do the comparison using the actual MAE's of the models.

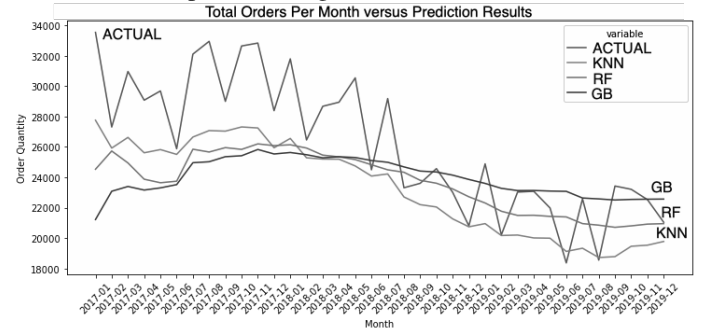


Figure 7. Prediction results and actual outcome for KNN, RF, GB models.

Our initial hypothesis suggested that GB would perform better than RF but when their MAEs compared, RF seemed to perform better. This may currently proves that RF is a better intermittent demand forecasting algorithm, although the preliminary results lack the usage of hyperparameter optimization. Therefore, this research can be iterated over optimizing the parameters and further obtaining metrics.

2) *Efficiency*: It is well-known that ensemble models are expected to work slower than the base estimators. As we calculated the time spent to fit models, we saw that KNN model fitted in 5.053 seconds, while RF fitted in 58.620 seconds, and GB fitted in 103.104 seconds. Nonetheless, we can increase or decrease fitting times of each model by changing parameter. The question we are interested in is "How can we shorten the duration without losing validation performance?"

B. Novelty of the Reserach

1) *Machine Learning Algorithms*: After a careful literature review, it is observed that there are few previous researches done on Gradient Boosting Regression on demand forecasting. This research proves that Gradient Boosting is also an applicable algorithm for this subject area.

2) *Data Preprocessing Steps*: We generally observed Time Series Analyses in the previous research. In addition to Time Series, we applied (i) clustering the categorical features, (ii) inserting zero orders for the days that there were no orders, and (iii) Binary Encoding of the rest of categorical features.

VI. CONCLUSION

In conclusion, based on the research conducted, we cannot accept the hypothesis that suggests GB performs better in terms of MAE than RF and KNN. Although this conclusion is correct for the moment, it is not trivial. Further modeling techniques such as feature selection and hyperparameter optimization need to be conducted in order to reach a trivial conclusion.

REFERENCES

- [1] I. RuiRui Xing, Xianliang Shi, Member, "A BP-SVM combined model for intermittent spare parts demand prediction," pp. 1085–1090, 2019.
- [2] Q. Xu, N. Wang, and H. Shi, "A Review of Croston's method for intermittent demand forecasting," *Proceedings - 2012 9th International Conference on Fuzzy Systems and Knowledge Discovery, FSKD 2012*, no. 61100009, pp. 1456–1460, 2012.
- [3] Z. Hua and B. Zhang, "A hybrid support vector machines and logistic regression approach for forecasting intermittent demand of spare parts," *Applied Mathematics and Computation*, vol. 181, no. 2, pp. 1035–1048, 2006.
- [4] H. Song, C. Zhang, G. Liu, and W. Zhao, "Equipment spare parts demand forecasting model based on grey neural network," in *2012 International Conference on Quality, Reliability, Risk, Maintenance, and Safety Engineering*, 2012, pp. 1274–1277.
- [5] Q. Qiu, C. Qin, J. Shi, and H. Zhou, "Research on demand forecast of aircraft spare parts based on fractional order discrete grey model," in *2019 IEEE 5th International Conference on Computer and Communications (ICCC)*, 2019, pp. 2212–2216.
- [6] A. Boukhtouta and P. Jentsch, "Support vector machine for demand forecasting of canadian armed forces spare parts," in *2018 6th International Symposium on Computational and Business Intelligence (ISCBI)*, 2018, pp. 59–64.
- [7] K. Wanchoo, "Retail Demand Forecasting: a Comparison between Deep Neural Network and Gradient Boosting Method for Univariate Time Series," pp. 5–9, 2019.
- [8] B. M. Pavlyshenko, "Machine-learning models for sales time series forecasting," *Data*, vol. 4, no. 1, pp. 1–11, 2019.
- [9] S. Khan, Z. A. Khan, Z. Noshad, S. Javaid, and N. Javaid, "Short Term Load and Price Forecasting using Tuned Parameters for K-Nearest Neighbors," *IIT 2019 - Information Technology Trends: Emerging Technologies Blockchain and IoT*, pp. 89–93, 2019.
- [10] A. Lahouar and J. Ben Hadj Slama, "Random forests model for one day ahead load forecasting," *2015 6th International Renewable Energy Congress, IREC 2015*, 2015.
- [11] M. Gumus and M. S. Kiran, "Crude oil price forecasting using XGBoost," *2nd International Conference on Computer Science and Engineering, UBMK 2017*, pp. 1100–1103, 2017.
- [12] R. Medar, V. S. Rajpurohit, and B. Rashmi, "Impact of Training and Testing Data Splits on Accuracy of Time Series Forecasting in Machine Learning," *2017 International Conference on Computing, Communication, Control and Automation, ICCUBEA 2017*, no. April 2020, pp. 1–6, 2018.
- [13] Tingting Huang, L. Wang, and T. Jiang, "Prognostics of products using time series analysis based on degradation data," in *2010 Prognostics and System Health Management Conference*, 2010, pp. 1–5.
- [14] M. D. Todd, *Sensor data acquisition systems and architectures*. Woodhead Publishing Limited, 2014, vol. 1. [Online]. Available: <http://dx.doi.org/10.1533/9780857099136.23>
- [15] S. Nuchprayoon, "Electricity load classification using K-means clustering algorithm," *IET Conference Publications*, vol. 2014, no. CP649, 2014.