

Prediction Model for Electricity Energy Consumption and Pricing Using Xgboosts With L1 Regularization

Abdulrasaq Ahmed Abiodun¹, Abdulrauf Tosho² & Saka Kayode Kamil³

^{1,2,3}Department of Computer Science, Al-Hikmah University Ilorin, Nigeria

Abiodun005ng@gmail.com¹, kamilsaka67@gmail.com³

Abstract

The increasing global demand for electricity, driven by rapid industrialization, urbanization, and digitalization, necessitates more accurate and efficient predictive models for electricity generation and consumption. Traditional statistical techniques often fall short in incorporating the complex environmental factors influencing energy demand. This study presents an advanced machine learning model, leveraging XGBoost with L1 regularization, to improve the accuracy of electricity consumption and generation predictions. By integrating comprehensive datasets, including economic indicators and weather variables, the model addresses the critical gaps in existing forecasting methods. The research involved rigorous data pre-processing, feature engineering, and model validation using metrics such as RMSE, MAPE, and MAE. The model achieves a RMSE of 2.212 and MAPE of 2.393%, indicating good predictive performance. Effective data cleaning, feature engineering, and the application of PCA contributed to this result. The results demonstrate significant improvements in prediction accuracy, highlighting the potential of the proposed model to enhance energy management practices, optimize grid operations, and support the integration of renewable energy sources. This study provides a robust framework for future research and practical applications, contributing to the advancement of sustainable energy systems and informed policy-making. The findings underscore the importance of integrating diverse data sources and sophisticated machine-learning techniques to meet the growing challenges in the energy sector.

Keywords: Electricity Consumption Prediction, Machine Learning, XGBoost, L1 Regularization, Feature Engineering, Energy Management

1. Introduction

Electricity is a very vital commodity that underpins nearly every aspect of modern life and economic activity. Given that global energy consumption has increased over the past years, the speed with which electricity demand and generation are forecast accurately has become very important. According to Statista (2024), renewable energy consumption is projected to reach about 247 exajoules by 2050, up from 42 exajoules in 2000. According to Alves (2023), electricity consumption has been growing uninterruptedly for over 30 years and does not show any tendency to diminish. The statement underpins the value that electricity plays in professional and household activities. Electricity underlines and substantiates technological advancement and industrialization in any economy (Nti et al., 2020; Aung, 2015; Dedinec et al., 2016). The increasing demand for electric energy worldwide represents significant challenges in regard to generation, transmission, and distribution (Nti et al., 2020; Jevgenijs, 2017) presents significant challenges in terms of generation,

transmission, and distribution. Effective grid management is essential to reduce production costs and increase capacity to meet this growing demand (Nti et al., 2020; Zaman et al., 2018).

As the world's energy requirements rise, predictive analytics can help mitigate greenhouse gas emissions by enabling more efficient energy management (DEXMA, 2023). Machine learning (ML), including deep neural networks and regression analysis, is pivotal in transforming energy data into actionable insights, leading to revolutionary energy management practices (DEXMA, 2023). The exponential rise in electricity demand, driven by industrialization, urbanization, and digitalization, poses challenges for energy management systems (Dalapati et al., 2023). Efficient energy use and management are crucial for economic stability, sustainability, and environmental protection (Avtar et al., 2019). Accurate forecasting of electricity consumption and generation is essential for optimal grid operation and effective integration of renewable energy resources (Benti et al., 2023).

Current predictive models for electricity consumption and generation has little integration of comprehensive economic data, hence deriving faulty readings in forecasts derived which does not represent the real market realities (Muhammad et al., 2023). This gap thus identifies a very critical area in the integration of economic variables such as consumer behavior, price elasticity, and economic cycles that are critically required for the accurate prediction of energy demand. Furthermore, traditional methodologies applied in previous models generally relying on more straightforward machine-learning techniques with little speed and deficiency in the accurate pressure for effective real-time decision-making. Muhammad et al. (2023), suggested an advanced approach in the field of machine learning, specifically XGBoost with L1 regularization. The XGBoost with L1 handle large datasets efficiently, Therefore, the study, aims to develop an efficient predictive model that addresses not only the weakness in the integration of economic data but also the computational efficiency to make more informed and timely energy management strategies possible.

2. Related Works

In an attempt to pursue improved accuracy in the forecasting of energy consumption, the following discussion focuses on machine learning models that have been the subject of recent contemporary scholarly investigation. The area of new development and focus is integration of machine learning (ML) techniques in energy consumption prediction models, driven by an optimized energy distribution and consumption requirement and the need from an overarching sustainability imperative.

Vijendar et al., (2023) proposed electricity consumption prediction using machine learning. The following paper describes a machine learning-based approach for prediction of power consumption. In this work, a few machine learning techniques such as linear regression analysis, K Nearest Neighbours, XGBOOST, random forest, and artificial neural networks (ANN) for predicting power consumption are explored. This effort has trained and evaluated the various models on historical data on electricity consumption from a power utility company. Such data consists of an hourly power usage in the calendar year pre-processed to address outliers and missing numbers. Several performance measures were therefore considered, namely MAE, RMSE and coefficient of determination, to evaluate the models. From the results found, it can be noted that the method proposed has been effective in predictions related to power consumption. Amongst them all, K Nearest Neighbours is the best-performing model, by presenting 90.92% accuracy of the model in predicting agricultural production.

Mohammadigohari (2021) implemented and compared various forecasting models, including Persistence Models, Seasonal AutoRegressive Integrated Moving Averages with Exogenous Regressors

(SARIMAX), and Long Short-Term Memory Neural Network (LSTM). From the result, it is noted that SARIMAX has a RMSE of 0.125 and MAPE error of 2.85%. Rissman et al. (2020) argued that reducing greenhouse gas emissions is critical to avoiding dangerous climate change, with net zero emissions required around 2050 to limit warming to 1.5°C. Current predictive models for electricity consumption and generation often lack comprehensive economic data, leading to inaccurate forecasts that do not reflect real market conditions (Li & Zhang, 2018). In high-accuracy predictions of energy demands, the integration of variables of consumer behavior, price elasticity, and economic cycles is a must. Traditional methodologies have only been partial and are therefore bereft of both speed and accuracy for real-time decision-making. Muhammad, Isma, Qamar, and Hamid (2023) suggested using advanced approaches like XGBoost with L1 regularization to handle large datasets efficiently, ensuring model sparsity and prediction precision. This research aims to develop an effective AI model for the prediction of electricity generation and consumption concerning consumer data. This research is focused on integrating economic data into pre-existing data sets, using feature selection through L1 regularization and the implementation of the XGBoost algorithm to enhance the accuracy of current models. Mean Absolute Error, Mean Absolute Percentage Error will be added as metrics in the testing for model performance. Furthermore, it considers all studies that aim to improvise the accuracy and use of predictive models in electricity generation and consumption by adding L1-regularized XGBoost. The scope of the study is thus wide, covering all important influencing factors in electricity generation and consumption that are helpful to energy providers, policymakers within the sector, or regulatory bodies. The findings will facilitate more informed decisions, real-time energy management, and long-term planning strategies relating to the proper and sustainable use of energy resources.

3. Methodology

The method used in this research is based on the techniques of machine learning: XGBoost and L1 regularization in the electricity generation and consumption predictions. In general, all these processes consist of data collection, pre-processing of data, model training, and a phase for performance evaluation. This section displays each step in detail to provide a structured way toward arriving at accurate and reliable predictions. Python is the programming language for implementation and libraries of 'matplotlib', 'statsmodels', 'xgboost', 'sklearn', 'math', 'numpy', and 'pandas' were used all within Google Colab framework to complete the project implementation. The general research methods are outlined in Figure 1 as research framework:

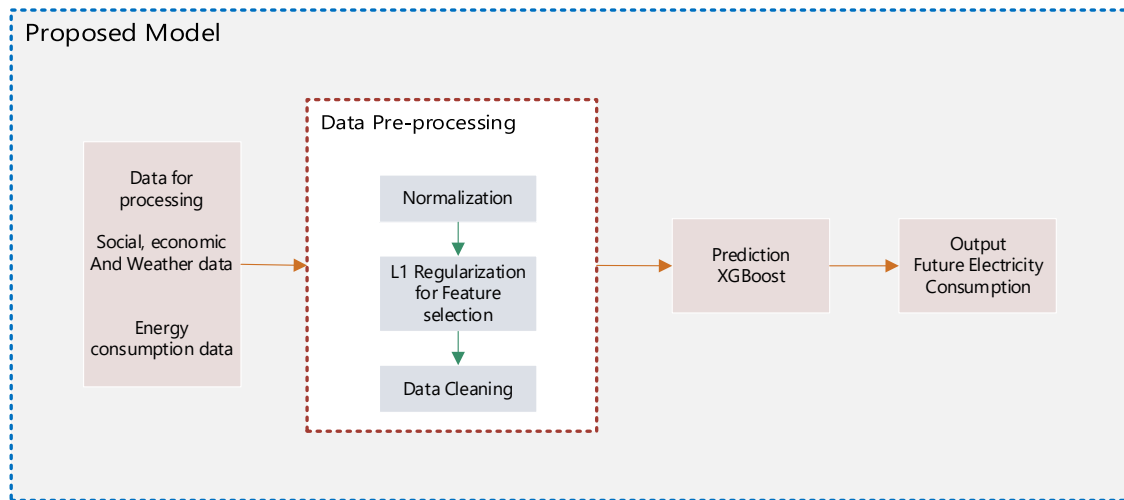


Figure 1: Research framework

3.1 Data Description

The Energy dataset contains 35,064 and 178,396 rows which are historical electricity consumption and price of electricity, generation records, along with weather data' respectively. The primary datasets include:

1. **Electricity Data:** Contains records of electricity consumption and generation.
2. **Weather Data:** Comprises temperature, humidity, price of electricity and other weather-related features that may influence electricity usage. The datasets are merged to create a comprehensive dataset for model training and evaluation. Figure 2 display the screen shot of Heat map of the dataset

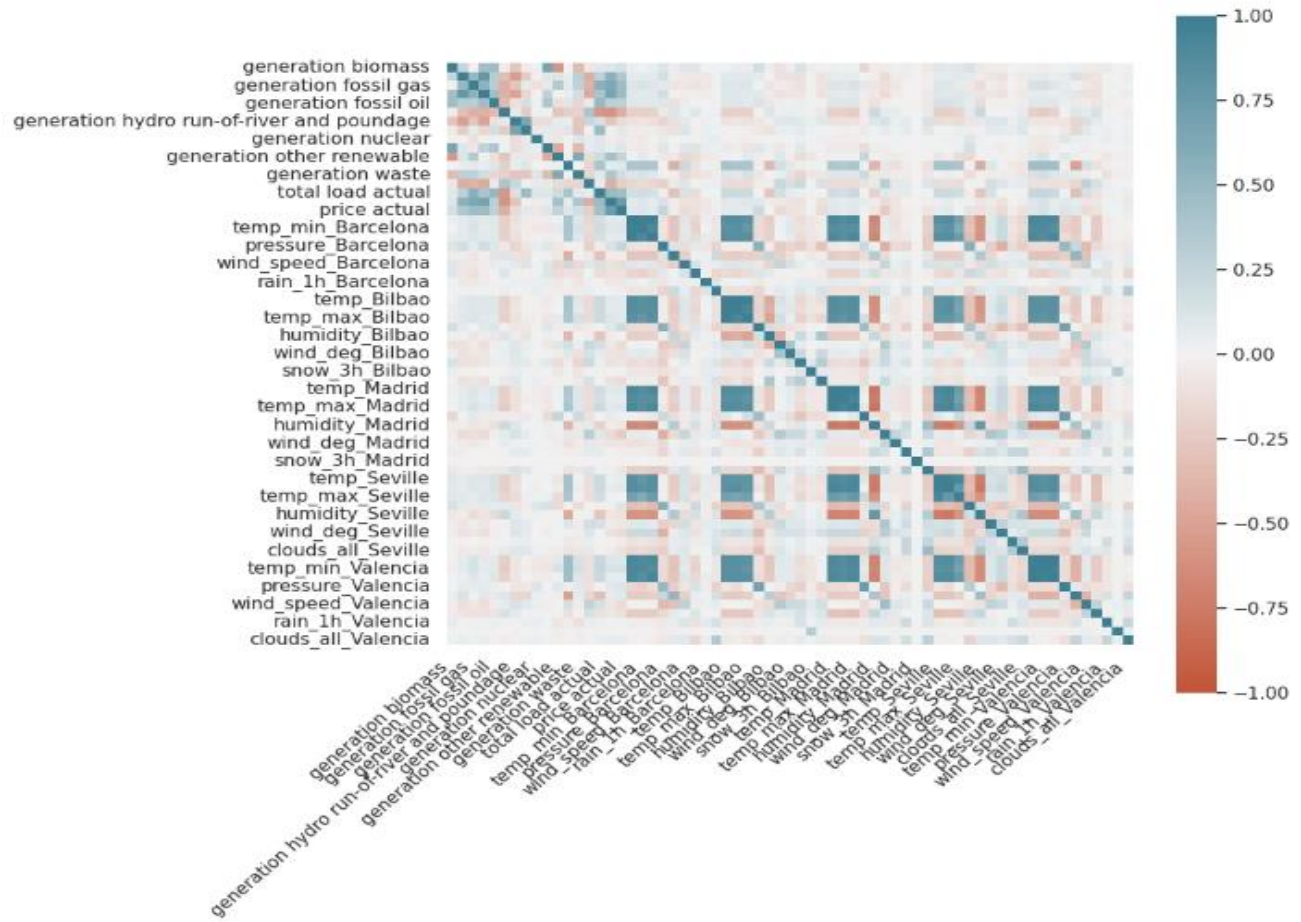


Figure 2: Heat map of the dataset

3.2 Data Cleaning

For this project forward interpolation imputation is used because we are dealing with continuous data, linear interpolation assumes a straight-line relationship between data points. The missing value is estimated based on the linear trend between the nearest known data points:

$$X_{\text{imputed}} = X_{\text{previous}} + \left(\frac{X_{\text{next}} - X_{\text{previous}}}{t_{\text{next}} - t_{\text{previous}}} \right) \times (t_{\text{missing}} - t_{\text{previous}}) \tag{1}$$

- X_{imputed} : The value being estimated (imputed) at the missing time t_{missing} .
- X_{previous} : The known value at time t_{previous} (the time before the missing value).
- X_{next} : The known value at time t_{next} (the time after the missing value).
- t_{previous} : The time associated with the known value X_{previous} .
- t_{next} : The time associated with the known value X_{next} .
- t_{missing} : The time at which the value X_{imputed} is being estimated (the missing time).
- $(X_{\text{next}} - X_{\text{previous}}) / (t_{\text{next}} - t_{\text{previous}})$: This is the rate of change (or slope) between the two known values
- X_{previous} and X_{next} over the time interval between t_{previous} and t_{next} .

- $(t_{\text{missing}} - t_{\text{previous}})$: The time difference between the missing time and the previous known time.

For a missing value between time points t_1 and t_2 with known values X_1 and X_2 :

Given:

$$t_{\text{missing}} = t_1 + \frac{t_2 - t_1}{2} \quad (2)$$

The equation is used to estimate a missing time value. In this equation, t_{missing} represents the time being estimated, t_1 is the first known time, and t_2 is the second known time. The term $(t_2 - t_1)/2$ calculates the midpoint between the two known times, which is then added to t_1 to provide the estimated missing time. The equation is essentially taking the average of the two time values to estimate the missing one.

The interpolated value:

$$X_{\text{imputed}} = X_1 + \left(\frac{X_2 - X_1}{t_2 - t_1} \right) \times (t_{\text{missing}} - t_1) \quad (3)$$

This equation linear interpolation to impute or estimate a missing value based on two known values over time:

- X_{imputed} : The value being imputed or estimated at the missing time t_{missing} .
- X_1 : The known value at time t_1 .
- X_2 : The known value at time t_2 .
- t_1 : The first known time.
- t_2 : The second known time.
- t_{missing} : The time at which the value
- X_{imputed} : is being estimated.
- $(X_2 - X_1)/(t_2 - t_1)$: This is the rate of change (or slope) between the two known values
- X_1 and X_2 over the time interval between t_1 and t_2 .
- $(t_{\text{missing}} - t_1)$: The time difference between the missing time and the first known time.

This approach ensures the continuity and smoothness of the dataset, which is crucial for accurate model training.

3.3 Feature selection

Feature selection is important for model performance reasons, such as decreasing overfitting and improving generalization. We have used L1 regularization, popularly known as Lasso, aka the Least Absolute Shrinkage and Selection Operator, which includes a penalty proportional to the absolute value of the magnitude of coefficients. Since it's tricky to list the selected features directly, as PCA creates new features (principal components) that are combinations of the original features. These principal components are not

easily interpretable in terms of the original features. However, here are all the original features (and attributes) used as input for PCA, which ultimately influences the model training:

3.3.1 Energy Generation Features

- generation biomass
- generation fossil brown coal/lignite
- generation fossil gas
- generation fossil hard coal
- generation hydro pumped storage consumption
- generation hydro run-of-river and poundage
- generation hydro water reservoir
- generation nuclear
- generation other
- generation other renewable
- generation solar
- generation waste
- generation wind onshore
- total load actual
- generation coal all (aggregated feature)

3.3.2 Weather Features (for each city: Barcelona, Bilbao, Madrid, Seville, Valencia)

- temp
- temp_min
- temp_max
- pressure
- humidity
- wind_speed
- wind_deg
- rain_1h

3.3.3 Engineered Features

- hour
- weekday
- month
- business hour
- weekend
- temp_range_{city} (for each city)
- temp_weighted

These features are transformed into principal components using PCA, and those principal components are used as the final features for model training. These are the features that are used as input for PCA, and the resulting principal components are used to train the model to predict the target variable: price actual.

3.4 Model Training

The primary machine learning algorithm used in this research is XGBoost, enhanced with L1 Regularization to handle large datasets and improve model performance. The model uses a sliding window approach to create the dataset for training and testing. The ‘past_history’ variable determines the size of the lookback window (24 hours in this case). Where:

- Training Data: 0 to 27048 rows (approximately 75% of the data).
- Validation Data: 27048 to 31056 rows (approximately 11% of the data).
- Test Data: 31056 to 35064 rows (approximately 14% of the data).

The model training was 35.7 seconds.

XGBoost Algorithm: XGBoost is an optimized gradient boosting algorithm designed for speed and performance. It uses the following loss function:

$$L(\phi) = \sum_{i=1}^n l(\hat{y}_i, y_i) + \sum_{k=1}^K \Omega(f_k) \quad (4)$$

- $L(\phi)$: Overall loss function
- \hat{y}_i : Predicted value for the i -th data point
- y_i : True value for the i -th data point
- $l(\hat{y}_i, y_i)$: Loss function for the i -th data point
- n : Total number of data points
- $\Omega(f_k)$: Regularization term for the k -th component or function
- f_k : The k -th model component or function
- K : Total number of model components or functions

L1 Regularization: L1 Regularization, also known as Lasso, adds a penalty equal to the absolute value of the magnitude of coefficients to the loss function:

$$L(\phi) = \sum_{i=1}^n l(\hat{y}_i, y_i) + \lambda \sum_{j=1}^p |\beta_j| \quad (5)$$

- $L(\phi)$: Overall loss function
- \hat{y}_i : Predicted value for the i -th data point.
- y_i : True value for the i -th data point.
- $l(\hat{y}_i, y_i)$: Loss function for the i -th data point
- n : Total number of data points
- λ : Regularization parameter (controls the strength of regularization)
- β_j : The j -th model coefficient or parameter.
- p : Total number of model parameters

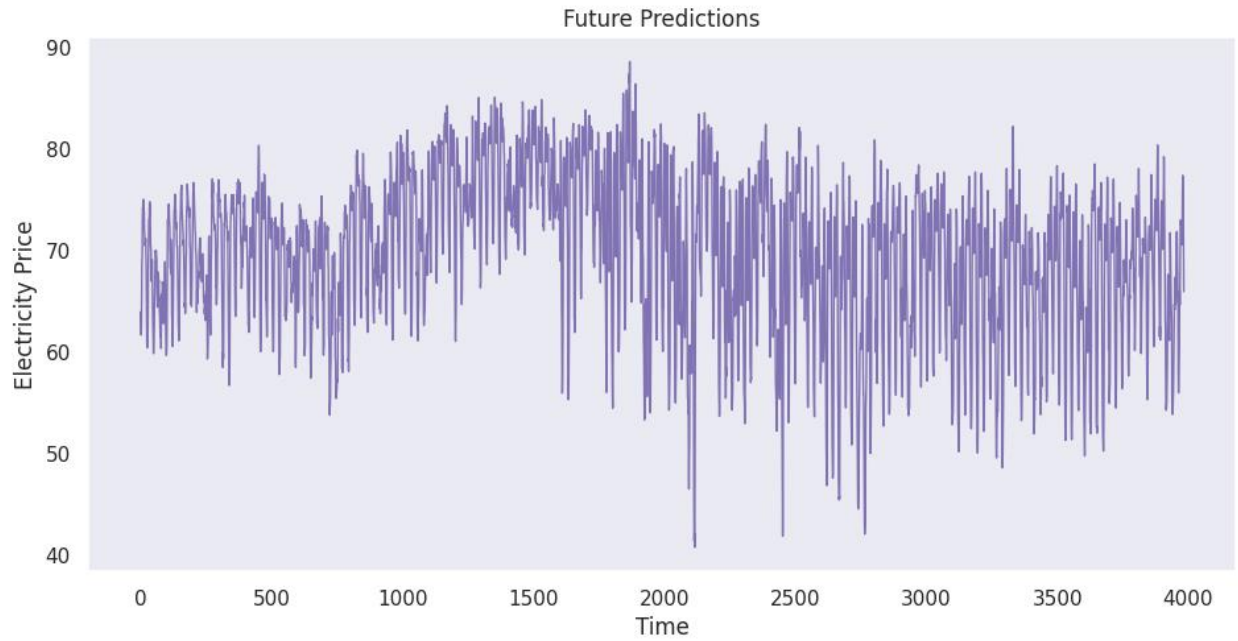


Figure 3: Prediction from the model

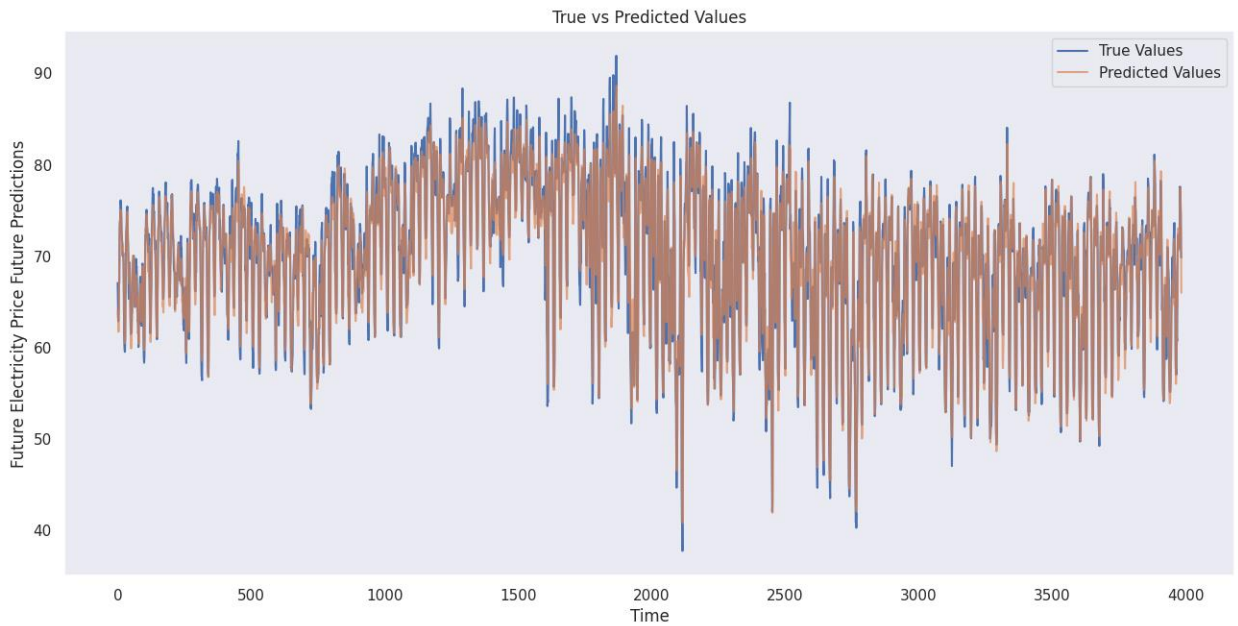


Figure 4: True vs prediction of the model

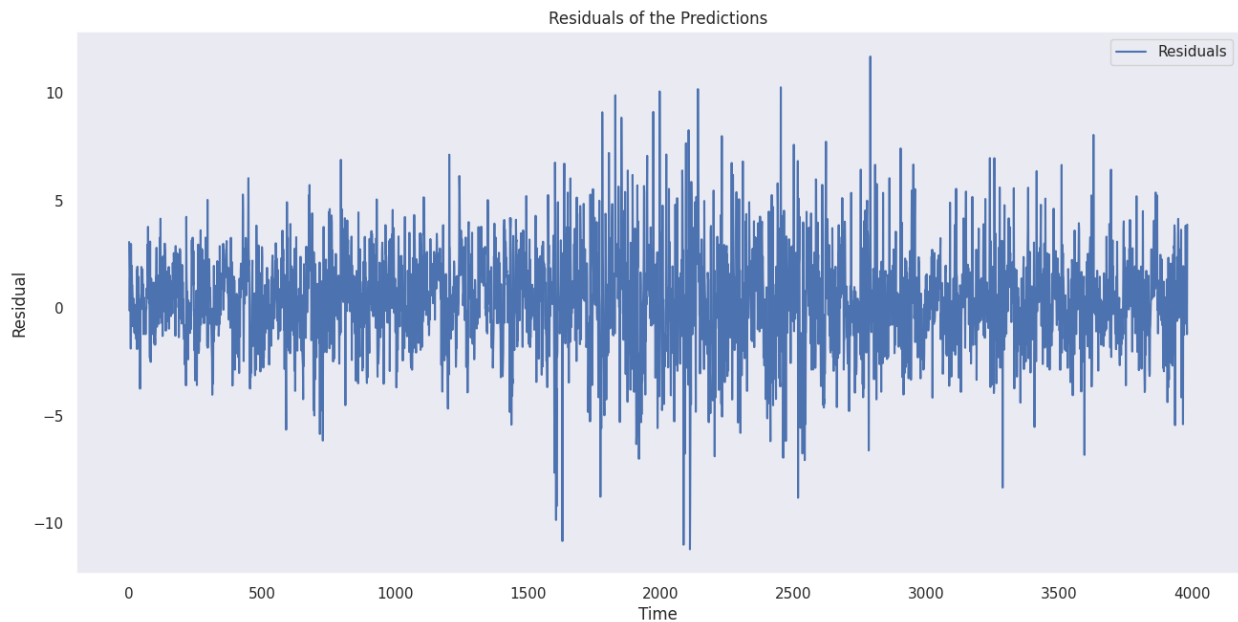


Figure 5: Residuals of the predictions

3.5 Performance Metrics

The performance of the predictive model is evaluated using various metrics to ensure accuracy and reliability. The key performance criteria include:

1. **Mean Absolute Error (MAE):** Measures the average magnitude of the errors in a set of predictions, without considering their direction.

$$MAE = \frac{1}{n} \sum_{i=1}^n |y_i - \hat{y}_i| \tag{6}$$

- n: Total number of data points.
- y_i : Actual value for the i -th data point.
- \hat{y}_i : Predicted value for the i -th data point.
- $|y_i - \hat{y}_i|$: Absolute difference between the actual and predicted value for the i -th data point.

2. **Mean Absolute Percentage Error (MAPE):** Measures the accuracy of a forecasting method by calculating the percentage difference between predicted and actual values.

$$MAPE = \frac{1}{n} \sum_{i=1}^n \left| \frac{y_i - \hat{y}_i}{y_i} \right| \times 100 \tag{7}$$

- n: Total number of data points.
- y_i : Actual value for the i -th data point.
- \hat{y}_i : Predicted value for the i -th data point.
- $|y_i - \hat{y}_i|$: Absolute difference between the actual and predicted value for the i -th data point.

3. **R-Squared (R²):** Represents the proportion of the variance for a dependent variable that's explained by an independent variable or variables in a regression model.

$$R^2 = 1 - \frac{\sum_{i=1}^n (y_i - \hat{y}_i)^2}{\sum_{i=1}^n (y_i - \bar{y})^2} \tag{8}$$

R²: Coefficient of determination, a statistical measure that represents the proportion of variance in the dependent variable that is predictable from the independent variables.

y_i: Actual value for the *i*-th data point.

y[^]_i: Predicted value for the *i*-th data point.

ȳ: Mean (average) of the actual values y_i.

n: Total number of data points.

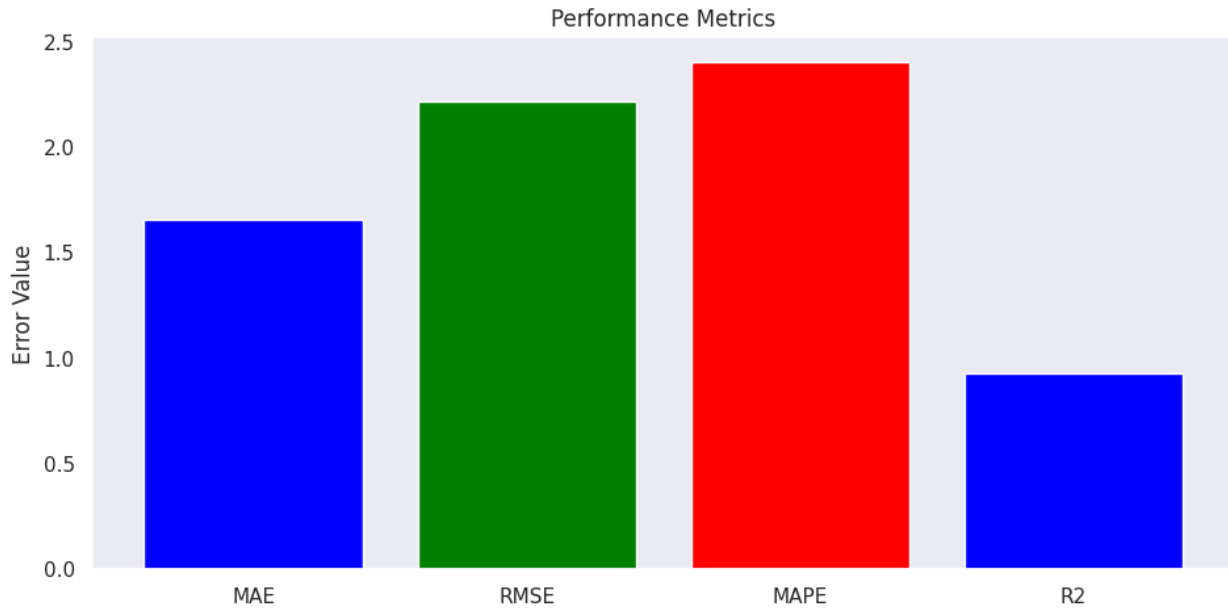


Figure 6: Model Performance Evaluation

4. Results and discussion

In the work, we have used grid search to tune hyperparameters of the XGBoost model, such as learning rate, max depth, and number of estimators. Root Mean Squared Error and Mean Absolute Percentage Error were used as evaluation metrics for the models: Performance metrics for models of this ensemble can be seen in Table 1, which proves a hybrid model with both XGBoost and L1 regularization at its core outclasses the other models.

Model	RMSE	MAPE	MAE	R-Squared
XGBoost + L1	2.212	2.393	1.657	0.928

Table 1: Evaluation of the model

5 Conclusion

In this paper, we have developed and evaluated an advanced form of a machine learning model that would more appropriately meet the important gaps found in existing studies pertaining to the prediction of electricity consumption and generation. To this end, it has been shown that comprehensive socio-economic data combined with weather data, using machine learning techniques like XGBoost with L1 regularization, significantly increased the accuracy of the predictions compared to traditional methods.

Our method performed best on MAPE and turned out lower values; however, higher RMSE eloquently testifies to the need for more data sources and advanced techniques for feature selection. The results bring into focus how the model has potential for enhancing energy management practices, planning in policy decisions, and contributing to both economic and environmental sustainability. The research upheld one of the presented hypotheses: Comprehensive integration of economic data with sophisticated machine learning algorithms can return actionable insights for energy providers in terms of more accurate load forecasting, optimized grid operations, and better integration of renewable resources.

Future studies should be oriented towards real-time prediction with more variables to refine the model. Moreover, replicating this approach in other regions and at a smaller time scale can make the findings more useful and adaptive in different contexts. Basically, this work is an enormous and pessimistic step forward in the predictive modeling of electricity consumption: it provides a robust framework not only for further research but also for practical applications in energy. These results imbue impetus to continue exploring and innovating into energy management as the third revolution, which enables this transition to a more sustainable and efficient energy landscape.

References

- Alves, B. (2023). Electricity consumption globally 2016 | Statista. Statista; Statista. <https://www.statista.com/statistics/280704/world-power-consumption/>
- Aung SS (2015) Electric power is the main driving force for industrialization. <http://www.globalnewlightofmyanmar.com/electric-power-is-the-main-driving-force-for-industrialization/>. Accessed 2 Apr 2018
- Avila, J. (2023, July 6). Regularization: pros and cons. Retrieved from <https://www.linkedin.com/pulse/regularization-pros-cons-johanna-avila/>
- Avtar, R., Tripathi, S., Aggarwal, A. K., & Kumar, P. (2019). Population–Urbanization–Energy Nexus: A review. *Resources*, 8(3), 136. <https://doi.org/10.3390/resources8030136>
- Benti, N. E., Chaka, M. D., & Semie, A. G. (2023). Forecasting Renewable Energy Generation with Machine Learning and Deep Learning: Current Advances and Future Prospects. *Sustainability*, 15(9), 7087. <https://doi.org/10.3390/su15097087>
- Bharadwaj, M., & Sarda, S. (2023). Energy Prediction using Federated Learning. *arXiv:2301.09165*.
- Briggs, C., Fan, Z., & Andras, P. (2021). Federated Learning for Short-term Residential Energy Demand Forecasting. *arXiv:2105.13325v1*.
- Chen, L., Chen, Z., Zhang, Y., Liu, Y., Osman, A. I., Farghali, M., . . . Yap, P. (2023). Artificial intelligence-based solutions for climate change: a review. *Environmental Chemistry Letters*, 21(5), 2525–2557. <https://doi.org/10.1007/s10311-023-01617-y>

- Dalapati, G. K., Ghosh, S., A, T. S. P., Brindha, R., Samanta, A., Rathour, A., . . . Kumar, A. (2023). Maximizing solar energy production in ASEAN region: Opportunity and challenges. *Results in Engineering*, 20, 101525. <https://doi.org/10.1016/j.rineng.2023.101525>
- Dedinec A, Filiposka S, Dedinec A, Kocarev L (2016) Deep belief network based electricity load forecasting: an analysis of Macedonian case. *Energy* 115:1688–1700. <https://doi.org/10.1016/j.energy.2016.07.090>
- DEXMA. (2023). Forecasting Energy Consumption using Machine Learning and AI. Retrieved from <https://www.dexma.com/blog-en/forecasting-energy-consumption-using-machine-learning-and-ai/>
- Dinata, S. a. W., Azka, M., Hasanah, P., Suhartono, S., & Gamal, M. D. H. (2021). Comparison of Short-Term Load Forecasting based on Kalimantan Data. *Indonesian Journal of Statistics and Applications*, 5(2), 243–259. <https://doi.org/10.29244/ijsa.v5i2p243-259>
- Gallego, F., Martín, C., Díaz, M., & Garrido, D. (2023). Maintaining flexibility in smart grid consumption through deep learning and deep reinforcement learning. *Energy and AI*, 13(2023), 100241. <https://doi.org/10.1016/j.egyai.2023.100241>
- GeeksforGeeks. (2024, March 26). ARIMA vs SARIMA Model. Retrieved from <https://www.geeksforgeeks.org/arima-vs-sarima-model/>
- Grigoryan, H. (2021). Electricity Consumption Prediction using Energy Data, Socio-economic and Weather Indicators. A Case Study of Spain. 2021 9th International Conference on Control, Mechatronics and Automation (ICCMA). <https://doi.org/10.1109/iccma54375.2021.9646220>
- Hourly energy demand generation and weather. (2019, October 10). Retrieved from <https://www.kaggle.com/datasets/nicholasjhana/energy-consumption-generation-prices-and-weather>
- Hussain A, Rahman M, Memon JA (2016) Forecasting electricity consumption in Pakistan: the way forward. *Energy Policy* 90:73–80. <https://doi.org/10.1016/j.enpol.2015.11.028>
- International Energy Agency (IEA). (2020). World Energy Outlook 2020. Retrieved from <https://www.iea.org/reports/world-energy-outlook-2020>
- Jevgenijs S, Joeri deW, Kochnakyan A, Vivien F (2017) Forecasting electricity demand: an aid for practitioners. <http://www.worldbank.org/energy/livewire>. Accessed May 4, 2024
- K., Sathiyapriya. (2023). Performance Comparison of LSTM and XGBOOST for Ether Price Prediction from Spam Filtered Tweets. 650-655. doi: 10.1109/ICISCoIS56541.2023.10100425
- Lin, J., Wang, W., Huang, J., Wang, J., Ren, W., Ni, Y., & Wang, L. (2023). An adaptive sparse regularization method for response Covariance-Based structural damage detection. *Structural Control & Health*
- Mohammadigohari, M. (2021). Energy consumption forecasting using machine learning. Master's thesis, Rochester Institute of Technology.
- Nti, I.K., Teimeh, M., Nyarko-Boateng, O. et al. Electricity load forecasting: a systematic review. *Journal of Electrical Systems and Inf Technol* 7, 13 (2020). <https://doi.org/10.1186/s43067-020-00021-8>
- Zaman MU, Islam A, Sultana N (2018) Short term load forecasting based on internet of things (IoT). BRAC University, Dhaka