

Short-Term Load Forecasting Using Machine Learning Models

Kanika Kuchinad*, Sangeeta Modi†

*Undergraduate Student, Dept. of Electrical and Electronics Engineering

Email: kanikakuchinad02@gmail.com

†Assistant Professor, Dept. of Electrical and Electronics Engineering

Email: smodi@pes.edu

PES University, Bangalore, India

Abstract—Machine learning has emerged as a powerful tool for solving complex real-world problems across various domains, including energy systems. The main objective of this paper is to offer a short-term electric load forecasting framework for solar-powered microgrids using three machine learning models—Random Forest, Support Vector Regression (SVR) and Linear Regression. Accurate load forecasting is very important due to the intermittent nature of solar energy. The system uses historical weather and load data to predict the load demand for the next 48 hours and also gives the comparison between the actual load versus the predicted load to study how efficiently the algorithm is working. Through comparative analysis of the three models, it is studied that the Random Forest model offers the best performance in terms of Root Mean Square Error (RMSE) and R-squared (R^2) metrics. The paper emphasises how data-driven techniques improve energy dependability, lower running costs, and allow more intelligent energy dispatch in distributed renewable systems.

Index Terms—Load forecasting, microgrid, machine learning, Random Forest, SVR, Linear Regression, renewable energy.

I. INTRODUCTION

Nowadays, the usage of solar microgrids is increasing due to the demand for clean and localized energy solutions. However, the variability and unreliability in solar energy generation introduces challenges in energy management. This research suggests making a predictive model with machine learning technologies and testing their capacity to forecast the electricity demand and compares their performance in predicting the electricity demand of a solar-powered microgrid. The dataset used in this study [11] is from a power station and contains real-time weather parameters such as solar irradiance, temperature, and wind speed, along with corresponding load data. This data is instrumental in forecasting energy demand and solar generation potential in microgrid systems..

II. LITERATURE REVIEW

Recent work in short-term load forecasting (STLF) has explored both classical and deep learning models.

Raju et al. [1] applied traditional machine learning models for STLF in smart grids, focusing on regression trees and random forests. However, they did not incorporate temporal features or multi-step forecasting, which limits applicability in dynamic grid environments.

Syed et al. [2] introduced a clustering-based DNN framework for grouped consumption profiles. While this reduced complexity, it lacked temporal inputs such as lagged loads, which are a key component of our feature set.

Masood et al. [3] developed a Quantile LSTM model integrated with clustering to deliver probabilistic forecasts. Although robust under uncertainty, it requires high computational resources and lacks real-time efficiency.

Azeem et al. [4] conducted a thorough review of forecasting models across generation types. Their analysis is broad but theoretical, with no direct experimentation or comparison like in our approach.

Alquthami et al. [5] benchmarked several ML models including enhanced decision trees. Their limited dataset and lack of exogenous weather inputs differentiate our broader and more practical setup.

Syed et al. [6] proposed a scalable global model for distribution network forecasting. However, it underutilized time-sensitive features like lag values, which we explicitly include.

Masood et al. [7] also proposed a deep LSTM encoder-decoder model for residential forecasting. While effective, it is less interpretable and computationally heavier than our ensemble-based Random Forest model.

Park et al. [8] proposed an attention-based deep learning framework combining GRU and LSTM units for STLF. While accurate, it is heavily reliant on deep architectures, limiting interpretability and energy efficiency.

Zhou et al. [9] applied XGBoost and LightGBM for load forecasting, showing improved accuracy over traditional models. However, their analysis was limited to single-household data and ignored long-term dependencies.

Wang et al. [10] introduced a hybrid CNN-LSTM model for demand forecasting. Their architecture captures spatial and temporal features but is prone to overfitting in small or noisy datasets—issues mitigated in our Random Forest-based approach.

III. DATASET AND PREPROCESSING

The dataset is taken from a powerstation that utilizes solar energy. It contains hourly weather and electric load records. Features include irradiance, temperature, dewpoint, specific humidity, and wind speed, along with electric load in MW.

A. Preprocessing Steps

- Extracted timestamp from year, month, day, and hour columns.
- Filled missing values using the previous value.
- Normalized all feature columns.
- Added derived features: Day of Week and Hour of Day.
- Introduced lag features (previous 3 hours of electric load).

IV. METHODOLOGY

A. Feature Selection

The final feature set includes irradiance, temperature, dew-point, specific humidity, wind speed, day of week, hour of day, and lagged electric load values.

B. Train-Test Split

80% of the data has been used for training and 20% for testing, with a forecast horizon of 48 hours.

C. Models Used

- **Linear Regression:** It is the most basic machine learning model for regression. It assumes linear relationship between input features and the target variable. It is fast and easy to interpret. However, it doesn't work very well with non-linear patterns. It works well only when there is a linear relationship between the input features and the output, i.e., the target. It is easy and fast and saves computing time. It may not be feasible for real world application.
- **Support Vector Regression (SVR):** SVR uses a Gaussian kernel to capture nonlinear relationships. It is effective for small data sets, but sensitive to feature scaling. Unlike Linear Regression, SVR doesn't penalize every small error—instead, it only adjusts the model when the prediction error exceeds a threshold, ε . SVR solves the following optimization problem:

$$\min \frac{1}{2} \|w\|^2 \quad \text{subject to} \quad |y_i - \hat{y}_i| \leq \varepsilon$$

where:

- w : Weight vector of the model, representing the orientation of the regression hyperplane.
 - $\|w\|^2$: Squared norm of the weight vector, minimized to ensure flatness of the regression function.
 - y_i : Actual target value for the i -th data point.
 - \hat{y}_i : Predicted value from the model for the i -th data point.
 - $|y_i - \hat{y}_i| \leq \varepsilon$: Epsilon-insensitive loss function, meaning small errors (within ε) are ignored.
 - ε : Tolerance margin that defines an acceptable range of error around the predicted value.
- **Random Forest:** Random Forest is a type of ensemble learning algorithm. It builds multiple decision trees during training using a technique called bagging, and gives the average output of all the decision trees as the final result. In this way, it handles complex interactions and avoids overfitting as taking the result of only one

decision tree might be too specific, thus it combines all the decision trees' output. Each tree is trained on a random part of the dataset, which means that each tree gets a different subset of the dataset to learn from. It handles large datasets and many features.

V. RESULTS AND DISCUSSION

A. Evaluation Metrics

To compare the performance of the machine learning models, two standard regression evaluation metrics are used: Root Mean Square Error (RMSE) and the R-squared (R^2) score. These metrics help in understanding the prediction accuracy and generalization ability of each model.

1) Root Mean Square Error (RMSE):

RMSE measures the average magnitude of the prediction error between the actual and predicted values. It is defined as:

$$\text{RMSE} = \sqrt{\frac{1}{n} \sum_{i=1}^n (y_i - \hat{y}_i)^2}$$

where y_i is the actual load, \hat{y}_i is the predicted load, and n is the number of predictions. RMSE is expressed in the same unit as the target variable (MW), making it directly interpretable. A lower RMSE indicates higher prediction accuracy. Since RMSE penalizes large errors more strongly due to squaring, it is sensitive to outliers.

2) R-squared (R^2 Score):

R^2 measures the proportion of variance in the actual values that is captured by the model. It is computed as:

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

where \bar{y} is the mean of the actual load values. An R^2 score of 1.0 represents perfect predictions, while a score of 0 means the model performs no better than simply predicting the mean. A higher R^2 score reflects better generalization to unseen data.

3) Why Both Metrics Are Important:

RMSE gives insight into the real-world magnitude of prediction errors, which is essential for operational planning in microgrids. R^2 , on the other hand, evaluates how well the model captures the trend and variability of energy consumption. Together, they provide a balanced view of model performance: one measures precision, the other measures explanatory power.

In this study, Random Forest achieved the lowest RMSE and highest R^2 score, indicating that it not only predicts accurately but also effectively models the load behavior patterns in solar-powered microgrid environments.

B. Model Performance Comparison

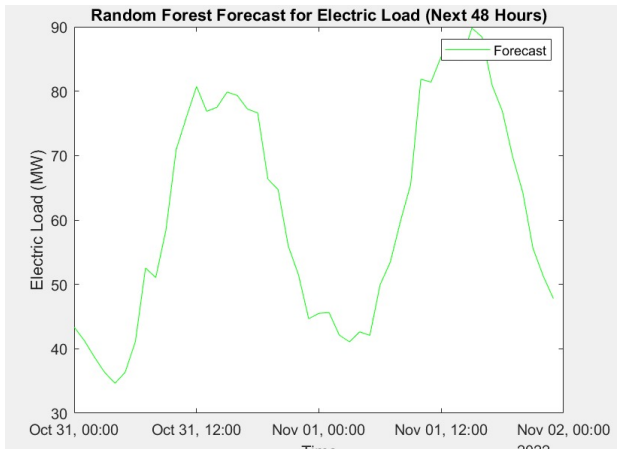


Fig. 1. Forecast using Random Forest

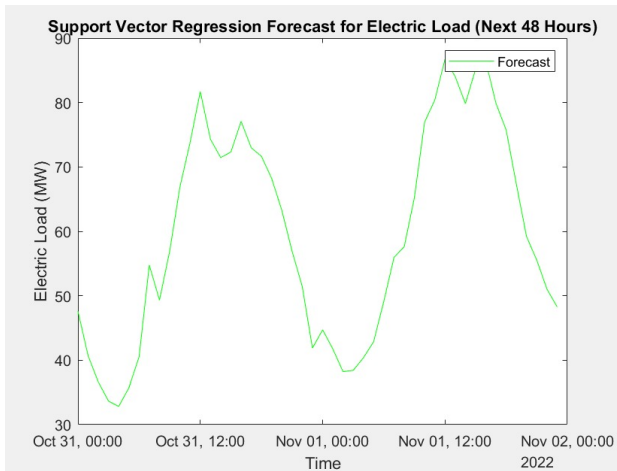


Fig. 2. Forecast using SVR

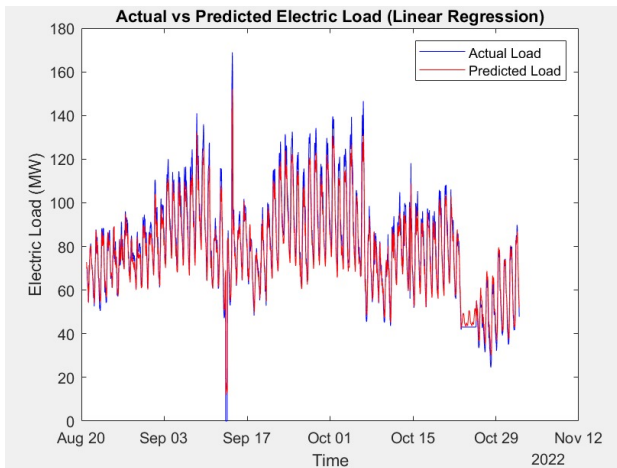


Fig. 3. Actual vs. Predicted Load (Linear Regression)

TABLE I
MODEL PERFORMANCE COMPARISON

Model	RMSE	R ² Score
Linear Regression	7.3567	0.8964
SVR	16.5911	0.4733
Random Forest	7.0669	0.9044

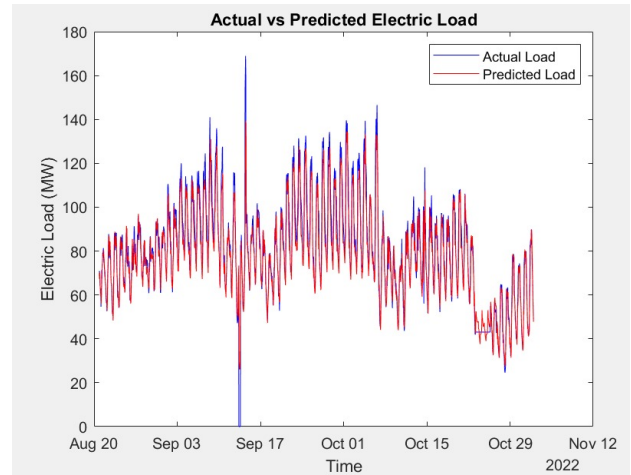


Fig. 4. Actual vs. Predicted Load (Random Forest)

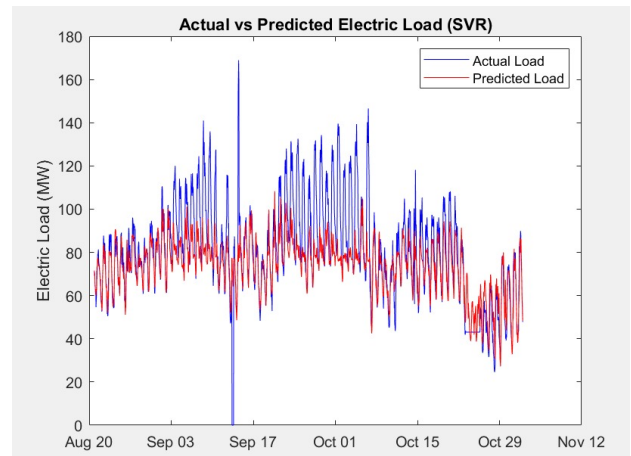


Fig. 5. Actual vs. Predicted Load (SVR)

C. Discussion

Three distinct methods of machine learning, Linear Regression, Support Vector Regression (SVR), and Random Forest, were employed to test the performance of the models using two different conventional metrics, the Root Mean Square Error (RMSE) and the R-squared (R^2) score. These metrics were computed using the predictions made on a 20% test split of the dataset.

1) Linear Regression: Linear Regression, which is the simplest among the three models, acted as a base line for comparison. Although the model was the simplest and made the linear assumptions, it gave a good account of itself since it reached an RMSE of 7.36 and an R^2 score of 0.896. In this regard, it can be said that the model follows the main direction of the data pretty well, especially when the relation between features (such as temperature, irradiance, and hour of day) and load is close to linear. Nevertheless, its result was a bit less than the Random Forest model. This model lagged because of its poor performance in dealing with nonlinear patterns and extensive feature interactions. This implies that although Linear Regression is able to quickly estimate the load, it lacks the ability to cover more detailed variations that may be present in a real-world load behavior.

2) Support Vector Regression (SVR):

SVR was predicted to perform better because of its capacity to model non-linear relationships via the Gaussian (RBF) kernel. Still, it hugely underperformed in comparison to both Linear Regression and Random Forest, yielding an RMSE of 16.59 and an R^2 of only 0.473. This steep decline in performance may be due to several reasons:

- SVR is feature-scaling sensitive. Even though normalization was used, it may not have been ideal for all features.
- The default hyperparameters were possibly not suited to the nature of this dataset. Tuning parameters such as the regularization parameter C , kernel coefficient γ , and epsilon ϵ might have enhanced performance.
- SVR might find larger or high-dimensional datasets difficult to model unless tuned properly, which may have affected its generalizability.

This indicates that SVR does show promise, but it requires more computational resources and significant hyperparameter optimization effort to match or surpass the performance of tree-based models like Random Forest on this task.

3) Random Forest: Random Forest emerged as the best-performing model, with an RMSE of 7.07 and an R^2 of 0.904. These results confirm that Random Forest is capable of handling the complex, non-linear relationships inherent in microgrid load forecasting. It leverages an ensemble of decision trees trained on randomly sampled subsets of the data and features, which helps it avoid overfitting and remain robust to noise. Unlike Linear Regression and SVR, Random Forest does not assume any specific functional form between input and output variables, making it suitable for real-world energy demand data that often exhibits erratic or seasonally varying patterns.

4) Comparative Insights:

Practically, the enhanced precision of Random Forest has important implications for microgrid operation. Precise short-term load forecasting enables improved scheduling of distributed energy resources and hybrid energy storage systems. With enhanced demand prediction, system operators can:

- Optimize battery and supercapacitor utilization within the HESS.
- Decrease over-reliance on backup generators.
- Reduce power losses and enhance energy dispatch strategies.
- Increase the overall sustainability and cost-effectiveness of the microgrid.

Conversely, the poor generalization capacity of SVR here renders it less appropriate without rigorous tuning. Linear Regression is still a good quick model for initial estimates or where interpretability is important, but its shortcomings are more evident when dealing with complex patterns.

5) Visual Observations:

The graphed forecasts further illustrate the differences in performance. The predicted values of the Random Forest model closely follow actual load profiles with minimal variance even under spikes in load. SVR, however, makes smoother predictions that do not pick up peaks or abrupt falls, resulting in greater residual errors. Linear Regression, though not as distant from actual as SVR, is still missing the finer trends because of its linearity.

Generally speaking, Random Forest offers the best-balanced and dependable method among the models experimented with, and hence is the most appropriate to be deployed in smart grid energy management systems.

VI. CONCLUSION

The study was designed to learn the impact of three machine learning models—Linear Regression, the Support Vector Regression (SVR) and Random Forest—for short-term load forecasting in a solar-powered microgrid. Out of the models that were examined, Random Forest consistently gave the best results with its lowest Root Mean Square Error (RMSE) of 7.07 and the highest R^2 score of 0.904. The fact that it is an ensemble model and its versatility to obtain non-linear interactions made it the most good at generalization and the most capable of tracking the fluctuations in energy consumption patterns. The findings establish the suitability of Random Forest for load forecasting in distributed energy systems, especially when there is a high demand for accuracy and reliability are the primary requirements. According to the study, Linear Regression, although being quick and easy to use, failed to satisfactorily model non-linear connections between input features and the load. It only achieved a good enough performance level so that for systems with less complex or variable datasets, it would still be a possible option. Nevertheless, it might not give consistent results when it comes to a significant change in the system inherent in the case of solar microgrids. The main reason for SVR's poor performance is that it is hard to tune the hyperparameters

and scaling the data is also an issue. The discrepancy in the RMSE and R^2 is a clear indicator that SVR will not be a realistic choice for problem-solving. The most significant point of this paper is that it emphasizes the importance of selecting the suitable forecasting algorithm for energy systems. With accurate load prediction, the efficiency of the microgrid is directly improved through the renewable resources and energy storage components' dispatch that is made to be optimized. It also results in a reduction in the occurrence of the system imbalance, the lower running of the operational costs as well as reliable support to intermittent solar energy. In conclusion, the Random Forest model is the most resilient, easy to understand, and realistic means of approach for load forecasting.

VII. FUTURE WORK

- Integration with Model Predictive Control (MPC) for real-time optimization.
- Implementation of deep learning methods like LSTM for temporal pattern recognition.
- Incorporation of features such as battery SOC and solar panel output.

REFERENCES

- [1] L. Raju, V. Easwaramoorthy, V. V., and K. M. Vimalan, "Application of Machine Learning Algorithms for Short term Load Prediction of Smart Grid," in *Proc. Int. Conf. Smart Electronics and Communication (ICOSEC)*, IEEE, 2020, pp. 1712–1716.
- [2] D. Syed, H. Ullah, M. Usama, M. A. Khan, and K. Muhammad, "Deep Learning-Based Short-Term Load Forecasting Approach in Smart Grid With Clustering and Consumption Pattern Recognition," *IEEE Access*, vol. 9, pp. 54989–55000, 2021.
- [3] Z. Masood, R. Gantassi, and Y. Choi, "Enhancing Short-Term Electric Load Forecasting for Households Using Quantile LSTM and Clustering-Based Probabilistic Approach," *IEEE Access*, vol. 12, pp. 77256–77268, 2024.
- [4] A. Azeem, M. Arif, M. A. Khan, A. W. Malik, and A. Paul, "Electrical Load Forecasting Models for Different Generation Modalities: A Review," *IEEE Access*, vol. 9, pp. 142239–142268, 2021.
- [5] T. Alquthami, M. Aljohani, M. A. Khan, and M. A. Jan, "A Performance Comparison of Machine Learning Algorithms for Load Forecasting in Smart Grid," *IEEE Access*, vol. 10, pp. 48420–48433, 2022.
- [6] D. Syed, M. Usama, K. Muhammad, M. A. Khan, and R. Amin, "A Global Modeling Framework for Load Forecasting in Distribution Networks," *IEEE Access*, vol. 9, pp. 93322–93335, 2021.
- [7] Z. Masood, R. Gantassi, Y. Choi, and M. Farooq, "A Deep Learning Method for Short-Term Residential Load Forecasting in Smart Grid," *IEEE Access*, vol. 12, 2024. (Accepted, in press).
- [8] D. Kaur, S. N. Islam, M. A. Mahmud, M. E. Haque, and Z. Dong, "Energy Forecasting in Smart Grid Systems: A Review of the State-of-the-art Techniques," *arXiv preprint arXiv:2011.12598*, 2020.
- [9] J. Guo, Y. Peng, Q. Zhou, and Q. Lv, "Enhanced LSTM Model for Short-Term Load Forecasting in Smart Grids," in *Proc. 9th EAI Int. Conf. Cloud Computing, Smart Grid and Innovative Frontiers in Telecommunications (CloudComp)*, 2020, pp. 461–476.
- [10] Y. Zhang, X. Zhang, and Y. Wang, "Deep ResNet-Based Ensemble Model for Short-Term Load Forecasting in Protection System of Smart Grid," *Sustainability*, vol. 14, no. 24, p. 16894, 2022.
- [11] Weather and Load Forecasting Dataset, Kaggle. [Online]. Available: <https://www.kaggle.com/datasets/itsaru/solar-power-generation-data>