

A Comparison of State of The Art Machine Learning Models and Auto-Regressive Generative AI for Time Series Forecasting Tasks on Electrical Grid Data

ADITYA UPADHYAY¹, UDAY SINGH², and ISHANI ARYA, ³

¹University of Malta, Msida, MSD 2080 MALTA (e-mail: aditya.upadhyay.23@um.edu.mt)

²University of Malta, Msida, MSD 2080 MALTA (e-mail: uday.singh.23@um.edu.mt)

³University of Malta, Msida, MSD 2080 MALTA (e-mail: ishani.arya.23@um.edu.mt)

This work was supported by the Faculty of ICT, University of Malta, Msida MSD 2080, Malta

ABSTRACT With the expansion of renewable energy sources, a new problem has come up. The problem of varying energy generation has posed a significant challenge to traditional grid management techniques. The implementation of smart grids can solve the problem. These grids, which have distribution systems that are more efficient, are able to overcome the limitations that were posed by the technology that was used in the past. These intelligent grids generate a substantial amount of data, which can subsequently be utilized to forecast future load and make adjustments to distribution. This project investigates the use of Apache Spark for the analysis of large amounts of data as well as a transformer model for the purpose of load estimation. This transformer model was tested against three machine-learning models: Random Forest, CatBoost, and LightGBM. The results reveal that this model is substantially more efficient and performs better across all evaluation metrics.

INDEX TERMS Big Data Analysis, Load Forecasting, Predictive Analytics, Smart Grids, Transformer Model

I. INTRODUCTION

A Smart grid is a network of electricity generators and consumers connected through the internet, efficiently matching demand and supply. This may be achieved using sensors, software, and advanced models.

Generic and traditional electric grids face several issues of imbalance in demand and supply, electrical outages, security threats, etc. Another issue with these grids is the inability to integrate renewable energy systems. Decentralized renewable energy systems require a bidirectional flow of electricity. Utilizing a Smart grid to optimize load prediction makes it possible to overcome the problems of electric cuts and changing consumption based on factors like trends and time of the day.

Utilizing a Smart grid to optimize load prediction makes it possible to overcome the problems of electric cuts and changing consumption based on factors like trends and time of the day. Using a Smart grid, we can collect data from various sources, including Sensors and meters. Also, certain data can be collected from external sources like weather and demographics.

A. DATASET

TABLE 1: Electrical Usage Data Sample.

LCLid	Usage (Kilowatt-hour (half-hourly))	Acorn
MA000084	0.42	Aflt
MA000084	0.31	Aflt
MA000084	0.25	Aflt
MA000084	0.30	Aflt
MA000084	0.26	Aflt

The dataset consists of energy consumption data from smart meters of households in London. A total of 5,567 house smart meters were used from November 2011 to February 2014. Preprocessing the data set included several steps that are standard to any project. Data cleaning, Data reduction, Normalising the data, and splitting it into test and train datasets. Data Cleaning involves removing null values and outliers. Further, Data reduction involves irrelevant details and then Normalising the data to make all values range between 0 and 1. Following this, the data needs to be split into 2 parts, test and train sets.

There are several softwares for big data pre-processing like Apache Spark and Hadoop. We choose Apache Spark for the pre-processing task due to its better performance for real-time data. The data collected from a Smart Grid is very massive and using Apache Spark is ideal. The data is generated and run through an ETL pipeline in real-time. The ability of Spark to use Resilient Distributed Datasets helps it achieve parallel processing.

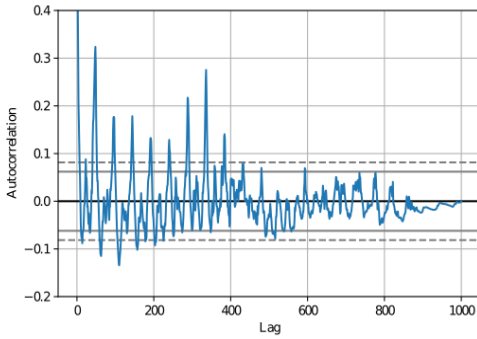


FIGURE 1: Auto-correlation on 500 hours data.

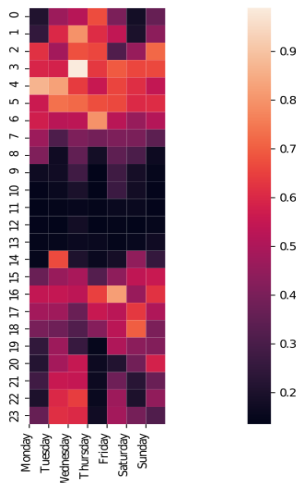


FIGURE 2: Heatmap of average hourly usage by the day.

B. MACHINE LEARNING MODELS USED

We have used a Transformer Model with an encoder-decoder structure. The data is input into the transformer model with 0s instead of missing values. This model gives us the predicted values replaced instead of the 0s for the unknown values. We have used three machine learning models: Random Forest, CatBoost and LightGBM on our dataset. These models are trained to use 12 hours of data and these predict data for half an hour.

Random Forest is a technique that is based on algorithms of decision trees. It can handle very large datasets with multiple dimensions. It can provide feature importance and is robust to overfitting. However, some of its limitations include larger

time that it takes to large time to train but we have no such limitation so we can make efficient use of it.

CatBoost is a gradient boosting on decision trees. It is designed to handle categorical data and famous for its efficient prototyping. It can be resource intensive with large datasets.

LightGBM is also a tree base learning algorithm. It uses a gradient boosting framework. It is very efficient for large datasets and handles missing values and categorical features. However like the other two models, LightGBM's complex models are very difficult to interpret

II. BACKGROUND

We have used three machine-learning techniques: Random Forest, CatBoost and LightGBM. Random Forest: It is an Ensemble learning model, which means it combines predictions of multiple machine learning models to make a much more accurate prediction that is better than the individual models, which are the decision trees.

A. MECHANISMS

1) Random Forest

The process of the random forest technique involves bootstrapping to the sample. Multiple subsets of the original dataset are created. This means that each subset may have repeated samples and might be missing some samples from the original dataset. Then the decision trees are constructed. When splitting a node during the construction of the tree, a Random Forest just considers a random subset of variables at each step. This randomness helps in reducing the correlation between them, making the trees more diverse, and therefore improving the overall model's accuracy and robustness. Each tree in the forest outputs a prediction, and the prediction with the most votes becomes the prediction of the model.

2) CatBoost

Short for "Categorical Boosting," is a machine learning algorithm used for classification and regression tasks. It is efficient for datasets with categorical features.

It works on the Gradient Boosting Framework. Boosting is an ensemble technique where new models are added sequentially to correct the errors made by existing models. Models are added until no significant improvements can be made. Gradient boosting involves building decision trees where each new tree helps to correct errors made by the previous tree. Each tree learns from a modified version of the dataset, which adjusts the weights of instances based on the previous tree's errors. CatBoost can naturally handle categorical features without the need for extensive preprocessing. CatBoost builds balanced trees, also known as symmetric trees, as opposed to depth-wise or leaf-wise trees used in other boosting methods. In symmetric trees, all leaves at the same level have the same depth, which makes the model more efficient in terms of memory and prediction speed. The value of each leaf in the decision trees is calculated using a gradient-based approach. This involves using the derivatives of the loss function to minimize the model's error.

3) LightGBM

Short for “Light Gradient Boosting Machine,” It is an efficient and scalable implementation of the gradient Boosting framework.

LightGBM creates a tree one at a time, where each new tree corrects the errors made by previously trained trees. It uses histogram-based algorithms to find the best split to grow trees. This significantly reduces the number of split points to consider, leading to faster computation and lower memory usage. Unlike many boosting frameworks that grow trees level-wise or depth-wise, LightGBM grows trees leaf-wise. LightGBM is optimized for efficiency and can handle large datasets that don’t fit into memory. It has efficient handling of categorical features. Unlike algorithms that require pre-processing (like one-hot encoding), LightGBM can handle categorical features natively. LightGBM includes L1 and L2 regularization techniques to prevent overfitting, which are typical in gradient boosting algorithms.

4) Transformers

Transformers are a type of neural network architecture that has become increasingly popular in recent years due to their success in natural language processing and computer vision tasks. The self-attention mechanism, upon which the Transformer architecture is built, enables the model to selectively focus on various segments of the input sequence during its processing. The encoder and decoder that make up the Transformer each have several layers of feed-forward and self-attentional neural networks. The output sequence is produced by the decoder, whereas the encoder processes the input sequence. Transformers have been applied to forecasting, anomaly detection, and classification, among other time-series analysis tasks.

The input sequence for a time-series analysis is made up of sequentially recorded samples, observations, or features that are arranged in a certain order. In many real-world applications, where data is recorded over a fixed sampling interval, time-series datasets naturally occur. Transformers have been shown to be effective in handling time-series data by incorporating positional encoding and multi-head attention mechanisms. The multi-head attention mechanism allows the model to attend to different parts of the input sequence simultaneously, while the self-attention mechanism allows the model to capture long-term dependencies in the input sequence..

B. RECALLING AND NORMALIZATION

Rescaling and normalization are two fundamental preprocessing techniques in data preparation for machine learning models.

Rescaling adjusts the range of data to a specific scale, often [0, 1] or [-1, 1]. The most common method of rescaling is min-max normalization. We have used min-max normalization for our project. The formula for Min-max normalization is:

$$X' = \frac{X - \min(X)}{\max(X) - \min(X)} \quad (1)$$

C. CROSS VALIDATION

Cross-validation is a method used in statistics to estimate the performance of any machine learning . It is primarily used to assess how the results of a statistical analysis will generalize to an independent dataset. This technique is essential in avoiding problems like overfitting and underfitting, ensuring that the model is robust and performs well on unseen data.

We have not used cross-validation because we used a small subset of the entire dataset due to time and resource constraints and using cross-validation would dilute the dataset and could result in samples with higher bias.

D. DIMENSIONALITY REDUCTION

Dimensionality Reduction is the process of obtaining a set of principal available by reducing the number of random variables under consideration. It transforms high-dimensional data into a lower-dimensional space, losing as little information as possible. Feature selection is the process of selecting a subset of relevant variables and predictors for use in developing a machine-learning model. It involves identifying those features that are most relevant to the prediction variable or output. Since the task is a univariate time series, there is no use for univariate analysis, dimensionality reduction, or feature selection for this task.

E. QUANTITATIVE MEASURES

1) Root Mean Squared Error (RMSE)

RMSE is the square root of the average of the squared differences between the predicted and actual values.

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

2) Mean Absolute Error (MAE)

MAE measures the average magnitude of the errors in a set of predictions. It does not consider the direction of the errors..

$$MAE = \frac{1}{n} \sum_{i=1}^n |y_i - \hat{y}_i| \quad (3)$$

3) Median Absolute Error(MedAE)

MedAE is the median of the differences between the predicted and actual values taken absolutely.

$$MedAE = \text{median}(|y_1 - \hat{y}_1|, |y_2 - \hat{y}_2|, \dots, |y_n - \hat{y}_n|) \quad (4)$$

4) Coefficient of Determination

R^2 , also known as the coefficient of determination, in a regression model, measures the proportion of the variance of the dependent variables that is explained by the independent variables.

$$R^2 = 1 - \frac{\sum_{i=1}^n (y_i - \hat{y}_i)^2}{\sum_{i=1}^n (y_i - \bar{y})^2} \quad (5)$$

III. DATA PREPARATION

In our study, smart meter energy consumption data from London households between 2011 and 2014 was used. The data underwent several pre-processing stages, including cleaning by removing outliers and null values, reduction by removing redundant details (LCLid, which is an ID assigned to each individual household, and Acorn, which is a categorical column representing the socio-economic status of the house, were dropped) and normalization which is done by making the values range with 0 to 1. When the preprocessing has been performed, the dataset has been divided into two parts which are the train and test set. The model is trained using the train set and then tested to compare the outputs using the test set. The dataset has been preprocessed on Apache Spark, which is very efficient for managing huge datasets which is a popular big data processing framework that can handle large datasets and perform distributed computing.

It is worth noting that Apache Spark is a scalable and efficient big data processing framework that can handle large datasets and perform distributed computing. It is widely used in various domains, including finance, healthcare, and transportation, to name a few. Apache Spark provides a unified analytics engine that can process data in batch, real-time, and interactive modes. It also supports various programming languages, including Python, Java, and Scala, making it accessible to a wide range of users. Apache Spark's distributed computing capabilities enable it to process large datasets in parallel, making it an ideal choice for big data processing tasks.

Once the data was processed and split into train and test sets it was converted into sequences in the format $x = (x_1, x_2, x_3, \dots, x_n)$ and $y = (y_m)$ of length $n = 24$ for input data and $m = 1$ for output data. In the context of the models the sequence x represents data from the $t - 12^{th}$ hour to the t^{th} hour and y represents data for the $t + 1^{th}$ hour where t is the current time.

IV. EXPERIMENTS

Our entire code system is based on Google Colab. We use it for a seamless multi-user experience. Its ability to connect Google Drive to the notebook is very helpful in making our task of handling huge data sets easy. We perform preprocessing using Apache Spark.

A. IMPLEMENTATION

Google Colab was used to implement the project. There are a total of six jupyter notebooks. The exploratory data analysis has been carried out in the EDA code file. The preprocessing has been carried out in the PySpark preprocessing file. There are three files for the transformer code. The other three models have been tested in the ML model training code file

1) Preprocessing

We perform the data processing and analysis using PySpark. We begin with a setup phase, where OpenJDK 8 and Spark 3.0.2 are installed, and environment variables are configured

for Java and Spark. This is followed by the initialization of findspark and the import of various PySpark modules, setting up a robust environment for distributed data processing. Then we load the dataset on London's household energy consumption, and utilize PySpark's DataFrame API for data manipulation. Key preprocessing steps include filtering the dataset by specific date ranges and 'Acorn-grouped' categories, and refining by unique household identifiers ('LCLid'). This pre-processing phase also involves data cleaning techniques such as handling missing values and dropping irrelevant columns, ensuring data quality for subsequent analysis.

In the analysis phase, we employ PySpark's MLlib for feature engineering, using VectorAssembler for feature vectorization and MinMaxScaler for feature scaling. The data is continuously filtered to focus on smaller subsets. Using Matplotlib and Seaborn libraries we generate heatmaps, providing visual insights into energy consumption patterns.

2) Transformer Model

We have developed a custom transformer. It is a standard transformer model with custom parameters. To construct the transformer model, we load adam optimizer and MSE loss function. We use this model with the preprocessed dataset.

For developing the model class, we define the MultiHead-SelfAttention class, a key element of the Transformer's architecture. This class enables the model to compute self-attention, which is critical for allowing the Transformer to focus on different parts of the input sequence. Following this, we make the TransformerBlock class, which creates a complete block of the Transformer model. This class integrates the attention mechanism with feed-forward networks, which is important for effectively processing sequential data.

Additionally, we implement the TokenAndPositionEmbedding class, which is responsible for embedding tokens and adding positional information. Then, we compile the model.

Finally, we test the model to find its accuracy, and the results are further discussed in the next section, which show that the modified transformer model is actually improvement over the other three tested models.

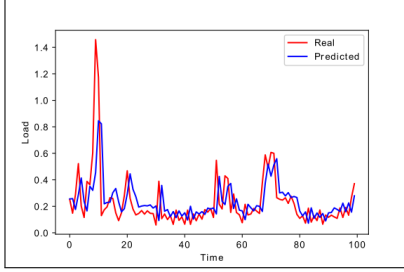
3) Random Forest

This code snippet describes the process of training and evaluating a regression model using the RandomForestRegressor from the sklearn.ensemble module. Here's a detailed explanation structured into a paragraph:

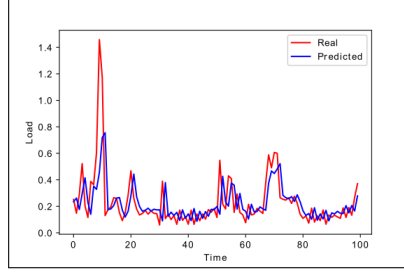
The code starts by creating a RandomForestRegressor model with a specified random-state of 0 to ensure the results are reproducible. This model is then trained using the training dataset (X-train, y-train). Following training, the model makes predictions on the test dataset (X-test). These predictions are then reshaped to match the expected format and are inverse transformed using a scaler, likely to convert them back to their original scale. Then using the sklearn package, we calculate the 4 quantifying measures.

TABLE 2: Various quantitative accuracy measures on the test set.

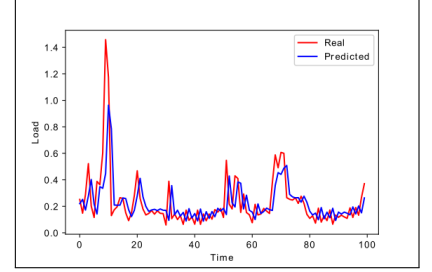
Model	RMSE	Mean AE	Median AE	R^2
Transformer	0.220	0.109	0.042	0.618
Light Gradient Boosting Model	0.224	0.106	0.037	0.604
Cat Boost	0.226	0.107	0.039	0.597
Random Forest	0.227	0.114	0.044	0.595



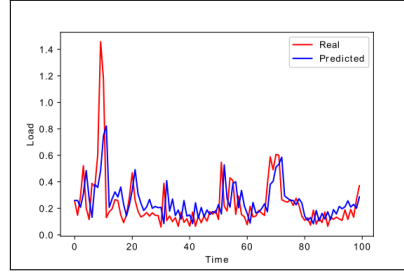
(a) Transformer



(b) Light Gradient Boosting Model



(c) Cat Boost



(d) RandomForest

FIGURE 3: Predicted values vs real values on a sample of the test set.

4) CatBoost

The code snippet utilizes the CatBoostRegressor from the CatBoost library to create and train a regression model. The model is initialized with specific hyperparameters: it's set to run for a maximum of 5000 iterations (trees), with a learning rate of 0.03 to control the speed of learning. The depth of each tree is capped at 3, and a random-state of 100 is used to ensure reproducibility of results. The model training is performed on the training dataset (X-train, y-train), and it evaluates its performance on a validation set (X-val, y-val). An important aspect of training is the use of early-stopping-rounds=1, which helps in preventing overfitting by stopping the training if the validation metric does not improve after one round, ensuring the model doesn't learn the noise in the training data.

5) LightGBM

Using the LGBMRegressor, we use this model. The settings for this model are set to as follows. Max-depth is set to 20, random state is 0, n-estimators is 1000, num-leaves is set to 25, and min-child-sample is set to 50. The evaluation matrix is rmse. Then using the sklearn package, we calculate the 4

quantifying measures.

B. RESULTS

After running all models, it is evident that the modified transformer has the best performance in comparison to Light gradient boosting, CatBoost and RandomBoost models. As shown in the Table 2, the Transformer model has the best coefficient of determination, which means that it explains the prediction the best. It also performed best in terms of RSME. But for MeanAE and MedAE, it has the lowest values. It can be concluded that this model is very efficient, and using the Apache Spark, it is very accurate. Apache Spark is very fast when compared to other frameworks used for big data analysis, which makes it ideal for our experiment. We have used univariate data, but further this model can be modified to make it a multivariate time-series which will be able to handle complex problems.

V. ETHICAL REVIEW

It is important to ensure that the data is collected, stored and used in an ethical manner. Due to this there arise several questions that one may need to evaluate. It is important to

analyse the ethical aspect of any project. We try to answer some of the important social and ethical review of this project. **Data source and Credibility:** The data has been taken from a renowned and credible source. This dataset was provided by UK Power Networks. There are no issues with the credibility and reliability of the dataset.

Data Bias: The data consists of energy consumption of 5,567 house holds from London. The model is targetted for household smart meters and represents the people it will serve. The data might be of one specific city but the energy consumption trends are typically the same across different cities. The data is collected by smart meters and has no human intervention. This means the model accuracy is not affected by the overlooks of certain populations.

Data Privacy: The data does not contain any personal information and is completely anonymous. The smart meters collect associate the consumption to a personal identity but for our analysis, the personal details are of no use.

Transparency and Accountability: The smart meters collect data in a complete transparent manner. The machine learning process is also completely transparent. However when deployed by a corporation, this might differ. The corporations usually government will be held accountable for the outcome predictions.

Fairness and Equity: The deployment of the model does not result in unfair or unequal treatment. Although the model predicts the load and the system will accordingly modify the electricity distribution. This is to create a better efficient system and is not unethical manner. This model will help increase the efficiency of the grid hence reducing the cost for the consumers as well.

Social Impact: There are no significant negative social concerns regarding this model.

Environmental Impact: The most significant positive aspect of the project is that it will assist in efficient energy management that will in turn help in reducing overall energy wasted because of the inefficient distribution system.

Accessibility: Yes, the web tool designed in an accessible manner, following universal design principles.

Consent and Opt-In: Yes, customer consent is required but the data collected is already essential for charging the customers according to their energy consumption.

Long-term Effects: The implementation of the model will help in reducing the overall carbon footprint. This will reduce the overall impact on environment due to the overproduction and unnecessary resource usage in electricity generation.

VI. WEB PORTAL USAGE GUIDE

The model has been hosted on GitHub Pages at the link <https://geareab.github.io/TimeSeriesGardio> where you can test the model live. We have integrated Gradio and HuggingFace for our deployment. Gradio provides seamless integration and allows us to make interactive web-based tools which are easy to use. Huggingfaces is a very useful tool to deploy the python gradio scripts. We embedded the Huggingface platform on a static site which then we further hosted on

GitHub pages. As GitHub pages aren't all dynamic sites, we simply embedded the Hugging faces panel in the Index.html file.

It is important to note that We have implemented only two models on the panel, LightGbm and CatBoost. Random forest and Transformer Models are too large to fit into an executable file. For such implementation, we will require Git LFS which is a paid service. As we are using only free services, the implementation of these models was not possible.

When you open the web application, there are 2 possible inputs.

If you want to enter custom values, you can use the second and third input values. In the sequence input text box you need to input a python list consisting of 24 float values. In the true output text box you need to enter the correct output value. If these are left blank then it will use the sequence from the test dataset, for you can use a slider to select the sequence from the test data set.

For the outputs, The first output is the sequence itself that represents the load during the 12 hours. The second and third outputs represent the catBoost and lightGBM predicted outputs respectively. The fourth output is the actual value itself. This can be used to compare the output results.

This way one can conveniently check the outputs of our model without the need to actually deploy the code themselves. However if one wishes to use the Random Forest or the Transformer model, they can use the code available in the project to check the outputs. These models were very large to be able to deploy them on free-to-run web applications.

...