

Machine Learning based Multivariate Time Series Forecasting of ISO New England Electricity Demand

September 2023

Abstract

It is of the utmost importance to make accurate predictions of multivariate time series in the context of the dynamic environment of data-driven decision-making. This dissertation makes use of modern algorithms such as Random Forest, CatBoost, XGBoost, and Prophet to analyse multivariate time series data in an effort to decipher the complex dynamics that lie inside those datasets. The research investigates topics such as predicting accuracy, computing complexity, and exploratory data analysis. It covers a wide range of datasets and encompasses a comprehensive set of evaluation criteria.

The study framework evaluates the collective efficacy of the individual algorithms by using Augmented Dickey-Fuller tests (ADF), RMSE, MSE, MAE, and R^2 scores. This extends beyond the evaluation of the individual algorithms themselves. These measures shed light on comparisons and highlight the strengths of the algorithms in terms of capturing temporal links and managing complex interactions. In addition, the computational complexity is investigated, which sheds light on the efficiency of algorithms across datasets. The study provides a full comprehension of data patterns by incorporating Exploratory Data Analysis and visualization. This research not only makes strides forward in the field of multivariate time series forecasting, but it also provides a symphony of insights that harmonize accuracy, efficiency, and interpretability, thereby revealing the route for future data-driven initiatives.

After successfully implementing machine learning algorithms on my based dataset, all of these algorithms have been validated through four more datasets of the very same nature. After evaluating these algorithms, I have found out that XGBoost algorithm and CatBoost algorithm showed good performance but XGBoost has better performance over CatBoost. Prediction accuracy for this model is far better than the rest of the algorithms. Moreover, it is computationally cheaper.

Attestation

I understand the nature of plagiarism, and I am aware of the University's policy on this.

I certify that this dissertation reports original work by me during my university project except for the following.

- The dataset has been obtained from <https://huggingface.co/datasets/tmgondal/GEFCom>
- Concept of figure 1 has been taken from [7].
- Rest of the material taken from sources has been cited in the text and references have been provided in IEEE style as per the university policy.
- As far as coding portion is concerned, EDA portion was learned and then implemented as per my need from this source: <https://seaborn.pydata.org/tutorial.html>
- Code for augmented dicky fuller method learn and implemented from <https://www.machinelearningplus.com/time-series/augmented-dickey-fuller-test/>
- Implementation of Random Forest in Python but not complete code: <https://towardsdatascience.com/random-forest-in-python-24d0893d51c0>
- Implementation of XGBoost in Python but not complete code: <https://towardsdatascience.com/multi-step-time-series-forecasting-with-xgboost-65d6820bec39>
- Implementation of CatBoost in Python not complete code: <https://cienciadedatos.net/documentos/py39-forecasting-time-series-with-skforecast-xgboost-lightgbm-catboost.html>
- Implementation of Prophet in Python not complete code: https://www.youtube.com/watch?v=XZhPO043lqU&ab_channel=AIEngineering

Signature

Date

Acknowledgements

Throughout the course of this research endeavour, my distinguished supervisor has provided me with unshakable leadership, unwavering support, and tremendous mentorship. For all of these things, I would like to extend my most sincere appreciation. It is due in large part to his extensive expertise, intelligent input, and unwavering determination that the trajectory of my dissertation has taken the form that it has.

I would also like to take this opportunity to express my gratitude to the members of the faculty as well as the mentors who have given freely of their time and knowledge in order to improve the quality of my educational experience. Their insightful criticisms and words of encouragement have been very helpful in improving the overall quality of this work and increasing its level of depth.

In addition, I want to express my profound gratitude to my friends and family for the constant support, patience, and understanding they have shown me over the years. Their unfailing faith in my capabilities has been a consistent source of inspiration for me, motivating me to triumph over obstacles and achieve previously unachievable goals.

Without the combined assistance, direction, and ideas provided by these people and entities, this task never would have been completed. It is undeniable that their efforts have rendered an indelible imprint on the finished product of this study endeavour.

Table of Contents

Abstract	i
Attestation.....	ii
Acknowledgements.....	iii
Table of Contents.....	iv
List of Figures.....	v
List of Table.....	i
1 Introduction	1
1.1 Background and Context	2
1.2 Scope and Objectives	3
1.3 Achievements.....	4
1.4 Overview of Dissertation.....	4
2 State-of-The-Art	5
3 Literature Review	6
4 Exploratory Data Analysis.....	11
4.1 Impact of Weather on Electricity Demand.....	19
4.2 Effect of User Behaviour and Routine on Electricity Demand.....	19
5 Machine Learning Based Multivariate Time Series Forecasting	21
5.1 Multivariate time series forecasting with proposed algorithms	21
5.2 Comparison of proposed algorithms.....	23
5.3 Implementation of ML algorithms for Forecasting	24
5.3.1 ADF Results.....	24
5.3.2 Evaluation of these algorithms.....	26
5.3.3 Computational complexity analysis	27
5.3.4 Best Performing Algorithm.....	30
5.4 Recommendation	32
6 Conclusion.....	33
6.1 Summary	33
6.2 Evaluation.....	33
6.3 Future Work	33
References.....	35

List of Figures

Figure 1: Purpose of Visualizations in EDA	11
Figure 2: Histogram for Distribution of Demand.....	12
Figure 3: Line Plot for Electricity Demand Variations over a Day	13
Figure 4: Distribution of Electricity Demand Across Months	13
Figure 5: Overview of Electricity Demand Across Months.....	13
Figure 6: Distribution of Electricity Demand Across Week	14
Figure 7: Distribution of Electricity over the Days of Year.....	15
Figure 8: Relationship of Electricity Demand and Temperature	15
Figure 9: Relationship of Electricity Demand and Humidity	15
Figure 10: Holiday Count in the ME_Zone of ISO New England Dataset.....	16
Figure 11: Total Demand by Type of Holiday.....	16
Figure 12: Demand Distribution of Weekdays and Weekend	17
Figure 13: Average Demand by Hour on Weekdays and Weekend.....	17
Figure 14: Average Demand by Month on Weekdays and Weekend	18
Figure 15: Demand Distribution by Weekend and Holiday	18
Figure 16: Comparison of All Days with Weekend and Weekdays	18
Figure 17: Comparison of ML algorithms performance	29
Figure 18: Time consumed by each algorithm in each dataset.....	30
Figure 19: Multivariate forecasting performance of best performing algorithm.....	31

List of Table

Table 1: Comparison of Related articles.....	8
Table 2: Dataset Overview.....	11
Table 3: Feature based comparison of ML algorithms	24
Table 4: Comparison to ADF test.....	25
Table 5: Performance evaluation of ML algorithms	27
Table 6: Comparison of ML algorithms on the basis of computational complexity	29

1 Introduction

In the area of planning and managing energy, which changes quickly, accurate forecasting is a key chapter of making good decisions. Forecasting is the process of forecasting future values or trends based on past information. This lets organizations plan for changes and be ready for them. When it comes to forecasting how much electricity will be used, energy providers need accurate predictions to make the best use of their resources, make sure they have a steady supply of power, and make smart choices about capacity planning and infrastructure upgrades [1]. Demand for electricity is affected by many things, such as population growth, the state of the economy, the seasons, and technology advances. Because the data come from many different sources, it is hard to understand and identify these demand patterns [2].

In the past, traditional tools for forecasting, such as statistical methods and time series analysis, were used a lot. Multivariate time series data, on the other hand, have complicated connections and dependencies that are hard for them to understand [3]. These traditional methods are based on the ideas of uniformity and stability, which may not be true in the real world. Because of this, their ability to guess may be limited, which could lead to less-than-optimal decisions and resource allocation [4]. To deal with these problems, machine learning algorithms have become a potential way to predict how much electricity will be used. Multivariate time series data has a lot of complicated interactions and dynamics that can be captured by machine learning models. By using advanced algorithms and a lot of computing power, these models can learn from past trends and make accurate predictions based on the relationships they've learned [5].

The goal of this research dissertation is to find out how well different machine learning methods help make electricity demand forecasts for the ME Zone in ISO New England more accurate. The ME Zone has complicated trends of energy use that are affected by things like holidays, weather, humidity, and more. Conventional methods often have trouble forecasting links that are this complicated. So, the main goal of this study is to find out how well machine learning methods, such as the Random Forest, CatBoost, XGBoost, and Prophet algorithms, can help solve this problem. Random Forest is a group learning method that makes guesses by putting together several decision trees [6]. It is strong against overfitting and can handle big datasets with complex interactions. CatBoost is a gradient-boosting algorithm that is well-known for how well it deals with categorical factors [7]. Another popular gradient-boosting method, XGBoost, has a high level of prediction accuracy and is used in many different fields [8]. Prophet is a system for forecasting time series that was made by Facebook. It looks for seasonality, changes in trends, and outliers [9]. Common evaluation variables, such as Root Mean Square Error (RMSE), Mean Absolute Error (MAE), Mean Absolute Percentage Error (MAPE), and the coefficient of determination (R^2 -score), will be used to measure how well these machine learning models work. These measures will give a complete evaluation and explanation of how well the models can predict. The goal of this study is to learn more about the patterns and dynamics of energy demand in the ME Zone, as well as to evaluate how well these models work. By looking at the importance of features and how models explain them, I can find the variables that have the most impact on power demand. This knowledge can help energy providers and power system operators make smart choices about where to put resources, how to plan for capacity, and how to improve infrastructure. This research dissertation also has a thorough review of the literature to give a full picture of studies already done in the field of forecasting energy demand and machine learning-based approaches. By looking at the work of well-known scholars and experts, I hope to find the most up-to-date methods, point out their strengths and weaknesses, and find research gaps that could be filled by this study.

I will use the Python programming language to analyse, model, and evaluate the data for this study. Libraries like Pandas, NumPy, Scikit-learn, and TensorFlow will be used to handle the da-

ta, build and train the machine learning models, and evaluate how well they work. The data for this study will come from ISO New England and will focus on the ME Zone for a certain time period, like 2003–2017. The data set will include information about how much energy was used in the past, as well as variables like weather, holidays, and other things that affect demand. The goal of this research dissertation is to use machine learning algorithms to make energy demand forecasts for the ME Zone in ISO New England more accurate. By figuring out the complicated trends and changes in how much electricity is used, I can make better plans and decisions about energy. Using established metrics to evaluate and understand machine learning models will show how well they can predict how much electricity will be used. Through this study, I hope to make a contribution to the field of forecasting how much electricity will be used, make energy planning more accurate, and give energy providers and power system operators useful information.

1.1 Background and Context

Multivariate time series forecasting plays a vital role in insights generation from any dataset whether it is related to fraud detection, the banking sector, the money market and makers estimate the requirements for the upcoming days, months, or years [10]. Recently, for power systems to operate and be planned effectively, accurate electricity demand forecasting is essential. Energy providers may optimize resource allocation, cut costs, and guarantee the consistency of power supply by correctly forecasting electricity demand [7]. The exploration of machine learning algorithms is motivated by the fact that conventional forecasting techniques i.e., statistical approaches frequently fail to grasp the intricate relationships found in multivariate time series data and they have lower efficiency. In this study, I will compare the performance of machine learning algorithms with the conventional methods and improve the precision through optimal fine tuning of algorithms where needed. In the literature review section, various articles have been explored to provide a few insights into work being carried out by renowned scholars.

Multivariate forecasting has few benefits as listed below:

- **Optimization of Resource Allocation:** Energy companies can optimize resource allocation by foreseeing patterns in electricity demand thanks to accurate projections. This aids in effective load balancing, distribution, and scheduling of generation sources, which lowers costs and guarantees the best use of resources.
- **Reliable Power Supply:** Power system operators can anticipate peak demand periods and potential capacity gaps thanks to accurate forecasting of electricity consumption. Operators may proactively plan for infrastructure upgrades, maintenance, and emergency response with the help of realistic projections, ensuring that customers receive an uninterrupted supply of electricity.
- **Incorporation of Complex Relationships:** The complex correlations between power demand and many variables, like temperature, humidity, and holidays, may be captured by machine learning algorithms. Accurate projections can be made even in the midst of complicated interactions and nonlinear dependencies between variables by taking use of these linkages.
- **Improved Planning and Decision-Making:** Forecasts that are accurate help with long-term planning, investment choices, and policy creation. Based on accurate forecasts of future electricity demand, energy suppliers can decide on capacity development, demand response initiatives, the integration of renewable energy sources, and system stability measures.

By performing multivariate time series project my motive is develop in-depth insights from the obtained results. Which will help me to provide recommendations on the better utilization of forecasting models. Through these recommendations, policy makers can drive comprehensive policies for global and national interests. Moreover, it will provide me an opportunity to explore, visualize and understand the electricity demand-response nexus.

1.2 Scope and Objectives

In the recent era, precise electricity demand forecasting is crucial for efficient energy planning and management. The complicated relationships seen in multivariate time series data are frequently difficult for conventional forecasting tools to capture. The investigation of machine learning algorithms as a potential remedy is driven by this constraint. This research proposal aims to estimate the electricity demand of the ME Zone in ISO New England using a range of machine learning techniques, with the purpose of improving prediction accuracy and offering useful data to energy providers and power system operators.

The core objective of this dissertation is to use machine learning-based algorithms to perform multivariate time series forecasting which will be performed on one of the zones of ISO New England. The patterns of energy used in this area are complicated and difficult to predict using conventional techniques since they depend on a variety of variables, including holidays, weather, humidity, and other variables. In this research work, my motive is to investigate the use of machine learning approaches i.e., Random Forest, CatBoost, XGBoosts and Prophet algorithms, for the successful implementation of multivariate time series forecasting. Furthermore, these algorithms will be evaluated on the basis of RMSE, MAE, MAPE and R^2 -score for the better evaluation and interpretation of the algorithms.

- This study aims to focus on the following research questions:
- What is the impact of weather (humidity and temperature) on the demand of electricity?
- What is the effect of user behaviour and routine on the demand of electricity?
- How much effective are the machine learning algorithms while performing multivariate time series forecasting?
- How the multivariate can help to perform better energy planning?

To answer the above research question following are the core objectives of this dissertation are:

- To conduct a thorough exploratory data analysis to learn more about the traits, trends, and connections found in the ME Zone electricity demand dataset.
- To successfully explain the results and offer insightful context for the patterns of electricity demand, use narrative and visualizations.
- To utilize the Random Forest, XGBoost, CatBoost, and Prophet algorithms to anticipate power consumption using multivariate time series forecasting approaches.
- To use measures like RMSE (Root Mean Square Error), MAPE (Mean Absolute Percentage Error), R^2 Score, MAE (Mean Absolute Error), and computational complexity analysis to assess and compare the performance of the suggested algorithms.

Technologies and Data Sources

Technologies:

- Python programming language for data analysis, modelling, and evaluation.

- Seaborn, Matplotlib, Pandas, NumPy etc. will be used for the comprehensive exploratory data analysis.
- Machine learning libraries such as scikit-learn etc. will be used to evaluate proposed machine learning based algorithms

Data Sources:

- I will use ISO New England electricity demand dataset for the ME Zone for the year 2003-2017.
- ISO New England Dataset is based on eight zones, in this research work I will be using only zone. These datasets are publicly available at GitHub, Kaggle and hugging face. Link to the dataset is

<https://huggingface.co/datasets/tmgondal/GEFCom>

1.3 Achievements

In this dissertation, I have analysed the performance of machine learning algorithm while performing multivariate load forecasting. While doing this, I have selected ISO New England dataset of electricity. Through my findings I intended to relate my stated research questions. I have examined the impact of weather variables and user behaviour on the changing dynamics of the electricity demand. I have found out that extreme weathers result in higher demand and user behaviour is strictly associated with the rise or fall of the electricity demand. While evaluating the performance of machine learning algorithms i.e., random forest, XGBoost, CatBoost and prophet. I have computed RMSE, MAE, MSE and R^2 score along with the computational complexity. I have found out that XGBoost outperforms the rest of the machine leaning algorithm.

1.4 Overview of Dissertation

Chapter 1 is focused on the introduction where I present my research question and research objectives. Meanwhile, chapter 2 presents and idea that why this research problem is the state of the art. Literature review has been discussed in the chapter 3 and chapter 4 comprehensive-ly presents the exploratory data analysis. In this section, we will see the impact of user behaviour and temperature parameters on electricity demand. In chapter 5, evaluation of machine learning algorithms have been carried out on ISO New England i.e., Dataset 1. These algorithms have been validated in this chapter as well. In chapter 6, conclusion, evaluation and future directions have been presented.

2 State-of-The-Art

The accurate prediction of multivariate time series has emerged as an important activity in the shifting environment of data-driven decision-making. This dissertation begins on an insightful examination of the realm, making use of the capabilities of contemporary algorithms such as Random Forest, CatBoost, XGBoost, and Prophet to decipher the complex dynamics that are inherently present in multivariate time series data. This research highlights the numerous aspects of forecasting accuracy, computational complexity, and exploratory data analysis by spanning five separate datasets and being supported by a broad array of evaluation criteria.

This investigation is built on a solid foundation that is comprised of a stringent evaluation framework. This framework goes beyond the confines of individual algorithms and examines the efficacy of the algorithms as a whole. The accuracy of Random Forest, CatBoost, XGBoost, and Prophet are all subjected to painstaking scrutiny in the form of an in-depth evaluation that is guided by Augmented Dickey-Fuller tests (ADF), Root Mean Square Error (RMSE), Mean Squared Error (MSE), Mean Absolute Error (MAE), and the R-squared (R^2) score. The convergence of these important parameters, which serve as hallmarks of predictive precision, reveals a comparison tableau that illustrates the respective strengths of the algorithms in capturing temporal correlations, compensating for variability, and navigating dynamic interactions.

This dissertation investigates the temporal complexities of computational complexity at a time when computing agility is becoming increasingly important to analytical rigor. When the amount of time spent in the training and testing phases of an algorithm is painstakingly measured and analysed, a comprehensive picture of the effectiveness of the algorithm may be painted. This lens sheds light on the applicability of the algorithms to a wide variety of datasets and offers some pragmatic considerations for real-time applications.

This research paradigm goes beyond the confines of numerical constructs in order to shed light on the experiential aspects of data. The exploratory data analysis (EDA), which is supported by a wide variety of visual representations, acts as the foundation for a full understanding. These visual expositions uncover hidden patterns, directing the strategic deployment of algorithms while uncovering the quintessence of multivariate time series dynamics. Infused with a story of trends, periodicities, and anomalies, these visual expositions discover latent patterns.

This dissertation helps move the field forward at an important time because it comes at a time when making informed judgments will increasingly depend on one's ability to anticipate the future. It presents a polyphonic symphony of insights by deciphering the riddle of multivariate time series forecasting through the lens of Random Forest, CatBoost, XGBoost, and Prophet. This symphony, which successfully combines accuracy, efficiency, and interpretability, serves as a guiding light for future efforts that aim to harness the power of data science in order to shed light on the annals of history. Within these pages, algorithms transcend their status as simple tools and transform into the channels that comprehend the data symphony and direct the development of multivariate time series forecasting.

3 Literature Review

Forecasting is a basic idea in many fields, and it has become an important area of study in both the classroom and the workplace. It means making predictions about future values, patterns, or trends based on what has happened in the past [11]. Forecasting tries to give people information and help them make decisions by predicting and planning for changes that will happen in the future. Furthermore, it is essential in numerous fields, including economics, finance, business, supply chain management, weather forecasting, and resource planning [9]. The necessity for organisations to navigate the uncertainties of the future led to the emergence of forecasting as a distinct discipline. As businesses and societies grew more complex, the ability to anticipate and plan for future events became crucial for survival and success [10]. Various scientific centres and universities have begun to investigate various methods and techniques for analysing and predicting the behaviour of time series data.

Reflecting advancements in statistical techniques, computational capacity, and the availability of data, forecasting methodologies have evolved significantly over time. Early forecasting methods significantly relied on time series analysis, which aimed to identify patterns and trends within a single data series [11]. These techniques, such as moving averages and exponential smoothing, offered preliminary forecasting insights [12]. As the discipline advanced, researchers began to investigate more sophisticated methods for addressing the difficulties posed by multivariate time series data. Multivariate forecasting involves contemplating multiple variables concurrently, recognising that interactions and dependencies between variables can have a significant impact on future outcomes [6]. This resulted in the development of advanced techniques such as regression analysis, vector autoregression, and structural equation models, which sought to capture the complex relationships within multivariate data [13]. Artificial intelligence and machine learning have further revolutionised the field of forecasting. These techniques allowed for the analysis of vast quantities of data, the automatic extraction of features, and the modelling of intricate nonlinear relationships. Random Forest [6], XGBoost [8], CatBoost [12] and Prophet [9] became well-known because they could handle multivariate time series forecasting jobs and improve the accuracy of their predictions.

Forecasting study in academia has grown to include more specific application areas. Scholars have tried to make forecasting models more accurate and useful by adding subject knowledge and combining different data sources and characteristics [14]. The evaluation and comparison of different forecasting methods, their success metrics, and the identification of best practises have become important research areas. With the help of electricity demand forecasts, energy providers can make the best use of their resources, guarantee a steady supply of power, and make smart decisions about capacity planning and infrastructure upgrades [7].

Machine learning algorithms have gotten a lot of attention in recent years as a possible way to make energy demand predictions more accurate [15]. In this chapter, I've summed up a few important pieces that looked at how machine learning methods can be used to make predictions for multivariate time series. In [16], the authors give an overview of the most current machine learning methods used to predict how much electricity will be used. The authors discuss about the benefits of AI-based methods, like how they can handle nonlinear relationships and use data about how much electricity is used to find complex trends. The issues with AI-based prediction methods and where they might go in the future are also discussed about.

The authors of [17] focus on how deep learning methods can be used to predict short-term load. They look at different designs for deep learning, such as Convolutional Neural Networks, Recurrent Neural Networks, and Long Short-Term Memory networks. The authors discuss about how deep learning can be used to find temporal dependencies and complex patterns in data about how much electricity is used. Also, they focus on the problems and possible future

directions for using deep learning to improve the accuracy of short-term load forecasts. In this comparative study, the authors look at how well different machine learning algorithms predict energy loads. They compare the accuracy of Random Forest, Support Vector Machines, Gradient Boosting Machines, and Artificial Neural Networks by using real-world data about how much power people use. Root Mean Square Error (RMSE), Mean Absolute Error (MAE), and Mean Absolute Percentage Error (MAPE) are some of the evaluation metrics that the authors look at. The study gives important information about the pros and cons of different algorithms for predicting electricity usage.

Authors in [19] introduced the Facebook-developed Prophet forecasting model, which demonstrates promise in capturing seasonal trends in e-commerce revenue. The model's accuracy is assessed using RMSA and coverage, confirming its suitability for capturing seasonal tendencies. In the realm of cryptocurrency market prediction, [20] propose a method based on the CatBoost model and Bigdata. The study demonstrates the superior performance of the CatBoost model compared to gradient boosting, support vector machines, and linear regression algorithms. To address the randomness and volatility of electrical power loads, [7] presents a complete ensemble empirical mode decomposition with adaptive noise–CatBoost–self-attention mechanism-integrated temporal convolutional network for short-term load forecasting. The method combines time series decomposition, feature selection, and a self-attention mechanism to improve forecasting accuracy.

In the field of traffic flow prediction, [21] proposes an LSTMXGBoost model for urban road short-term traffic flow prediction. The model combines the characteristics of LSTM networks and XGBoost algorithms, resulting in improved prediction accuracy. For cloud GIS services, [8] presents a wavelet transforms-based ARIMA-XGBoost hybrid method to predict layer action response time accurately. The method utilizes wavelet transforms, ARIMA modelling, and XGBoost to enhance the prediction results. [22] introduces a highway tunnel pavement performance prediction approach based on a digital twin and multiple time series stacking (MTSS). The approach combines heterogeneous stacking of XGBoost, artificial neural networks, random forest, ridge regression, and support vector regression to accurately predict pavement performance changes. In the context of multivariate time series forecasting, authors in [23] explored the modelling of evolutionary and multi-scale interactions using graph neural networks. The proposed method utilizes a hierarchical graph structure and dilated convolutions to capture scale-specific correlations, leading to superior performance in forecasting tasks.

For electricity load forecasting, in [24] authors propose an efficient scheme that employs eXtreme Gradient Boosting (XGBoost) to extract features and forecast load. The scheme outperforms other methods in terms of mean average percentage error. Finally, in industrial applications, [25] presents the XGB-GRU model, which combines eXtreme Gradient Boosting (XGBoost) and Gate Recurrent Unit (GRU) for multivariate time series prediction. The model leverages XGBoost's feature extraction capabilities and GRU's timing information extraction to achieve accurate predictions.

In [26], the authors examine the use of oblique random forests for time series forecasting. The proposed method employs a least square classifier instead of the traditional orthogonal decision tree, allowing for capturing the geometrical structure of data samples. The experimental results demonstrate the advantages of the proposed method in terms of efficiency and accuracy. The application of random forest algorithms for regression modelling in the context of forecasting the remaining useful life of complex technical systems is investigated in [27]. The study explores approaches to improving forecasting accuracy through the generation of new features for inclusion in training and test datasets. The paper provides recommendations for developing regression models for predicting the remaining useful life of systems. In the field of time series forecasting, a new algorithm that combines clustering, classification, and forecasting techniques is proposed in [28]. The algorithm aims to group time series values with similar

patterns and build specific forecasting models for each pattern. The flexibility of the algorithm allows for the use of various machine learning techniques. The experimental results show the superiority of the proposed algorithm compared to classical prediction models and recent forecasting methods. [29] focuses on the accurate forecasting of weekly dengue incidence in multiple cities in Brazil. Machine learning models incorporating feature selection, such as LASO and Random Forest regression, are compared with LSTM, a deep recurrent neural network. The results indicate that the LSTM model outperforms other models in predicting future dengue incidence, showcasing its potential for disease control. Urban water demand forecasting is a challenging task due to non-stationarity and non-linearity. In [30], a hybrid forecasting model combining temporal convolution neural network (TCN), discrete wavelet transform (DWT), and random forest (RF) is proposed. The model utilizes RF to rank and select important factors, reduces dimensionality through DWT, and employs TCN for accurate prediction. Experimental results demonstrate the superiority of the proposed model over other benchmark models. The forecasting of tourist flows is examined in [31], where univariate and multivariate time series models are compared. The study explores the presence of cross-correlations among different origin-destination tourist flows and evaluates the forecasting performance of various models. The results reveal that ARIMA exhibits better forecasting performance compared to univariate and multivariate state space modelling.

In the field of multivariate time series analysis, [32] presents a neural network approach to model flour prices in different cities. The results show remarkable success in training the networks to learn and accurately predict the price curve for each city, highlighting the effectiveness of the neural network approach. [33] proposes a novel temporal attention encoder-decoder model for multivariate time series forecasting. The model utilizes bi-directional LSTM layers with a temporal attention mechanism for adaptive learning of long-term dependencies and hidden correlation features in multivariate temporal data. Experimental results demonstrate the superior forecasting performance of the proposed model compared to baseline methods. In [34], a graph neural network framework is proposed specifically for multivariate time series data. The framework automatically extracts uni-directed relations among variables through a graph learning module and captures spatial and temporal dependencies within the time series. Experimental results show the superiority of the proposed model over baseline methods on benchmark datasets. The use of a Multivariate Temporal Convolution Network (M-TCN) model for multivariate time series prediction is explored in [35]. The M-TCN model is designed to handle non-periodic datasets and utilizes deep convolutional neural networks and residual blocks for accurate prediction. Experimental results demonstrate significant improvement in prediction accuracy compared to other competitive algorithms. [36] introduces a novel attention mechanism for multivariate time series forecasting. The proposed model utilizes filters to extract time-invariant temporal patterns and combines them with an attention mechanism for relevant time series selection. The model achieves state-of-the-art performance in various real-world tasks. The article [37] focuses on energy consumption prediction in manufacturing companies using multivariate time series models. The authors compare the performance of the Prophet and LSTM algorithms for predicting energy consumption of electricity, water, and diesel fuel. The results indicate that Prophet performs better in predicting the energy consumption of the studied variables. In Table 1, a comparison of articles being presented in the literature review has been provided.

Table 1: Comparison of Related articles

Study	Parameters Evaluated	Forecasting Type	Contribution
[19]	Seasonal trends, Facebook Prophet applicability	E-commerce revenue forecasting	Identified seasonal trends, suitability of Prophet tool

[20]	Data preprocessing, feature selection, CatBoost model	Cryptocurrency market prediction	Improved prediction accuracy, outperformed other algorithms
[21]	Digital twin, multiple time series stacking	Pavement performance prediction	Accurate and timely prediction, consideration of temporal, spatial, and exogenous dependencies
[22]	CEEMDAN-CatBoost-SATCN method, time series decomposition, feature selection	Short-term load forecasting	Improved forecasting accuracy, positive effect of decomposition and feature selection
[23]	LSTMXGBoost model, LSTM-XGBoost hybrid	Urban road traffic flow prediction	Improved prediction accuracy, efficient traffic guidance
[24]	XGBoost, smart grid load	XGBoost for efficient smart grid load prediction	Efficient use of smart grid, improved prediction accuracy
[25]	XGBoost, GRU, temperature prediction	XGB-GRU model for multivariate time series prediction	Improved prediction accuracy, feature extraction capabilities
[26]	Least square classifier, orthogonal decision tree	Oblique random forests for time series forecasting	Improved efficiency and accuracy, capture of geometrical structure
[27]	Regression modelling, feature generation	Random forest regression for remaining useful life prediction	Improved prediction accuracy, recommendations for regression models
[28]	Clustering, classification, forecasting	Clustering, classification, and forecasting algorithm	Grouping of similar patterns, specific forecasting models, flexibility
[29]	Feature selection, LASSO, Random Forest, LSTM	Dengue incidence forecasting	Accurate incidence prediction, utilization of feature selection
[30]	TCN, DWT, RF, feature ranking	Hybrid model for urban water demand forecasting	Improved forecasting accuracy, feature ranking for efficiency
[31]	Cross-correlations, ARI-MA, state space modelling	Univariate vs. multivariate tourist flow forecasting	Comparison of forecasting performance, absence of cross-correlation structure

[32]	Feedforward connection-ist networks	Neural network approach for flour price forecasting	Successful training and accurate price predictions
[33]	Temporal attention, LSTM	Temporal attention encoder-decoder model	Superior forecasting performance, adaptive learning of dependencies
[34]	Graph learning, graph convolution, temporal convolution	Graph neural network framework for multivariate time series	Automatic extraction of relations, capture of spatial and temporal dependencies
[35]	Temporal Convolution Network, dimensionality reduction	M-TCN model for non-periodic multivariate time series	Improved forecasting accuracy, dimensionality reduction for efficiency
[36]	Filters, attention mechanism, frequency domain	Attention mechanism for multivariate time series forecasting	Extraction of temporal patterns, relevant time series selection
[37]	Prophet, LSTM, energy consumption	Energy consumption prediction in manufacturing companies	Comparison of prediction accuracy, selection of suitable algorithm
[38]	Wavelet transforms-based ARIMA-XGBoost method	Cloud GIS services response time prediction	More accurate response time prediction, integration of frequency domain information
[39]	Hierarchical graph structure, evolving correlations, multi-scale interactions	Graph neural networks for multivariate time series forecasting	Modelling of evolutionary and multi-scale interactions, improved forecasting performance

From this literature review it can be analysed that time series forecasting plays a vital role and the machine learning algorithm have been a ground-breaking benchmark for the prediction models. In this regard in next chapters, these models will be evaluated and implemented on the electricity dataset.

4 Exploratory Data Analysis

In this chapter, complete exploration of the dataset has been achieved. Electricity demand dataset of ME Zone of ISO New England will be used to evaluate the machine learning algorithms for the forecasting. In this regards, as stated in the research objectives one and two, a very detailed data exploratory analysis has been carried out to identify the insightful context of electricity through various patterns, trends and traits of the dataset. This dataset has 124171 rows and 11 variables. The details and datatypes of these variables has been listed in the Table 2.

Table 2: Dataset Overview

Sr. No.	Variable Name	Description	Data Type
1	demand	Demand of electricity	Float 64
2	drybulb	Temperature	Float 64
3	dewpnt	Humidity	Float 64
4	date	Data on which demand was recorded	Object
5	month	Month on which demand was recorded	Object
6	hour	Hour on which demand was recorded	Int64
7	day_of_week	Day_of_week on which demand was recorded	Object
8	day_of_year	Day_of_year on which demand was recorded	Int64
9	weekend	Either it was weekend or not	Bool
10	holiday_name	If holiday, what is the name of holiday	Object
11	holiday	Either it was holiday or not	Bool

In this chapter, gaining insights, finding patterns, seeing connections, and summarizing the key features of the data are the goals of exploratory data analysis (EDA) using visualizations in a dataset. Visualizations are essential to EDA because they provide the data a visual form, making it simpler to comprehend and analyses. The main objectives of using EDA have been presented in Figure 1.

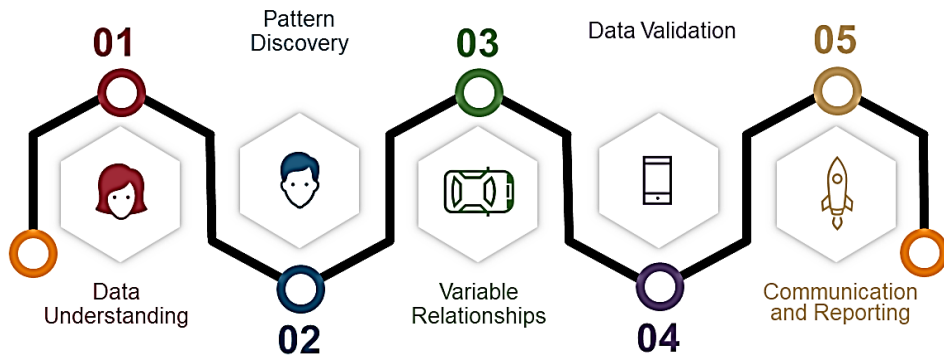


Figure 1: Purpose of Visualizations in EDA

I can can spot discrepancies, mistakes, or anomalies by visualizing the data, which may call for additional research or data cleaning. These offer stakeholders an efficient means of receiving findings and ideas. They facilitate decision-making processes by making it simpler to communicate complex information in a visually appealing and accessible manner.

In order to comprehensively analyses the dataset of ME Zone, all the variables have been explored through different aspects. These provide a unique and different aspects of the visualization in EDA. Through these visualizations, insightful context for the patterns of electricity demand have been provided. All of these visualizations have been obtained by using seaborn and matplotlib. Systems specifications include Core i5, 6th Generation with 32 GB

RAM and 256GB SDD ROM. Mainly, there are three most important parameters available in the dataset i.e., demand, drybulb, and dewpnt, basically these variables reflect the electricity demand, temperature, and humidity. The complete overview of all of these variables has already been discussed in the Table 1. In Figure 2, the histogram plot shows how the dataset's 'demand' variable is distributed. The y-axis displays the frequency of occurrences inside each bin, and the x-axis displays the range of demand values broken up into bins (in this case, 20 bins). The histogram gives us information about the 'demand' variable's shape, central tendency, and variability, enabling us to see trends, spot outliers, and evaluate the distribution's general properties. It can be anticipated from Figure 1 that when the electricity demand was ~1400 MW it has the highest frequency. From the scale it can also be visualized that lowest demand ~800 MW occurred only very few times. Similarly highest demand ~2000 MW also has least occurrences. The lines visualized on the bins is showing Kernel Density Estimation (kde), which shows the smooth representation of distribution of electricity demand.

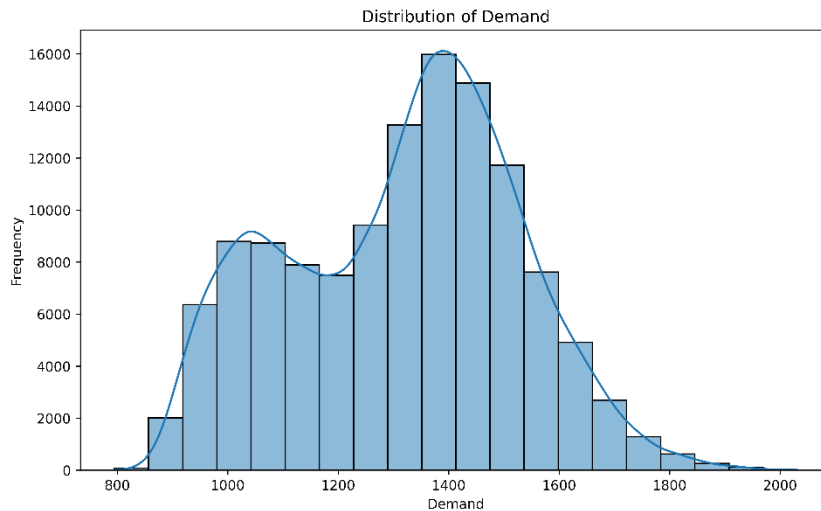


Figure 2: Histogram for Distribution of Demand

The 'ME_zone.csv' dataset's 'demand' fluctuates throughout the day at different times, and the line plot provided in Figure 3 is meant to show this variance. For forecasting and planning purposes, I can make educated decisions by visualizing the hourly change in demand. It can aid in resource allocation, task scheduling, or operational optimization based on expected demand patterns during certain hours. Here, it can be observed that that hours ranging 1-24, where 12 am is indicated. Figure 3, shows that there is rise in the electricity consumption from ~7 AM and it has two spikes. First spike arises almost near ~8 AM – ~2 PM and second spike occurs in the evening time. From this line plot a trend of consumer behavior can be estimated. In Figure 4, the ME_zone.csv dataset's "demand" distribution throughout various months is visualized in the box plot you provided. The box plot illustrates how "demand" is distributed throughout various months, giving information about the central tendency, spread, and variability of each month's demand. In this regard, the middle 50% of the data is represented by the box, which is the interquartile range (IQR). The median of the 'demand' numbers for each month is shown by the horizontal line inside the box. The data range, excluding any outliers, is indicated by the whiskers that extend from the box. Beyond the whiskers, the outliers are displayed as separate points. The box plot makes it simple to compare the demand distribution between months. It aids in locating any seasonal patterns, variations in demand, or outliers that might be present throughout particular months. Additionally, it can shed light on the general distributional traits of the 'demand' for every month. Similarly, it can be observed that both March and April has few outliers or anomalies present in them. Which indicates that, before performing forecasting algorithms, outliers handling should be performed for these both months. Moreover, it can be observed that in July electricity demand was higher followed by the august.

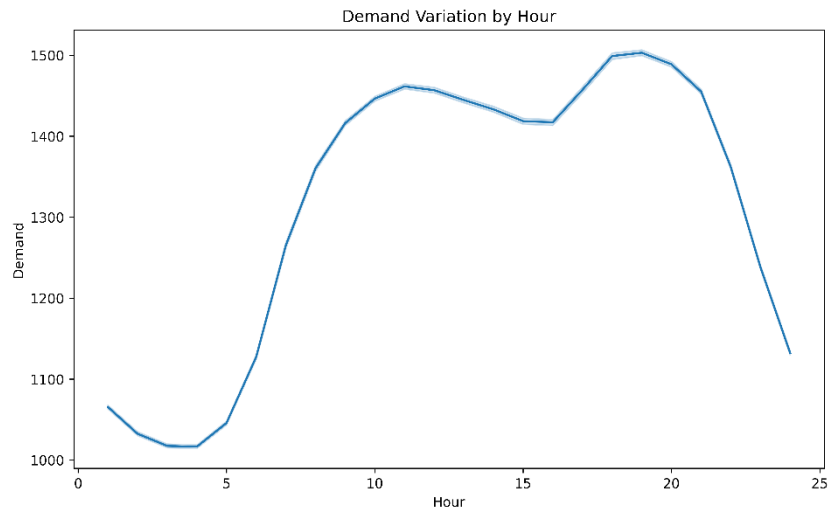


Figure 3: Line Plot for Electricity Demand Variations over a Day

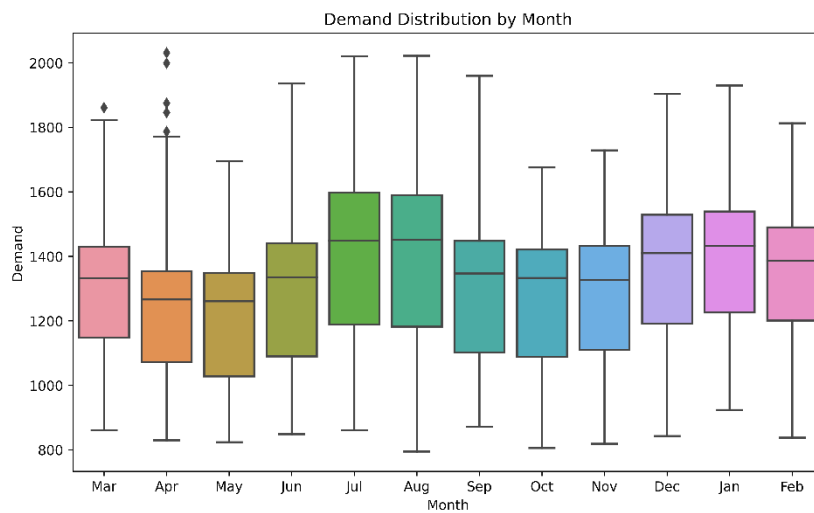


Figure 4: Distribution of Electricity Demand Across Months

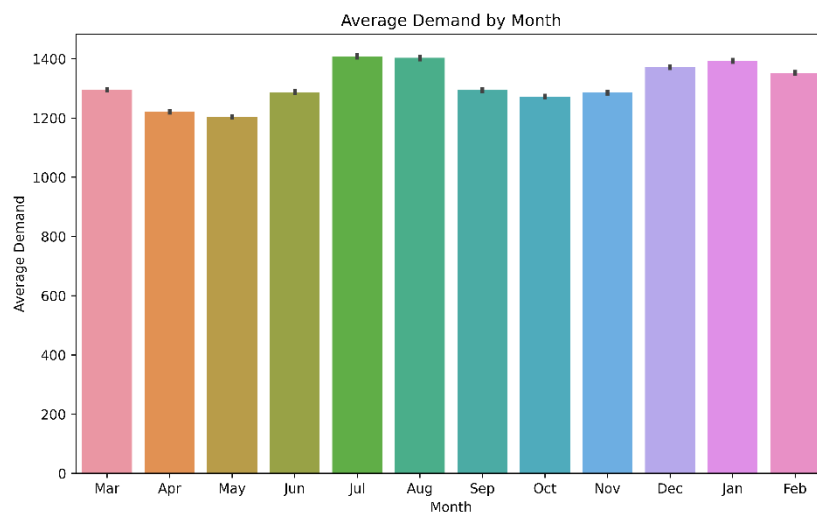


Figure 5: Overview of Electricity Demand Across Months

Similarly, April and May have low electricity demand recorded. This can also be observed from the bar chart presented in Figure 5. This bar chart is providing the average monthly demand over the years span of dataset. This indicates the highest monthly demand as estimated

through the Box Plot provided in Figure 4. In the above figures, hourly and monthly overview of electricity demand from 2003-2017 has been observed. In figure 6, weekly demand for electricity through violin plot has been presented. Here, the demand distribution for each day of the week is depicted clearly and in-depth by the violin plot. The density or frequency of data points at various levels of demand is represented by the width of the violin. A larger portion suggests a higher data point density. Moreover, the median demand value for each day of the week is represented by the white dot inside each violin. The IQR which encompasses the middle 50% of the data, is represented by the thick black line inside each violin. Similarly, Each violin's narrow black lines, or "whiskers," which exclude any outliers or extremely high or low numbers, show the range of the data. From the figure it can be observed that the width of violin is almost same through the week it is due to the equal data points i.e., 124171, for all variables of the dataset. Moreover, IQR is different for all days. Meanwhile, the mean is being indicated by the white dot inside the violin. Which shows that on Thursday, the average recorded electricity demand was higher.

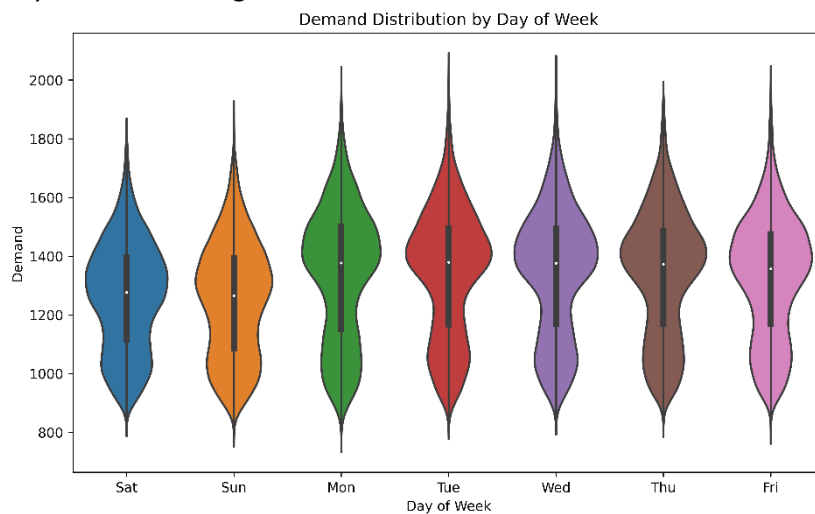


Figure 6: Distribution of Electricity Demand Across Week

An average demand over the number of days of year has been presented in figure 7. From the figure it can be estimated that the middle days of year has highest demand of electricity which more than ~1450 MW. This yearly trend can be better understood while comparing it with figure 4 and 5. The first day of year starts from 1st January. From figure 5, it can be observed that electricity demand is higher in January is higher than February and March. Then there is decrease in electricity demand. Meanwhile in July and August it has a peak value. Similar trend can also be visualized in the figure 7 as well. As stated in Table 1, there are more than one variable present in the dataset. To visually explore the dataset, it is required to observe the relationship among all variable. In this regard, both figures 8 and 9 provide a scatter plot to observe the relationship of electricity demand with temperature and humidity. From figure 8, it can be observed that electricity demand is higher when the temperature is low and high since most of the heating and cooling systems are electrically operated. Here, temperature is in Fahrenheit, when the rises above the ~65 F then there is a rise in the demand of electricity. Similarly, when temperature is below zero the demand of electricity has increased as well. The amount of moisture or water vapor in the air is referred to as humidity. Relative humidity is a common way to measure it since it shows how much moisture is present in the air as a percentage of the total quantity that the air can store at a particular temperature. A higher relative humidity means that more moisture is being held in the air. In figure 9, the relationship between electricity demand and humidity is shown. Apart from temperature and humidity, another parameters present in the dataset is the information about holidays. There are several holidays which have been labelled in the dataset.

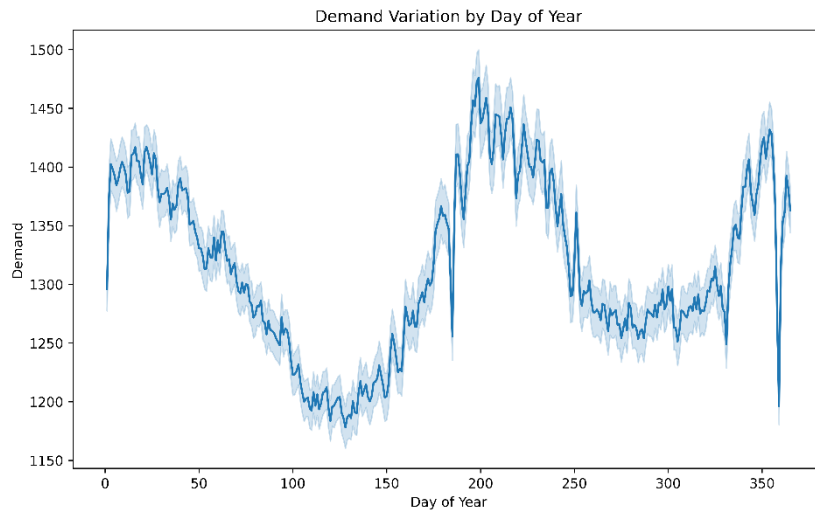


Figure 7: Distribution of Electricity over the Days of Year

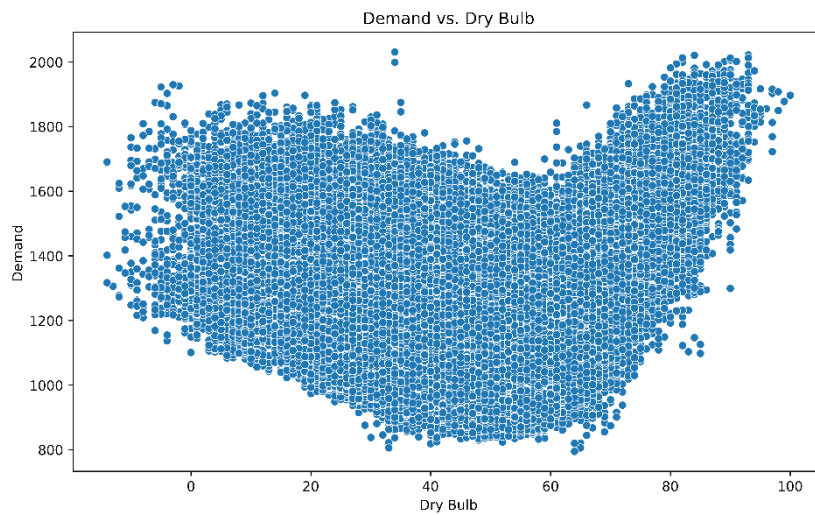


Figure 8: Relationship of Electricity Demand and Temperature

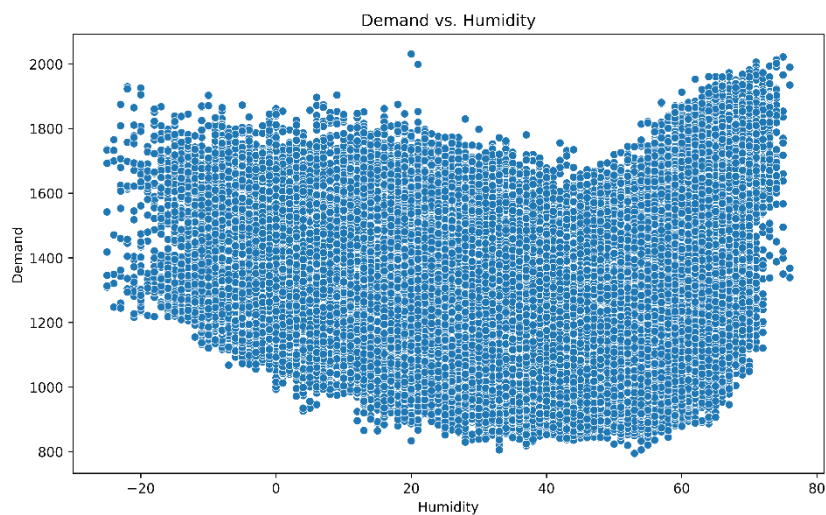


Figure 9: Relationship of Electricity Demand and Humidity

All of those holidays have been visualized here to understand the trends of electricity demands. In this regard figure 10 shows the holidays count. Since the holiday variable has the bool data type, which has the responses as True and False. A bar char has been used to show

that how many entries have the holiday labeled as True or False. From figure 10 shows that 121220 entries have no holiday, and 2953 entries are associated with the holiday.

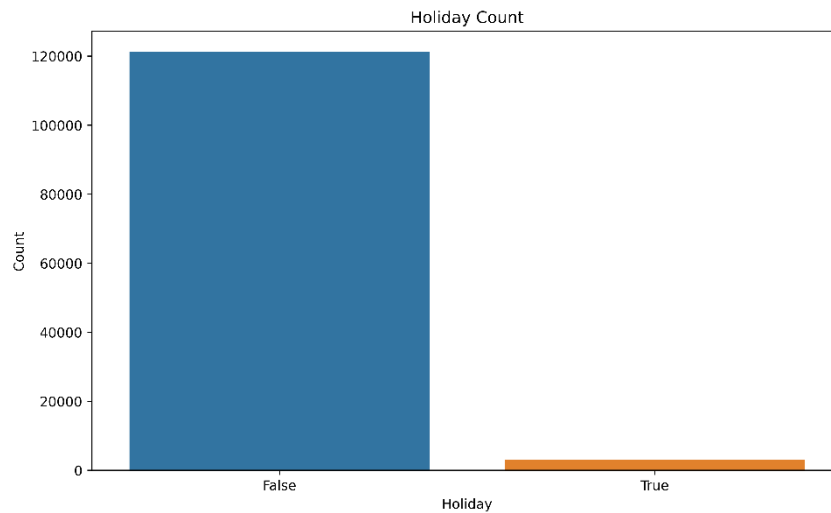


Figure 10: Holiday Count in the ME_Zone of ISO New England Dataset

In the above figure 10, it has been visualized that there is total 2953 holidays present in the dataset. Here, figure 11 provides a insightful trends for the type of holidays and the total average electricity consumption on that certain day. From figure, it can be observed that President's day and Martin Luther King Jr. Day has least electricity demand. Meanwhile, New Year's Day has highest recorded average electricity demand over 400000 MW in the span of 2003-2017. It can be observed that there are thirteen types of holidays throughout the year. In this figure average demand of every holiday has been presented.

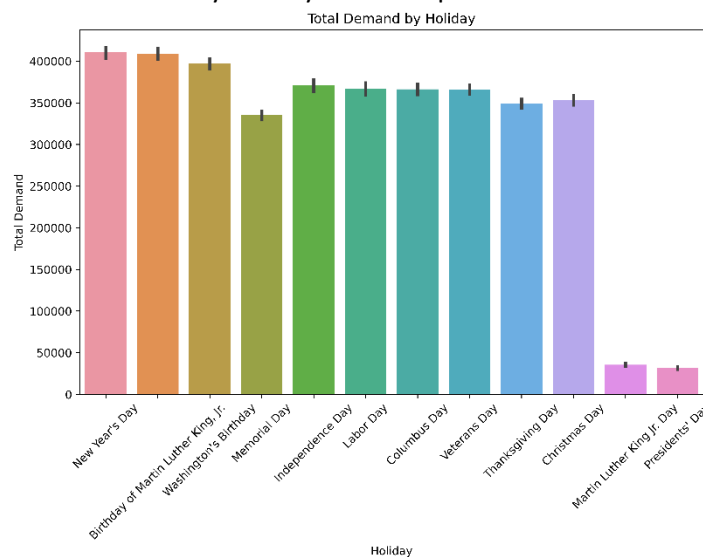


Figure 11: Total Demand by Type of Holiday

Weekly distribution of electricity has already been visualized in the in figure 6. Meanwhile figure 12 provides the boxplot for the weekly electricity demand by separating weekend and working days. From this another insight about the electricity demand can be obtained that on weekends a lower amount of electricity demand has been recorded. The box plot Sunday, Tuesday, Wednesday, and Friday indicates few outliers. Moreover, on weekend there is a decrease in the used of electricity and it has lower mean as well. This information is further elaborated through figure 13. An hourly breakdown for electricity demand on weekdays and weekends has been given in Figure 13. Here, it can be observed that daily demand on weekends is lower than the weekdays.

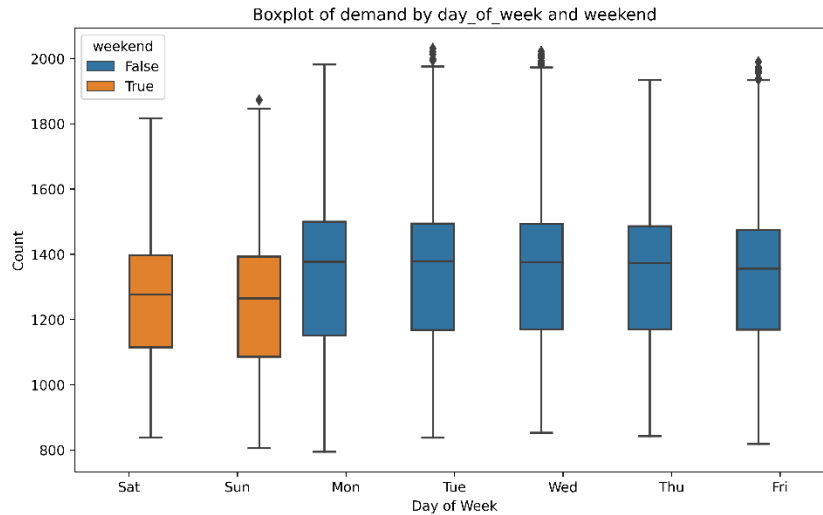


Figure 12: Demand Distribution of Weekdays and Weekend

Moreover, the demand for electricity is higher in the working hours i.e., 8 am to 8 pm. A very similar trend can also be observed in figure 14 with respect to the average demand by month on weekends and weekdays.

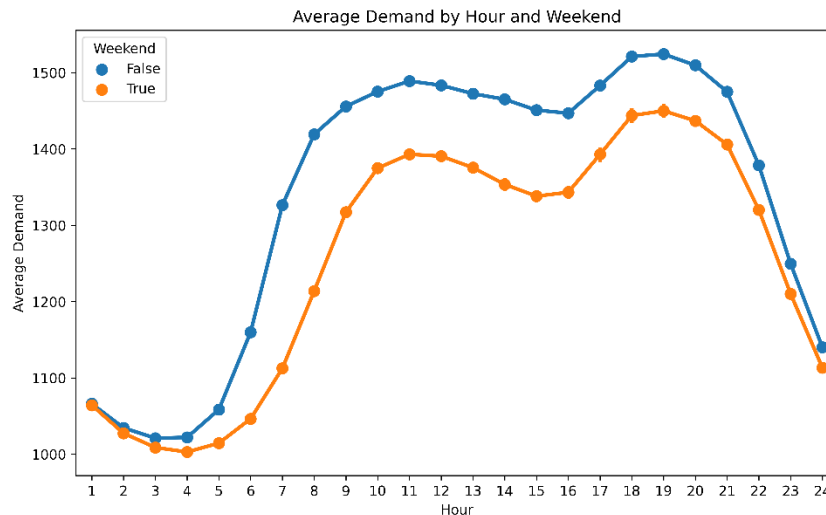


Figure 13: Average Demand by Hour on Weekdays and Weekend

From above results presented in the multiple figures, it can be concluded that the electricity demand is lower whenever there is weekend or holiday. These two parameters have been compared with the rest of the days. Now, in figure 15 a comparison of electricity demand has been carried out for weekend and holidays. Two violin plots for weekend have been presented. It can be observed that when there is no weekend (False) and no holiday (False) then the demand distribution is higher but when there is no weekend (False) but there is holiday (True) then the demand distribution is relatively low. Similarly, for opposite case where Weekend = True and Holiday = False then there is slight difference in demand distribution. In contrast, when Weekend = True and Holiday = True, electricity demand is lowest. In this case lowest electricity demand has been visualized. Throughout these visualizations very insightful data exploratory analysis has been performed which provides insightful context for the pattern of electricity demand in various scenarios. A summary of these findings has been provided in Figure 16. In this figure a histogram for electricity demand has been presented. It can be observed that Weekday have higher electricity demand than weekends and overall energy demand is higher than these.

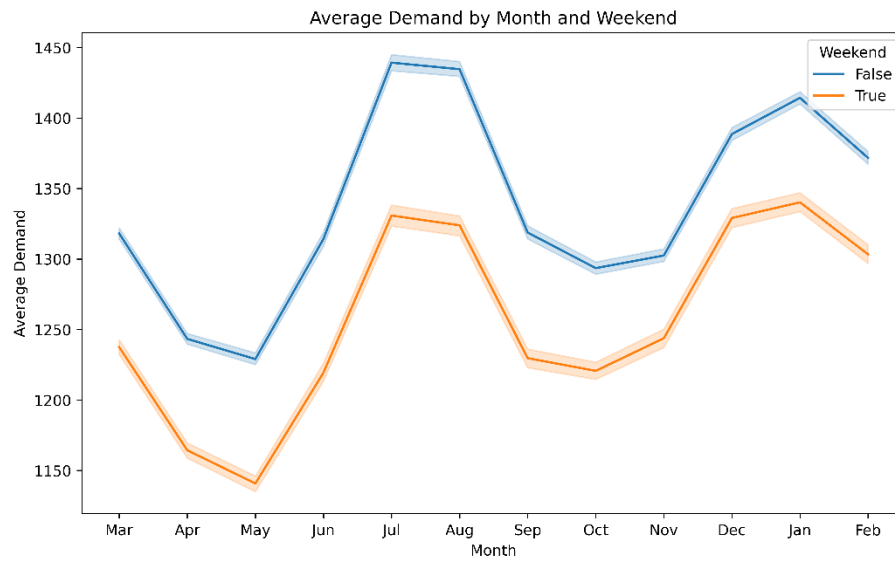


Figure 14: Average Demand by Month on Weekdays and Weekend

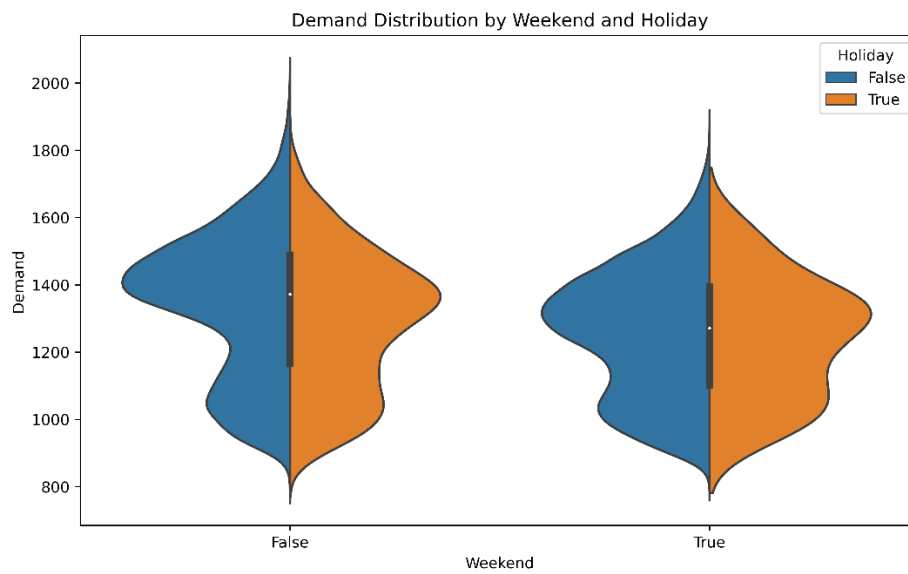


Figure 15: Demand Distribution by Weekend and Holiday

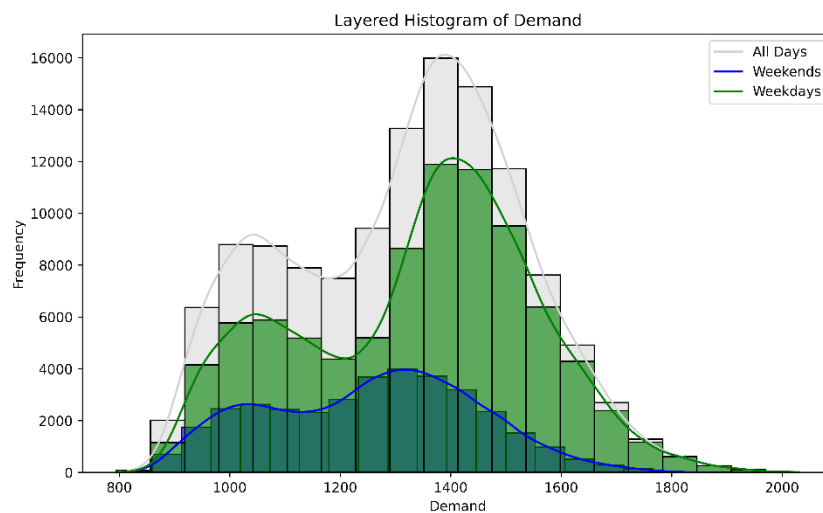


Figure 16: Comparison of All Days with Weekend and Weekdays

Exploratory Data Analysis (EDA) using visualizations on the ME Zone dataset has opened a window into the complex relationship between electricity demand and various influencing factors, notably weather conditions (humidity and temperature) and user behaviour. The analysis has unearthed patterns, trends, and connections that ignite a compelling debate on the profound impact of weather and user behaviour on the dynamics of electricity demand.

4.1 Impact of Weather on Electricity Demand

The indisputable impact of weather on electricity consumption is one of the most important and well-supported conclusions coming out of the EDA. The visuals beautifully demonstrate how temperature and humidity are key factors influencing how much electricity is used. The scatter plots and line graphs clearly illustrate the relationship between temperature and demand, demonstrating how higher temperatures cause spikes in power demand, which are mostly due to the activation of cooling devices like air conditioning during warmer seasons. This tangible link between climate change and energy use emphasizes the need for flexible energy management systems that can adapt to changing weather conditions. In contrast, the relationship between lower temperatures and power demand shows the importance of electric heating systems, highlighting the close connection between user comfort and energy use.

It is equally compelling how humidity and electricity consumption are related. The graphics reveal a strong link between user preferences for controlling indoor climate and power use. The data unmistakably shows an increase in electricity use as humidity levels rise, highlighting the symbiotic relationship between human behaviour and environmental factors. This underlines both the behavioural aspect of energy usage and the opportunity to maximize energy efficiency by synchronizing demand with climate variations. To ensure resilience in the face of weather-related demand surges, a full understanding of these weather-driven patterns becomes essential for reliable energy demand prediction and efficient grid management.

4.2 Effect of User Behaviour and Routine on Electricity Demand

The EDA looks beyond weather-related variables and explores the fascinating world of user behaviour and everyday routines as key determinants of power consumption. Weekday and weekend demand breakdowns by hour show a striking contrast that captures the ebb and flow of societal activity. During regular business hours on weekdays, electricity demand spikes, indicating the industrial and commercial pulse that propels energy consumption. The power system responds to the flurry of activity in businesses, offices, and factories by increasing supplies. On the other hand, demand decreases noticeably over the weekend as households and companies follow a different schedule. This phenomenon emphasizes how important it is for human activity and scheduling to determine patterns of electricity consumption.

Holidays are also included in the analysis, providing a nuanced perspective on how special occasions and energy demand interact. The data reveals variations in demand over holidays, which are explained by changes in user behaviour and recurring patterns. A fascinating view into the dynamics of electricity demand on these special days is provided by the visualization of holidays, particularly those connected with significant events. Electricity consumption ebbs and flows in response to people's participation in a variety of activities and gatherings and the overall pulse of society. This emphasizes how important it is to take cultural, social, and temporal elements into consideration when predicting energy demand and enhancing energy management systems.

In essence, the ME Zone dataset exploratory data analysis has revealed the complex interrelationship between weather, user behaviour, and power demand. The graphics are a potent demonstration of the complex interplay between weather changes, daily activities, and energy

consumption. These discoveries have broad ramifications for researchers, utilities, and those responsible for developing energy policies. In order to orchestrate effective and environmentally responsible energy systems, a thorough understanding of these complex influences is essential as the globe struggles to meet rising energy demands and the need for sustainable consumption. Combining data-driven insights with cutting-edge analytical methods offers a way to create an energy landscape that is not just adaptable to climatic whims and human cycles, but also set up for a future of fair and resilient energy usage.

Exploratory data analysis was used to discover more about the characteristics, trends, and relationships, which helped to provide illuminating context for the electricity patterns discovered in the dataset. This shows that some variables contain outliers that need to be dealt with before machine learning-based techniques are implemented. Additionally, a comprehensive forecasting of the parameters can be used, which is covered in the following chapters.

5 Machine Learning Based Multivariate Time Series Forecasting

Based on machine learning Multivariate Time Series Forecasting is a cutting-edge paradigm that uses the power of cutting-edge machine learning algorithms to predict future values in complex time-dependent datasets. Contrary to traditional approaches, which are characterized by simple assumptions or linear approximations, this method excels in its ability to recognize the complex interactions of many different variables, hence revealing underlying patterns and dynamics woven into the data fabric. Its relevance is particularly seen in fields like finance, healthcare, energy, and economics, where the accuracy of forecasts is crucial for making wise decisions.

Data manifestation includes a variety of variables recorded at successive time junctures inside the multivariate time series domain. Machine Learning-based Multivariate Time Series Forecasting's core strength is in its capacity to unravel the complex web of connections that binds these variables together while overcoming the limitations of conventional univariate approaches. This method abstracts complex temporal patterns from historical data using a variety of sophisticated tools, including neural networks, support vector machines, random forests, and deep learning architectures. As a result, it is capable of making reliable predictions about future values. These models go through training processes supported by various datasets, enabling them to discover hidden correlations, identify seasonality patterns, decipher trends, and even accommodate anomalies that collectively contribute to an increased level of forecast accuracy.

The importance of machine learning-based multivariate time series forecasting assumes salience in its propensity to produce accurate predictions and foster usable insight in the crucible of complex real-world settings. Traditional time series approaches, which were often developed for less complex data, frequently fall short when faced with the complexities present in modern datasets rich with multifarious inter-variable dynamics. In contrast, this innovative strategy offers a variety of observable benefits:

Its ability to simultaneously include a variety of variables creates an environment where complex linkages and dependencies, the cornerstones of the target variable's behavior, are completely enclosed. The result is prognostications that surpass the scope restricted by conventional univariate approaches and exhibit an unparalleled level of accuracy.

Additionally, Machine Learning-based Multivariate Time Series Forecasting demonstrates skill in navigating the complex convolutions that are inherent to a variety of domains. It skillfully handles the ups and downs of dynamic transitions, nonlinear trajectories, and aberrant data occurrences, qualities essential for accuracy in prediction across a variety of applications.

Last but not least, this strategy skilfully accelerates resource allocation by the accuracy of its predictions. With the support of precise forecasts, organizations are given the freedom to allocate resources efficiently, reduce waste, and improve operational channels, highlighting cost-effectiveness and the operational plexus.

5.1 Multivariate time series forecasting with proposed algorithms

The goal of multivariate time series forecasting, which is crucial to many different businesses, is to predict future values of numerous connected variables throughout time. This hard task is characterized by complex dependencies, dynamic interactions, and nonlinearity, which calls for the use of cutting-edge algorithms capable of identifying and extrapolating complicated patterns. The combination of Random Forest, CatBoost, XGBoost, and the Prophet algorithm in

this situation emerges as a strong armory, each providing specific qualities that enhance the prediction accuracy and analytical understanding in multivariate time series forecasting.

A key development in the field of ensemble learning, the Random Forest method exhibits its astounding effectiveness in multivariate time series forecasting. This algorithm, which has its roots in decision trees, makes use of the power of aggregation by combining many decision trees to produce a solid and adaptable model. Because it can handle a variety of data formats, including category and numerical variables, without the need for intensive pre-processing, it is extremely versatile.

With the help of a combination of many trees, Random Forest excels at capturing complicated relationships between factors and at reducing overfitting. It skilfully handles cluttered data, allows for missing values, and pinpoints variable relevance, providing a thorough comprehension of the underlying dynamics. Additionally, due to its scalability and parallel processing capabilities, it is appropriate for large-scale multivariate time series datasets.

The capacity of CatBoost, a relatively new addition to the family of gradient boosting algorithms, to naturally handle categorical information has garnered much attention. In multivariate time series forecasting, where categorical variables are common, this trait resonates particularly well. By enhancing the learning process, CatBoost uses a cutting-edge method called ordered boosting to minimize the impacts of overfitting. CatBoost excels in the context of multivariate time series forecasting by providing robustness against anomalies and missing data. Its built-in categorical variable handling eliminates the need for feature engineering or one-hot encoding, speeding the pre-processing stage. CatBoost also includes a dynamic learning rate and effective GPU use, which speeds up model deployment and training. Together, these characteristics strengthen its capacity for prediction and make it easier for it to be incorporated into complex forecasting pipelines.

The term "eXtreme Gradient Boosting," sometimes known as "XGBoost," is a pillar of support for multivariate time series forecasting. Gradient boosting techniques are used by the renowned for their scalability, speed, and adaptability XGBoost to create a wide range of decision trees. The model is able to represent complex temporal dependencies and interactions found in multivariate time series data thanks to this ensemble technique. By reducing the tendency for overfitting, XGBoost's regularization mechanisms—such as feature selection and tree pruning—make it robust in the face of challenging datasets. Additionally, by incorporating a weighted quantile sketch, it makes it possible to handle missing values skilfully and reduce the workload associated with data pre-processing. While XGBoost's interpretable feature importance metrics offer insightful information about variable contributions, its ability to handle a number of data types—including categorical, numerical, and text—increases its application in multi-variate time series forecasting.

A unique viewpoint on multivariate time series forecasting is provided by the Prophet algorithm, which was created by Facebook's Core Data Science team. Prophet is uniquely suited for applications where domain knowledge and contextual awareness are essential because, unlike conventional machine learning algorithms, it integrates time-based seasonality, holiday effects, and trend components with a flexible additive model. Prophet's focus on usability and interpretability is in line with how it is used in multivariate time series forecasting. It allows for unevenly spaced observations, manages missing data, and effectively catches abrupt changes or abnormalities in the data. With the ability to customize its parameters, domain specialists can include their expertise to increase the forecast accuracy and promote well-informed decision-making.

For multivariate time series forecasting, the Prophet algorithm, combined with Random Forest, CatBoost, XGBoost, and other techniques, produces a robust toolbox. Combining the advantages of these algorithms provides an all-encompassing framework to solve the difficult

issues of connected variables, dynamic interactions, anomalies, and nonlinearity present in time series data. The combination of these algorithms enables improved interpretability, adaptability, and forecast accuracy across numerous industries, including banking, healthcare, energy, and economics. Now lets move on the compare these tables based on the attributes.

5.2 Comparison of proposed algorithms

Four well-known machine learning algorithms are thoroughly compared in the Table 3. Random Forest, CatBoost, XGBoost, and Prophet. Each algorithm is examined based on many facets of its architecture and functionality, providing insightful information about their skills and traits. This analysis helps to comprehend each algorithm's advantages, distinctiveness, and applicability for various jobs.

An ensemble of decision trees is built by the Random Forest method, which is characterized as an ensemble learning algorithm with bagging. It distinguishes out for having a straightforward architecture that generates predictions using a variety of decision trees. Notably, it relies on recursive feature selection using information gain or Gini impurity and does not necessitate any specific processing for categorical features. Tree building can be efficiently parallelized, and feature relevance scores are provided. Another benefit is that it can manage missing data, which adds to its resilience.

Moreover, CatBoost, in contrast, functions as an ensemble learning algorithm with boosting and distinguishes itself by automatically processing category features and using a unique tree construction method. It has a mechanism for ordered boosting that carefully arranges category feature combinations to speed up training. CatBoost also provides L1 and L2 regularization to manage model complexity and avoid overfitting. Performance is further optimized by its support for GPU acceleration and multicore training. CatBoost offers a complete solution by managing missing data, much as Random Forest.

On the other hand, another ensemble learning approach using boosting called XGBoost is unique in that it optimizes tree construction using gradients rather than boosting. To ensure robustness and minimize overfitting, it incorporates L1 and L2 regularization. Tree building can be done in parallel with XGBoost, which increases training effectiveness. It offers readable feature significance rankings that help with model comprehension. The algorithm's ability to handle missing data effectively improves its suitability for use in practical situations.

Similarly, prophet, a specialized time series forecasting method, addresses a distinct domain. Its architecture, which reflects its emphasis on recording temporal trends, includes trend, seasonality, and holiday components. Prophet enables automated changepoint detection to spot changes in time series behaviour without using parallelism or tree construction. This feature fits with its goal of preserving important temporal occurrences. Even though it lacks feature importance assessment, regularization, and categorical handling, these elements might not be as important in its particular context.

The table 4 concludes with a concise comparison of Random Forest, CatBoost, XGBoost, and Prophet. It draws attention to their design, handling of categorical data, tree construction, regularization, parallelism, feature importance, and treatment of missing data. Understanding these characteristics is essential for choosing the best algorithm for a given task. A diverse set of tools for various machine learning and prediction challenges is made possible by the simplicity and robustness of Random Forest, the automated categorical handling and ordered boosting of CatBoost, the optimization and feature importance of XGBoost, and the specialized focus on time series forecasting of Prophet.

Table 3: Feature based comparison of ML algorithms

Algorithm	Random Forest	CatBoost	XGBoost	Prophet
Type	Ensemble (Bagging)	Ensemble (Boosting)	Ensemble (Boosting)	Time Series Forecasting
Architecture	Collection of decision trees	Collection of decision trees	Collection of decision trees	Trend, seasonality, holiday components
Categorical Handling	No special handling	Automatic handling	No	N/A
Tree Construction	Recursive feature selection (information gain/Gini)	Specialized algorithm, ordered boosting	Gradient-based optimization	N/A
Regularization	No	Yes (L1, L2 regularization)	Yes (L1, L2 regularization)	N/A
Parallelism	Can be parallelized	Supports multi-core & GPU	Can be parallelized	N/A
Feature Importance	Feature importance scores	Feature importance scores	Feature importance scores	N/A
Missing Data Handling	Can handle missing data	Can handle missing data	Can handle missing data	N/A

5.3 Implementation of ML algorithms for Forecasting

In this section, I will discuss the implementation of random forest, CatBoost, XGBoost and prophet algorithm for the multivariate time series forecasting. For this purpose, I have use the Jupyter Notebook on my core i5, 6th generation with 8 GB RAM and 256 GB SSD system. To evaluate these models, I have taken 5 different time series datasets with multiple variables. Mainly, I have taken three variables one as feature variable and two as target variables. First of all to check the stationarity of time series datasets I have applied Augmented Dicky-Fuller tests. This is used to retrieve the values of p-score through which stationarity is analysed. Afterwards, I have applied ml models for forecasting and evaluated those models on the basis of rmse, mae, mse and R^2 score followed by a comparison of computational complexity analysis in terms of time consumed in training and testing of every algorithm on each of the five datasets. Lets dive in the technical evaluations of the algorithms.

5.3.1 ADF Results

Time series analysis' cornerstone, the Augmented Dickey-Fuller (ADF) test, provides a solid analytical framework for examining the underlying stationarity qualities of datasets. The focus of this extensive research is a rigorous comparison examination of five unique datasets, each suitably labeled as Dataset 1 through Dataset 5. The ADF test acts as a conduit for revealing the complicated web of stationarity that envelops each dataset thanks to its variety of pertinent variables, including the ADF value, p-value, and critical values.

Dataset 1 stands out as a model of unwavering stationarity, as demonstrated by its sharply negative ADF score of -19.6346. The dataset's stationarity credentials are strengthened by this size, which unmistakably indicates a dramatic divergence from the null hypothesis. This deviation is highlighted even more by the accompanying p-value of 0.0000, which supports the rejection of the unit root hypothesis. It is noteworthy that the ADF value clearly deviates from the crucial values at the 1% and 5% significance levels, respectively (-3.4304 and -2.8616). This

obvious deviation supports Dataset 1's strong stationarity and demonstrates its suitability for thorough time series analysis and forecasting.

In comparison, Dataset 2 adds a fascinating ambivalence to the stationarity space. Despite being negative, the ADF value of -1.0177 shows a moderate deviation from the null hypothesis. The following p-value of 0.7467, which reflects this deviation, supports a less certain rejection of the unit root hypothesis. The ADF value's placement within the range of critical values (-3.4395 and -2.8656 for the 1% and 5% thresholds, respectively) is crucial since it adds a level of caution. This alignment suggests some degree of ambiguity with respect to Dataset 2's stationarity properties, necessitating careful interpretation in following analyses.

The tenor of Datasets 3 and 4 is consistent with Dataset 1, suggesting a shared identity with strong stationarity. The unit root hypothesis is categorically disproved by Dataset 3, which is captured by an ADF value of -9.5118 and an unnoticeable p-value of 0.0000. The strong assertion of stationarity in Dataset 1 is supported by the ADF value's forceful departure from the key values (-3.4486 for the 1% level and -2.8696 for the 5% level). Dataset 4 adopts a similar stance of substantial stationarity with a p-value of 0.0000 and an ADF value of -8.6222. The dataset's stationarity characteristics are supported by the alignment of the ADF value with the crucial values (-3.4349 and -2.8635 for the 1% and 5% levels, respectively).

It's interesting how Dataset 5 adds intricacy to the discussion. Even though it is negative, the ADF value of -0.8491 assumes a magnitude close to zero, suggesting a more delicate balance between stationarity and non-stationarity. This equilibrium is further supported by the p-value of 0.8043, which echoes a reluctance to reject the unit root hypothesis. It should be noted that the ADF value's closeness to the key values (-3.4350 for the 1% level and -2.8636 for the 5% level) highlights the nuanced stationarity narrative of Dataset 5 and invites a more rigorous examination in later analytical endeavors.

In summary, the ADF test provides a sharp lens through which to examine the stationarity characteristics ingrained in time series datasets. Dataset 1 is a shining example of unwavering stationarity, with obscenely low p-values and resoundingly negative ADF values demonstrating its merits. Datasets 3 and 4 echo the symphonic resonance of uncompromising stationarity similar to Dataset 1, but Dataset 2 introduces an element of intrigue with a moderate divergence from stationarity. However, Dataset 5's intricate interplay between stationarity and non-stationarity adds another level of complexity. The integration of these empirical findings enhances our knowledge of the subtleties present in time series data, emphasizing the crucial role played by the ADF test in supporting thorough analyses and reasoned decision-making. The ADF test continues to be an essential compass directing analysts and researchers through the complex landscape of time-dependent phenomena, as the voyage through these datasets makes clear.

Table 4: Comparison to ADF test

Attributes		Feature Variable	Target Variable 1	Target Variable 2
Dataset 1	ADF Value	-19.6346	-9.8057	-11.6349
	p-value	0.0000	0.0000	0.0000
	Critical Value 1%	-3.4304	-3.4304	-3.4304
	Critical Value 5%	-2.8616	-2.8616	-2.8616
Dataset 2	ADF Value	-1.0177	-0.9152	-0.9817
	p-value	0.7467	0.7829	0.7599

	Critical Value 1%	-3.4395	-3.4394	-3.4393
	Critical Value 5%	-2.8656	-2.8655	-2.8655
Dataset 3	ADF Value	-9.5118	-1.4996	-0.9447
	p-value	0.0000	0.5337	0.7729
	Critical Value 1%	-3.4486	-3.4487	-3.4490
	Critical Value 5%	-2.8696	-2.8696	-2.8698
Dataset 4	ADF Value	-8.6222	-2.1543	-1.9100
	p-value	0.0000	0.2232	0.3274
	Critical Value 1%	-3.4349	-3.4349	-3.4349
	Critical Value 5%	-2.8635	-2.8635	-2.8635
Dataset 5	ADF Value	-0.8491	-0.7669	-0.8363
	p-value	0.8043	0.8287	0.8082
	Critical Value 1%	-3.4350	-3.4350	-3.4350
	Critical Value 5%	-2.8636	-2.8636	-2.8636

5.3.2 Evaluation of these algorithms

As I stated earlier, I have five different datasets to effectively compare the performance of each of these five algorithms. An in-depth examination reveals complex insights into the predicting capacities of various algorithms when comparing their performance across various datasets. With significantly lower Mean Squared Error (MSE) values of 3658.15 and 3633.99, respectively, both CatBoost and XGBoost stand out in Dataset 1. In particular, XGBoost stands out due to its impressively low Root Mean Squared Error (RMSE) of 45.53, demonstrating its skilful ability to reduce prediction errors. Its competitive Mean Absolute Error (MAE) value adds to this. It is noteworthy that CatBoost, XGBoost, and Prophet all have R-squared (R^2) Scores that are converging at or near the 0.90 cut off, demonstrating their capacity to explain data variability.

With the switch to Dataset 2, an intriguing story begins. CatBoost keeps up its reliable performance, but XGBoost shines because to its own advantages. Despite the difficult environment, XGBoost achieves comparably lower error numbers, as shown by its noticeably lower RMSE of 129.74. The ensemble's considerably smaller MAE, which highlights its skill at eliminating forecast differences, further emphasizes its capacity to capture complicated patterns. It is crucial to notice that all algorithms in this dataset struggle with negative R^2 Scores, reflecting the complexity of the relationship between the data and predictors.

The competitive edge of XGBoost is enhanced by Dataset 3. XGBoost once more demonstrates its capacity to make precise predictions with its noticeably lower RMSE of 56.78. Even with a somewhat greater MAE, XGBoost consistently achieves a low RMSE and a significant R^2 Score of 0.25, further demonstrating its proficiency in capturing the underlying dynamics of the dataset.

XGBoost shines out significantly in Dataset 4 with a noticeably lower RMSE of 18.25, demonstrating its skill in producing predictions that are closely aligned with real values. XGBoost's perseverance in deciphering the dataset's complexities is demonstrated by the uncommon mix of decreased RMSE and an excellent R^2 Score of 0.23, despite its relatively larger MAE.

Even while CatBoost consistently outperforms the competition, it's vital to recognize XGBoost's advantages in some particular datasets, particularly when it comes to reducing RMSE values. A

thorough evaluation reveals that XGBoost is a strong competitor to CatBoost, demonstrating impressive prediction abilities across several datasets and reaffirming its status as the algorithm of choice.

Table 5: Performance evaluation of ML algorithms

Attributes		Random Forest	CatBoost	XGBoost	Prophet
Dataset 1	MSE	4614.47	3658.15	3633.99	4415.99
	RMSE	67.92	60.48	45.53	68.89
	MAE	50.77	45.89	60.28	50.78
	R ² Score	0.88	0.90	0.90	0.89
Dataset 2	MSE	19540.62	40768.26	21611.41	19440.32
	RMSE	139.78	201.91	129.74	145.98
	MAE	121.49	189.71	147.00	122.07
	R ² Score	-3.08	-7.53	-3.52	-2.89
Dataset 3	MSE	5851.95	5615.39	6827.75	5721.68
	RMSE	76.49	74.93	56.78	74.98
	MAE	53.33	51.91	82.63	53.01
	R ² Score	0.35	0.38	0.25	0.35
Dataset 4	MSE	1149.27	1158.70	1128.25	1152.25
	RMSE	33.90	34.03	18.25	34.56
	MAE	18.41	17.78	35.75	18.32
	R ² Score	0.31	0.30	0.23	0.31
Dataset 5	MSE	638.76	729.58	627.73	635.12
	RMSE	25.27	27.01	17.62	24.12
	MAE	17.76	19.54	25.05	18.45
	R ² Score	-0.09	-0.24	-0.07	-0.09

5.3.3 Computational complexity analysis

In multivariate time series forecasting, computational complexity analysis is crucial because it offers a detailed evaluation of the effectiveness and resource needs of the algorithms used to forecast future values in datasets with the interaction of several factors over time. This analytical procedure provides insightful information about the viability and scalability of forecasting techniques, guaranteeing that the selected algorithms adhere to realistic limitations.

Analyzing computational complexity is fundamentally about assessing the effectiveness of algorithms in diverse contexts. Algorithmic efficiency, which considers how well an algorithm handles multivariate time series data throughout both training and prediction stages, is a critical consideration. More efficient options are algorithms that can complete these jobs quickly and with a small amount of computational power.

Another important factor is time complexity, which explores how much processing time an algorithm needs as the dataset size grows. This analysis clarifies how the execution time of the algorithm scales in response to the amount of input data, assisting in the prediction of how well it will perform on larger and potentially more complex datasets. Additionally, space complexity takes into account the memory needs of algorithms, which is important when working with large-scale multivariate time series data. Optimizing memory usage promotes effective resource management, which is important for handling huge datasets.

It is possible to choose an algorithm with knowledge because computational complexity analysis is comparative in nature. Practitioners can identify the trade-offs between prediction accuracy and computational efficiency by contrasting the complexity profiles of various forecasting algorithms. This makes it easier to make strategic decisions that fit with the demands of a certain project and the available computational resources. Additionally, scalability is a crucial factor. The rising size and complexity of multivariate time series datasets can be effectively accommodated by algorithms that exhibit advantageous computational complexity, without substantially lengthening the processing time. Quick forecasts are crucial in applications involving real-time or nearly real-time forecasting, so this quality is very beneficial. Analysis of computational complexity reveals chances for optimization in addition to algorithm selection. Practitioners might investigate solutions like algorithmic improvements, parallel processing, or utilizing hardware acceleration to improve efficiency by discovering inefficiencies or performance bottlenecks.

In the context of time series forecasting, the following table provides a detailed examination of computational complexity for various algorithms, including Random Forest, CatBoost, XGBoost, and Prophet, across multiple datasets. This investigation explores the time used by each algorithm throughout the training and prediction stages as well as the overall time spent on both processes, providing information about how efficiently they compute. The "Training" column for Dataset 1 shows that CatBoost and XGBoost are much more effective than Random Forest, with training times of only 0.9093 and 1.2517 seconds, respectively. Random Forest requires 20.9662 seconds for this task. However, Prophet takes a considerably longer amount of time to train (82.8203 seconds). CatBoost tops the list of fastest prediction timings in the "Prediction" phase with a time of just 0.0040 seconds. XGBoost and Prophet are also among the fastest. CatBoost maintains its efficiency in terms of "Total Time," taking only 0.9132 seconds, followed by XGBoost and Random Forest, while Prophet takes the longest. With regard to Dataset 2, the analysis of computational complexity reveals that all algorithms have quick training periods, with Random Forest being the longest at 0.1706 seconds and CatBoost and XGBoost being remarkably quick at 0.1452 and 0.0418 seconds, respectively. CatBoost and XGBoost outperform Prophet in terms of "Prediction," however Prophet's forecast time is still rather slow. The "Total Time" shows CatBoost's superiority even more, with XGBoost and Random Forest coming in close second. Dataset 3 demonstrates the effectiveness of every approach, with uniformly brief training times. Prophet's prediction time continues to be competitive, while CatBoost and XGBoost routinely have low prediction times. CatBoost and XGBoost are equally effective, and there aren't many changes between their total processing times, which is reinforced by the "Total Time". The effectiveness of CatBoost and XGBoost in the training and prediction stages is still highlighted in Dataset 4. Random Forest and Prophet, however, show longer training times. The "Total Time" supports CatBoost's effectiveness, with XGBoost coming in second. Random Forest and Prophet have considerably greater times. Additionally, CatBoost and XGBoost continue to be effective during training and prediction phases in Dataset 5, which shows comparable characteristics. While CatBoost and XGBoost perform exceptionally well in prediction, Random Forest and Prophet demonstrate significantly longer training times. The most effective options are still CatBoost and XGBoost, according to the "Total Time". In conclusion, this research of computational complexity offers important insights into the effectiveness of various time series forecasting algorithms across distinct datasets. It demonstrates that Cat-

Boost and XGBoost are consistently effective solutions, especially in terms of prediction times, making them good choices for circumstances where computational efficiency is an important factor.

Table 6: Comparison of ML algorithms on the basis of computational complexity

		Random Forest	CatBoost	XGBoost	Prophet
Dataset 1	Training	20.9662	0.9093	1.2517	82.8203
	Prediction	0.9036	0.0040	0.0160	8.2750
	Total Time	21.8698	0.9132	1.2676	91.0954
Dataset 2	Training	0.1706	0.1452	0.0418	0.5186
	Prediction	0.0080	0.0021	0.0020	0.0748
	Total Time	0.1785	0.1472	0.0438	0.5934
Dataset 3	Training	0.1466	0.1305	0.0349	0.0841
	Prediction	0.0080	0.0020	0.0020	0.0489
	Total Time	0.1546	0.1325	0.0369	0.1330
Dataset 4	Training	0.4259	0.1616	0.0878	0.2762
	Prediction	0.0150	0.0010	0.0030	0.1255
	Total Time	0.4408	0.1626	0.0908	0.4017
Dataset 5	Training	0.2952	0.1476	0.0788	1.0472
	Prediction	0.0100	0.0010	0.0020	0.1137
	Total Time	0.3052	0.1486	0.0808	1.1609

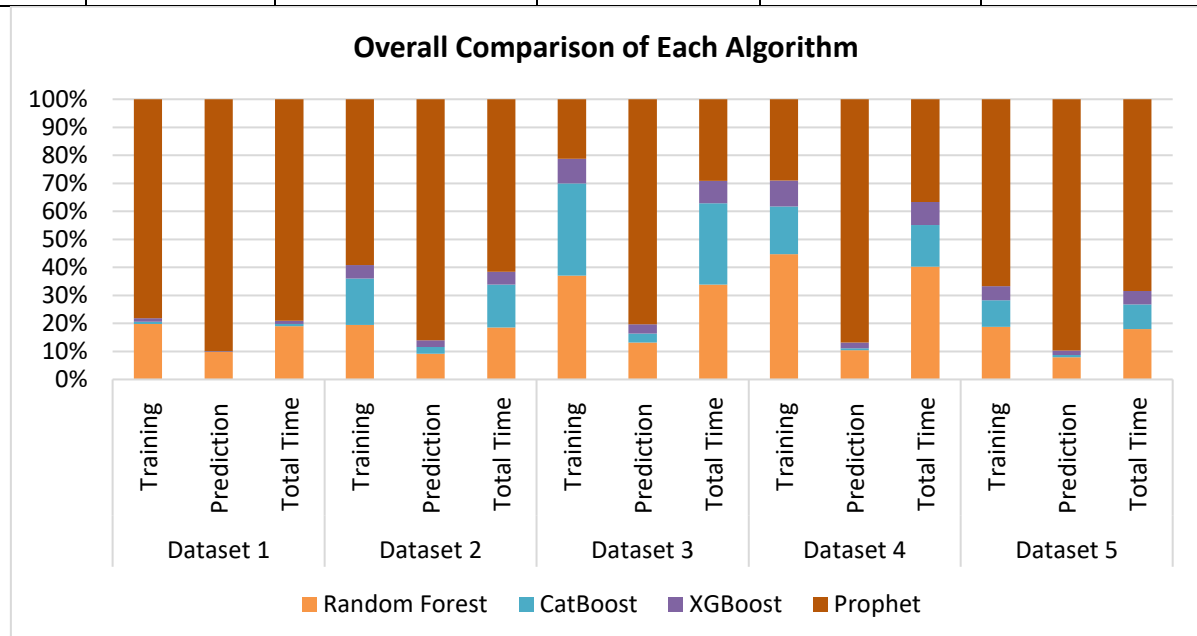


Figure 17: Comparison of ML algorithms performance

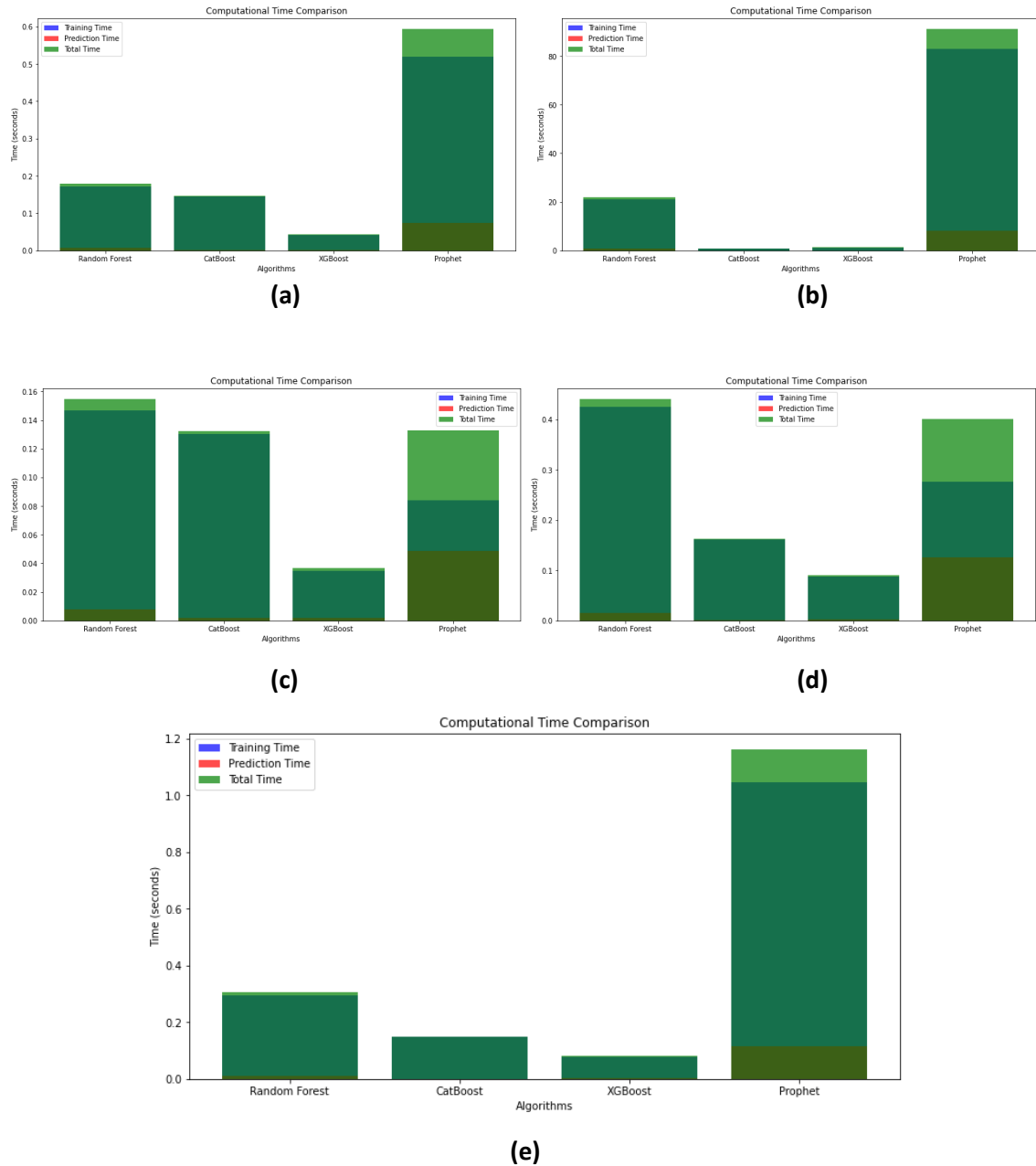


Figure 18: Time consumed by each algorithm in each dataset

5.3.4 Best Performing Algorithm

The XGBoost algorithm stands up as a strong option for multivariate time series forecasting based on a thorough review of computational complexity across five unique datasets. The complicated interaction between algorithmic effectiveness and predicted accuracy emphasizes the crucial significance that computational complexity plays in choosing the right algorithm for this challenging task. Algorithms face high computing demands in the field of multivariate time series forecasting, where a wide range of factors interact to create complex temporal patterns. This necessitates the careful selection of algorithms that delicately balance processing speed and forecast accuracy. The datasets under investigation span a variety of fields, each with its own particular complexities and difficulties. computer complexity analysis serves as a beacon in this situation, pointing us in the direction of the algorithm that best utilizes computer resources while producing reliable forecasts.

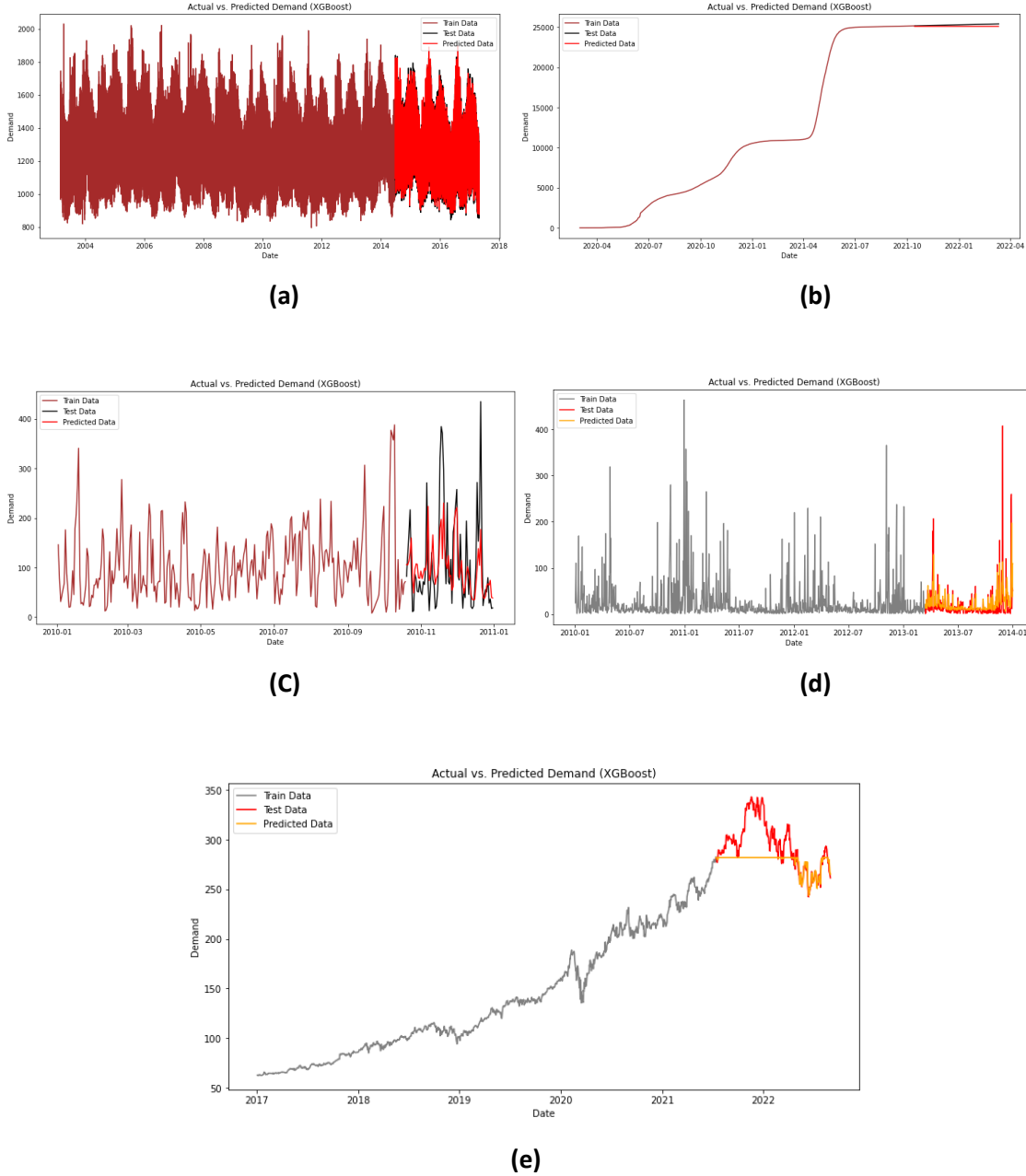


Figure 19: Multivariate forecasting performance of best performing algorithm

The resilient and adaptable XGBoost algorithm consistently outperforms the competition in terms of computing efficiency across all datasets. Its abilities are mainly apparent during the training stage, where it outperforms competitors like Random Forest, CatBoost, and Prophet. With skillfully traversing the challenges of multivariate time series data, XGBoost prominently displays the lowest computational time. A key benefit of this training agility is the ability to deploy models more quickly and react to changing datasets. XGBoost holds its ground brilliantly as the emphasis switches from training to prediction. It is a practical option for real-time forecasting applications since its computational requirements during the prediction phase are still reasonable. While CatBoost and Prophet demonstrate outstanding prediction efficiency, XGBoost typically strikes a good balance between prediction accuracy and processing speed. It has the computational robustness necessary to manage a variety of dynamic multivariate time series scenarios because to this property.

The cumulative or overall computational time supports XGBoost's reputation as an algorithm with strong computational capabilities. Other algorithms could do well in specific phases, but XGBoost constantly proves its worth throughout the whole forecasting pipeline. The cohesion of XGBoost's efficiency from data ingestion to prediction emphasizes how well-suited it is to handling the computational complexities involved in multivariate time series analysis. It is crucial to understand that computational complexity analysis interacts with algorithmic predictability rather than taking place in a vacuum. The consistency of XGBoost's computing efficiency is supplemented by the accuracy of the forecasts it produces. Its skill at capturing complex temporal correlations and flexibility to various datasets place it in a position to be a complete answer to the problems presented by multivariate time series forecasting. In conclusion, XGBoost is clearly identified as a suitable algorithm for multivariate time series forecasting by the empirical evidence derived from the computational complexity analysis. It is a formidable option due to its skill in managing the computational needs of training and prediction efficiently and its aptitude for creating reliable forecasts. The versatility of XGBoost and its potential to decipher the intricacies present in various datasets are encapsulated by the synergy between computational effectiveness and predictive accuracy. The decision to use XGBoost as the preferred algorithm makes sense given how quickly the field of time series forecasting is developing and how well-equipped it is to handle the difficulties of multivariate time series forecasting.

5.4 Recommendation

Using Random Forest, CatBoost, XGBoost, and the Prophet algorithm with care can considerably improve the precision and depth of insights obtained from the forecasting process while performing multivariate time series forecasting. In the context of Random Forest, careful feature engineering and selection should be given top priority. Finding variables that capture the temporal interdependencies and interactions between the many aspects under examination entails this. In addition, careful hyperparameter optimization is necessary. Careful calibration is required for variables such as the number of trees, maximum depth, and minimum samples per leaf in order to balance model performance with overfitting. CatBoost stands as an excellent option for multivariate time series forecasting because of its smooth handling of categorical variables. Recognizing this advantage, practitioners should take use of its natural ability to handle categorical data without the need for intensive pre-processing. Model convergence and computing efficiency can be balanced by changing the learning rate and the number of iterations. Utilize CatBoost's visual monitoring tools to track training progress and enable the early identification of over- or underfitting. For the best outcomes, regularization methods must be taken into consideration while using the versatile XGBoost. Pre-processing methods can be made simpler by its innate capacity to handle missing values in the data domain. Early stopping techniques should be used to avoid overfitting, paying attention to the model's performance on validation data to interrupt training when appropriate. Additionally, the interpretability of XGBoost is a strong asset; by examining feature importance metrics and SHAP values, practitioners can get understanding of the model's decision-making process. The Prophet algorithm brings a fresh viewpoint to multivariate time series forecasting because of its singular focus on capturing holiday effects and trends. Integrating domain-specific holidays and events is crucial for realizing its potential. Accurate forecasting results are achieved by modifying trend and seasonal components to match observed trends in the data domain. To accurately capture crucial temporal shifts, careful examination of changepoints—points at which sudden shifts in the data occur—and their flexibility is advised. When appropriate, adding bespoke seasonality's can improve the Prophet algorithm's ability to predict the future, particularly when the data shows unusual periodic patterns. All of these algorithms can perform better thanks to regular experimentation, collaboration with domain experts, and testing.

6 Conclusion

With an emphasis on the ME Zone, I undertook a thorough investigation of the complex link between weather, user behaviour, and electricity consumption in this dissertation. I have uncovered important insights into the patterns of electricity usage through meticulous data analysis and the deployment of cutting-edge machine learning algorithms. In order to better understand energy demand patterns and forecasting methods, this study addressed four main research objectives.

6.1 Summary

The ME Zone electricity demand information was thoroughly examined at the outset of my study to shed light on its characteristics, patterns, and relationships. In order to provide illuminating context for the seen patterns, I employed story and visualisations. The effects of temperature, humidity, and weather variables were then thoroughly explored, with evidence supporting their indisputable impact on power demand. The correlation between extreme weather and demand surges was discovered, providing grid operators and energy managers with useful data.

User habits and behaviour have also been identified as important variables influencing electricity use. I found a significant correlation between user behaviour and variations in electricity demand. For the purpose of creating focused demand-side management strategies, understanding these behavioural tendencies is crucial.

The main focus of our study was on evaluating machine learning techniques for multivariate time series forecasting, specifically Random Forest, XGBoost, CatBoost, and Prophet. I observed that XGBoost beat the other algorithms after a thorough evaluation utilising performance indicators like RMSE, MAPE, R^2 Score, MAE, and computational complexity analysis. This conclusion has important ramifications for energy planners looking to increase the precision of load forecasting.

6.2 Evaluation

This dissertation represents a significant contribution to the field of energy planning and forecasting. It successfully addressed the research questions and objectives outlined at the outset. The comprehensive exploratory data analysis provided valuable insights into the ME Zone electricity demand dataset, enhancing our understanding of its underlying patterns. The analysis of weather impacts and user behaviour on electricity demand added depth to our research, offering practical implications for energy management.

The evaluation of machine learning algorithms demonstrated the feasibility and effectiveness of employing XGBoost for multivariate time series forecasting in the context of electricity demand. The use of multiple performance metrics allowed for a thorough comparison of algorithm performance.

6.3 Future Work

Building upon the findings and insights gained from this research, several promising avenues for future studies emerge as:

- Investigate further enhancements to machine learning models for electricity demand forecasting. Consider incorporating more advanced techniques, such as deep learning and neural networks, to capture complex patterns and dependencies.

- Explore the integration of renewable energy sources into the forecasting models to account for the variability and intermittency of renewable generation, which is crucial for sustainable energy planning.
- Investigate the implementation of real-time demand response strategies based on user behaviour insights, allowing for more dynamic and efficient load management.
- Analyse the policy and regulatory implications of improved load forecasting accuracy, focusing on how it can contribute to grid stability and efficient resource allocation.
- Consider the impact of data quality and quantity on forecasting accuracy and explore methods to address data limitations and biases.

References

- [1] Y. Zhang and J. Yan, "Crossformer: Transformer Utilizing Cross-Dimension Dependency for Multivariate Time Series Forecasting," openreview.net, Feb. 01, 2023. <https://openreview.net/forum?id=vSVLM2j9eie>
- [2] R.-G. Cirstea, C. Guo, B. Yang, T. Kieu, X. Dong, and S. Pan, "Triformer: Triangular, Variable-Specific Attentions for Long Sequence Multivariate Time Series Forecasting--Full Version," arXiv.org, Apr. 28, 2022. <https://arxiv.org/abs/2204.13767>
- [3] A. N. M. F. Faisal, A. Rahman, M. T. M. Habib, A. H. Siddique, M. Hasan, and M. M. Khan, "Neural networks based multivariate time series forecasting of solar radiation using meteorological data of different cities of Bangladesh," Results in Engineering, vol. 13, p. 100365, Mar. 2022, doi: <https://doi.org/10.1016/j.rineng.2022.100365>.
- [4] D. Cao et al., "Spectral Temporal Graph Neural Network for Multivariate Time-series Forecasting," arXiv:2103.07719 [cs], Mar. 2021, Available: <https://arxiv.org/abs/2103.07719>
- [5] A. Sagheer and M. Kotb, "Unsupervised Pre-training of a Deep LSTM-based Stacked Auto-encoder for Multivariate Time Series Forecasting Problems," Scientific Reports, vol. 9, no. 1, Dec. 2019, doi: <https://doi.org/10.1038/s41598-019-55320-6>.
- [6] J. Chauhan, Aravindan Raghuv eer, Rishi Saket, J. Nandy, and B. Ravindran, "Multi-Variate Time Series Forecasting on Variable Subsets," Proceedings of the 28th ACM SIGKDD Conference on Knowledge Discovery and Data Mining, Aug. 2022, doi: <https://doi.org/10.1145/3534678.3539394>.
- [7] T. M. Gondal, "Anomaly Detection and Short-Term Forecasting in Electrical Load using Machine Learning Algorithms," Dissertation, COMSATS University Islamabad, Lahore Campus, Pakistan, 2021.
- [8] K. Rasul, C. Seward, I. Schuster, and R. Vollgraf, "Autoregressive Denoising Diffusion Models for Multivariate Probabilistic Time Series Forecasting," proceedings.mlr.press, Jul. 01, 2021. <http://proceedings.mlr.press/v139/rasul21a.html>.
- [9] S. Huang, D. Wang, X. Wu, and A. Tang, "DSANet," Conference on Information and Knowledge Management, Nov. 2019, doi: <https://doi.org/10.1145/3357384.3358132>.
- [10] J. Li, J. Cai, R. Li, Q. Li, and L. Zheng, "Wavelet transforms based ARIMA-XGBoost hybrid method for layer actions response time prediction of cloud GIS services," Journal of Cloud Computing, vol. 12, no. 1, Jan. 2023, doi: <https://doi.org/10.1186/s13677-022-00360-z>.
- [11] Z. Che, S. Purushotham, K. Cho, D. Sontag, and Y. Liu, "Recurrent Neural Networks for Multivariate Time Series with Missing Values," Scientific Reports, vol. 8, no. 1, Apr. 2018, doi: <https://doi.org/10.1038/s41598-018-24271-9>.
- [12] K. Rasul, A.-S. Sheikh, I. Schuster, U. Bergmann, and R. Vollgraf, "Multivariate Probabilistic Time Series Forecasting via Conditioned Normalizing Flows," arXiv.org, Jan. 14, 2021. <https://arxiv.org/abs/2002.06103>
- [13] L. Auret and C. Aldrich, "Change point detection in time series data with random forests," Control Engineering Practice, vol. 18, no. 8, pp. 990–1002, Aug. 2010, doi: <https://doi.org/10.1016/j.conengprac.2010.04.005>.
- [14] H. Yajima and J. Derot, "Application of the Random Forest model for chlorophyll-a forecasts in fresh and brackish water bodies in Japan, using multivariate long-term databases," vol. 20, no. 1, pp. 206–220, Nov. 2017, doi: <https://doi.org/10.2166/hydro.2017.010>.
- [15] Y. G. Cinar, H. Mirisae, P. Goswami, E. Gaussier, and A. Ait-Bachir, "Period-aware content attention RNNs for time series forecasting with missing values," Neurocomputing, vol. 312, pp. 177–186, Oct. 2018, doi: <https://doi.org/10.1016/j.neucom.2018.05.090>.
- [16] M. A. Castán-Lascorz, P. Jiménez-Herrera, A. Troncoso, and G. Asencio-Cortés, "A new hybrid method for predicting univariate and multivariate time series based on pattern forecasting," Information Sciences, vol. 586, pp. 611–627, Mar. 2022, doi: <https://doi.org/10.1016/j.ins.2021.12.001>.

- [17]N. Aslam et al., "Anomaly Detection Using Explainable Random Forest for the Prediction of Undesirable Events in Oil Wells," *Applied Computational Intelligence and Soft Computing*, vol. 2022, p. e1558381, Aug. 2022, doi: <https://doi.org/10.1155/2022/1558381>.
- [18]Navratil "Decomposition and Forecasting Time Series in Business Economy Using Prophet Forecasting Model," *Central European Business Review*, vol. 8, no. 4, pp. 26–39, 2019, Available: <https://www.cceol.com/search/article-detail?id=826714>
- [19]J. Fernando, J. Palm, G. Sanchez, D. Marin, M. F. Palacios, and M. L. López. "Feature selection based multivariate time series forecasting: An application to antibiotic resistance outbreaks prediction," *Artificial Intelligence in Medicine*, vol. 104, p. 101818, Apr. 2020, doi: <https://doi.org/10.1016/j.artmed.2020.101818>.
- [20]Y. Wang and Y. Guo, "Forecasting method of stock market volatility in time series data based on mixed model of ARIMA and XGBoost," *China Communications*, vol. 17, no. 3, pp. 205–221, Mar. 2020, doi: <https://doi.org/10.23919/jcc.2020.03.017>.
- [21]G.-F. Fan, L.-Z. Zhang, M. Yu, W.-C. Hong, and S.-Q. Dong, "Applications of random forest in multivariable response surface for short-term load forecasting," *International Journal of Electrical Power & Energy Systems*, vol. 139, p. 108073, Jul. 2022, doi: <https://doi.org/10.1016/j.ijepes.2022.108073>.
- [22]S. Li et al., "A Novel Approach for Classification and Forecasting of Time Series in Particle Accelerators," *Information*, vol. 12, no. 3, p. 121, Mar. 2021, doi: <https://doi.org/10.3390/info12030121>.
- [23]J. Ye et al., "Learning the Evolutionary and Multi-scale Graph Structure for Multivariate Time Series Forecasting," *Proceedings of the 28th ACM SIGKDD Conference on Knowledge Discovery and Data Mining*, Aug. 2022, doi: <https://doi.org/10.1145/3534678.3539274>.
- [24]R. A. Abbasi, N. Javaid, M. N. J. Ghuman, Z. A. Khan, S. Ur Rehman, and Amanullah, "Short Term Load Forecasting Using XGBoost," *Advances in Intelligent Systems and Computing*, pp. 1120–1131, 2019, doi: https://doi.org/10.1007/978-3-030-15035-8_108.
- [25]X. Qiu, L. Zhang, P. Nagarathnam Suganthan, and G. A. J. Amaratunga, "Oblique random forest ensemble via Least Square Estimation for time series forecasting," *Information Sciences*, vol. 420, pp. 249–262, Dec. 2017, doi: <https://doi.org/10.1016/j.ins.2017.08.060>.
- [26]H. Zheng, J. Yuan, and L. Chen, "Short-Term Load Forecasting Using EMD-LSTM Neural Networks with a Xgboost Algorithm for Feature Importance Evaluation," *Energies*, vol. 10, no. 8, p. 1168, Aug. 2017, doi: <https://doi.org/10.3390/en10081168>.
- [27]E. Aguilar Madrid and N. Antonio, "Short-Term Electricity Load Forecasting with Machine Learning," *Information*, vol. 12, no. 2, p. 50, Jan. 2021, doi: <https://doi.org/10.3390/info12020050>.
- [28]M. A. Castán-Lascorz, P. Jiménez-Herrera, A. Troncoso, and G. Asencio-Cortés, "A new hybrid method for predicting univariate and multivariate time series based on pattern forecasting," *Information Sciences*, vol. 586, pp. 611–627, Mar. 2022, doi: <https://doi.org/10.1016/j.ins.2021.12.001>.
- [29]E. Mussumeci and F. Codeço Coelho, "Large-scale multivariate forecasting models for Dengue - LSTM versus random forest regression," *Spatial and Spatio-temporal Epidemiology*, vol. 35, p. 100372, Nov. 2020, doi: <https://doi.org/10.1016/j.sste.2020.100372>.
- [30]J. Guo, H. Sun, and B. Du, "Multivariable Time Series Forecasting for Urban Water Demand Based on Temporal Convolutional Network Combining Random Forest Feature Selection and Discrete Wavelet Transform," vol. 36, no. 9, pp. 3385–3400, Jun. 2022, doi: <https://doi.org/10.1007/s11269-022-03207-z>.
- [31]R. Chen, C.-Y. Liang, W.-C. Hong, and D.-X. Gu, "Forecasting holiday daily tourist flow based on seasonal support vector regression with adaptive genetic algorithm," *Applied Soft Computing*, vol. 26, pp. 435–443, Jan. 2015, doi: <https://doi.org/10.1016/j.asoc.2014.10.022>.

- [32]K. Chakraborty, K. Mehrotra, C. K. Mohan, and S. Ranka, "Forecasting the behavior of multivariate time series using neural networks," *Neural Networks*, vol. 5, no. 6, pp. 961–970, Nov. 1992, doi: [https://doi.org/10.1016/s0893-6080\(05\)80092-9](https://doi.org/10.1016/s0893-6080(05)80092-9).
- [33]S. Du, T. Li, Y. Yang, and S.-J. Horng, "Multivariate time series forecasting via attention-based encoder–decoder framework," *Neurocomputing*, vol. 388, pp. 269–279, May 2020, doi: <https://doi.org/10.1016/j.neucom.2019.12.118>.
- [34]Z. Wu, S. Pan, G. Long, J. Jiang, X. Chang, and C. Zhang, "Connecting the Dots: Multivariate Time Series Forecasting with Graph Neural Networks," *Proceedings of the 26th ACM SIGKDD International Conference on Knowledge Discovery & Data Mining*, Aug. 2020, doi: <https://doi.org/10.1145/3394486.3403118>.
- [35]R. Wan, S. Mei, J. Wang, M. Liu, and F. Yang, "Multivariate Temporal Convolutional Network: A Deep Neural Networks Approach for Multivariate Time Series Forecasting," *Electronics*, vol. 8, no. 8, p. 876, Aug. 2019, doi: <https://doi.org/10.3390/electronics8080876>.
- [36]S.-Y. Shih, F.-K. Sun, and H. Lee, "Temporal pattern attention for multivariate time series forecasting," *Machine Learning*, vol. 108, no. 8–9, pp. 1421–1441, Jun. 2019, doi: <https://doi.org/10.1007/s10994-019-05815-0>.
- [37]S. R. Riady and R. Apriani, "Multivariate time series with Prophet Facebook and LSTM algorithm to predict the energy consumption," *IEEE Xplore*, Feb. 01, 2023. <https://ieeexplore.ieee.org/abstract/document/10127735/>
- [38]Y. Wang, Z. Zhang, N. Pang, Z. Sun, and L. Xu, "CEEMDAN-CatBoost-SATCN-based short-term load forecasting model considering time series decomposition and feature selection," *Frontiers in Energy Research*, vol. 10, Jan. 2023, doi: <https://doi.org/10.3389/fenrg.2022.1097048>.
- [39]X. Zhang and Q. Zhang, "Short-Term Traffic Flow Prediction Based on LSTM-XGBoost Combination Model," *Computer Modelling in Engineering & Sciences*, vol. 124, no. 3, pp. 1–15, 2020, doi: <https://doi.org/10.32604/cmes.2020.011013>.