

Exploratory Data Analysis of Global Power Plants using Various Machine Learning Algorithms

Sheikh Adil Habib¹, and Dharminder Kumar²

¹ M. Tech Scholar, Department of Electrical Engineering, RIMT University, Mandi Gobingarh, Punjab, India

² Professor, Department of Electrical Engineering, RIMT University, Mandi Gobingarh, Punjab, India

Correspondence should be addressed to Sheikh Adil Habib; sneharajput9393@gmail.com

Copyright © 2022 Made Sheikh Adil Habib et al. This is an open-access article distributed under the Creative Commons Attribution License, which permits unrestricted use, distribution, and reproduction in any medium, provided the original work is properly cited.

ABSTRACT- Nuclear plants' rewards and prices, etc and severe negative costs, are determined by their technology and the amount of electricity they create. Most nations, especially emerging ones where electricity output is expected to grow significantly, do not disclose plant-level generating statistics. The Global Power Plant Database uses this technical information to estimate the yearly energy generation of power plants. For several forms of fuels, including airflow, renewables, freshwater (hydro), as well as gas power generation, we employ different estimating models. Statistical regression and machine learning techniques are used in the process. Predictive factors include foliar data like as seed size and fuel type, as well as state characteristics also including total GDP per megawatt of installed capacity. We indicate that fossil modelling would provide more high accuracy for wind, renewable power, and hydropower is produced. Natural gas plant estimates are also improving, although the margin of error remains considerable, especially for smaller facilities.

KEYWORDS- Global Power plants, Machine learning, STLF, MAPE.

I. INTRODUCTION

Modern civilisation runs on electricity. Despite technological advances, info on physical power production by reactors is often kept proprietary by factory and service operators, difficulty in obtaining, to get by others. The Global Power Plant Database (GPPD) was produced by the World Resources Institute (WRI) and its as a fully accessible, open-access register of said west's electricity production with collaborations[1]. It covers data upon gas turbines major components such as capability (gigawatt hours), position, and diesel, which was assembled from numerous of public records. The data containing the stations' electric power (megawatts (mw) when it is publicly available. For 33 jurisdictions, verified methods of actual tree productivity are known.as of June 2019. The maximum electric power rate of a facility is described by its power plant capacity, which is commonly defined in megawatts (MW).[2] A 100 MW plant will create 100 megawatt-hours of power if it works at full capacity for one hour. In other words, capacity refers to the plant's size and potential production rate, whereas generation refers to the plant's actual power output over time. Thermal power plants generate energy using a variety of inputs, as well as the fuels used to start and run mills and the moisture used

to cool them. (Warm) waters released to a watercourse and water vapour that dries are common by-products. as well as contaminants to the air, water, and soil. Energy planners may utilize past plant generating data to track emissions and determine the best way to fulfil changing energy demand over time. The frequency and intensity with which a power plant operates vary by plant type.[3]

Annual power plant generation can be calculated using statistical models or approaches based on electrical grid optimization (also known as optimum dispatch).

Optimal dispatch models may be computationally demanding and require to create high findings, facts on performance and efficiency is required, which may vary based on operational strain and growth stage.[4] Information on efficiency is presently unavailable worldwide. Statistical models are used to evaluate the relationship between yearly unit and production variables utilizing data from power companies with stated yearly power, such as output, diesel engines, and inauguration year. These estimated correlates are then used to the characteristics of vegetation with no known population to determine their yearly multiplication.[5] Than a conveyance models, a methodology forecasts asserts on how closely an electricity resembles plants those have registered generated., rather than a system optimization. Statistical models combined with machine learning approaches are employed in our approach to estimate yearly plant generation as precisely as feasible.

Correlations between generation, technological qualities, and system factors are captured by machine learning techniques. Ummel (2012) took a unit-level approach to the problem.[5] We estimate plant-level generation since generation data is frequently given at a worldwide level.

Only a few governments make annual statistics on power plant generation public. Even when data is published, it is not always in a uniform format. We've gathered reliable publicly available data over the years. Based upon the number of sites and thermodynamic efficiency for which we have annual producing statistics for 2016, the most recent year for which we obtain data. in which the generation estimating study and modelling were conducted. Nearly half of the plants have generation statistics; however, the majority are cantered in the United States and other affluent nations.[6]. Global power plant era is depicted in Table 1.

Table 1: Reported In 2016, era was broken down by geographical region

	PLANTS IN GLOBAL POWER PLANT DATABASE	PLANTS WITH REPORTED GENERATION DATA	PLANTS WITH REPORTED GENERATION DATA (%)	CAPACITY WITH REPORTED GENERATION DATA (%)	SOURCE
UNITED STATES	8644	7,944	91.9	97.5	EIA
INDIA	8611	427	49.6	93.2	CEA
AUSTRALIA	4299	248	57.8	69.1	NGER
EUROPEAN UNION	9846	679	6.9	39.4	ENTSO-E & JRC-PPDP
OTHER*	10,304	272	2.6	1.0	Multiple sources
TOTAL	30,084	9,570	31.9	--	--

New plant kinds produce in diverse situations. Because nuclear, coal, and even certain nuclear plants have good efficiency once up and running, they normally run continuously, or at generating units, but typically require funds and efforts to set on, shut off, or modify working levels. Maintenance crews are able to minimize the incidence of runaway and ramp-down events as a consequence. Heating systems with a multi or reached its peak rating have shorter start-up and lock times and are utilized once usage surges. [7] Seedlings that rely on green energy sources such as the sun or wind generate only when they are available. Daily generating position of prominence in annual total generation that varies each marijuana plant, which also is commonly estimated and expressed by the load factor. The volume quality is a fundamental consideration.[8]

Equation 1 $effcy = tgfcy / (tcfcy \cdot \# \text{ hours in a year})$
 tgfcy is the full capacity of all fuel f reactors through producing a good in kw the year y; tcfcy would be the feature based of all fuel f units in producing a good in megapixels for year y. The thermal efficiency of a nuclear plant is mostly among 0 and 1. A score of 0 signifies that the factory was not operational during the year, whereas a coefficient of zero means that the factory was fully operational across the year.[9] Factories cost money all year and so go offline on a regular basis, hence yearly capital costs never touch 1. The most realistic annual energy capacity are around 90 percent.

A. Analysis of Wind, Solar, Hydro, and Natural Gas

To increase precision, we divide the wind turbines by fuel type. Plant type depends on the kind of plant. We focus on estimating generating for facility fed by wind and solar, hydro, and gas and oil, when data is scarce, but much knowledge is known to create quantitative estimations. We don't talk about nuclear energy. Because the International Atomic Energy Agency publishes data on all nuclear power stations. We additionally remove coal plants from our study since estimates of coal-plant generation are sought by other specialized initiatives, such as Gray et al. (2018) released by Carbon Tracker.[10] We don't have enough data on other sorts of plants: Oil plants account for just Biomass, trash, lava, wave, etc tidal generators

account for 4.7 percent of GPPD generation. all have a modest number of plants, even though these fuel sources may account for a considerable share of generation in particular areas. Statistical analysis would not produce correct conclusions with such scant data, and it would be difficult to quantify and express mistakes for such fuels.. As a result, we use typical capacity factors for those fuel sources to impute generation for these sorts of facilities.[11] It is acceptable to believe that a power plant's generation is unaffected In some cases, the generating rates from other hydroelectric forms in the system. Solar and wind energy, for examples, are only generated when the sun comes and whenever the air is blown, and at very zero net costs. Emission is possible. Restrictions, these facilities are included in the generating mix whenever they are available. This assumption is less accurate for thermal power plants, as the system operator would often dispatch units that are the cheapest among all available plants. The generation of a plant will thus be determined not only by its own qualities and however, the features and prices of competing plants, as well as total power consumption, have a role.[12]

B. Baseline Model for Most Countries

We have total yearly generating data broken down by fuel type. We can develop a consistent measure of yearly The European Environmental Agency (IRENA), which provides both institutional framework and governmental supply for renewable facilities, was used to calculate the throughput. For non - renewable energy sources in United nations Conference on trade Co-operation and Advancement (OECD) countries, the Imf (IEA) provides energy and capacity by country and fuel, as well as an estimate of productivity per kilogram placed.[13] Since declared generating by location by fuel came from the IEA, but overall volume comes by regional statistics via the GPPD, we provide a more imperfect estimation for non - renewable energy sources in non-OECD countries. This information was utilized by Byers et al. [14] in the initial recreation of the Global Power Plant.

II. OBJECTIVES

- Load data from a variety of sources has been used to encourage electric utilities in developing nations to use STLF techniques based on machine learning for more trustworthy voltage regulation
- To pick the input variables for the first times from the fresh unknown sample, analysis of data, graphical insights, and analyzed using spss version such as engine interpretation, iqr commentary, and block investigation are utilized.
- A complete predictors matrix is generated utilizing predefined transient response for sequential and customizable and non - parametric STLF equations. The predictor matrix is not very complicated theoretically, and data outside of historical trends is not required.
- This article evaluates responses were taken forecasting methods, linear and univariate product lines, using a range of analytical parameters such as MAPE, RMSE, MSE, R-square, and confidence interval. techniques. Furthermore, we performed a comprehensive seasonal study to assess and compare the effectiveness of the proposed methods.

III. LITERATURE REVIEW

Because it is reliant on a large number of multimodal and semi periodic and climatic factors, also including important festivals, dryness, and cold, the STLF is challenging and demanding. Several forecasting approaches have been developed based on substantial study on the STLF problem during the last two decades. Several regression models have been developed in [13], including ARIMA (Alarm Embedded Line Graph) and festive are two types of ARIMA (SARIMA). Exponential smoothing and SARIMA employ the latent parameter estimates of STLF data's to turn pro data into t - statistics. Auto-correlation (ACF) and fractional auto-correlation (PAC) are two types of auto-correlation (PACF) studies can be used to detect this seasonality [15]. In addition, discusses an overview of numerous different statistical regression models, as well as their variables and methodologies, single coefficient of determination and multivariate regression, for example. Is from the other hand, univariate garch techniques fail to capture temporal variations and non-linear load profile structures [14].

To improve the performances of analytical regression equations, Supervised Classification (PCA), analyses, and notch iteration can be used. to address the above-mentioned weaknesses for STLF. By using correlation analysis to determine eigenvalues of multi-variate electrical load data, PCA applies a dimension reduction approach premised on invertible approach

However, with PCA, selecting the coefficients of the covariance matrix is time consuming, and this can lead to the loss of important seasonal effects of temperature on electrical demand data . Singular Value Decomposition (SVD), on the other hand, is more robust than PCA in extracting both seasonal and random components. SVD, on the other hand, works with a complicated unitary matrix that is computationally intensive.

IV. METHODOLOGY

Our primary aim is to properly estimate yearly past plant-level generation. Models of machine learning are ideally adapted to the task (Olden and colleagues, 2008). We use certified predictive models to define the link between of different variables. e—in this example, plant-level yearly generation—and the researchers' chosen independent factors, or predictors. Given a collection of data that contains The model includes both relationship between variables. is optimized across many iterations with carefully calibrated parameters to minimize the estimation error (called the labelled training data). A huge number of algorithms may be used in machine learning models. As the model algorithm, we used the GBT vulnerable and prone (transfer learning tree). The GBT is a suitable modelling method in this case cos of the foregoing:

- Regression trees can capture nonlinear interactions (for example, there is no linear connection among both weather conditions and power). As the wind speed increases, a turbine will eventually run out of power to create.
- Tree-based models make it simple to determine which predictors contribute the most to the outcome.
- Tree-based designs foretell by identifying like instances in the test dataset, so that the guess remains inside that instructional set's target variable's area.

This method creates a succession of judgement trees again and over, seeking to explain variability by addressing the space containing plant-level features. Using a huge tress that are created systematically improves the model's classification precision, even when per tree is fairly bad on its own (Elith et al. 2008). The dependent variable in our situation might be yearly capacity factor or plant-level annual generation. As seen in Equation 1, they are inextricably linked. For two reasons, we chose capacity factor as the dependent variable. First, employing generation would place a greater For bigger plants, usually yield most use across the year, the accent is on reducing mistakes. Second, kwh values are simple to interpret. By inverting Equation 1, we can easily compute yearly generation for a plant if we have an estimate for the capacity determinants One of the risks of training sets is that they "overfit" to the test dataset, resulting in a small random variable when matched to the training sets but a large coefficient of determination when performed to new unfamiliar data. This risk is magnified if the learning algorithm do not cover all potential events, which is the case in our case because bulk of the labeled assertions in the currently GPPD are for plants. We lower the risk by splitting the United States into communities depending on the North American Electric National Corp (NERC) categorization and utilizing distinct region-level throughput, resulting in more variance in petroleum electricity production; and test the programs using non-real datasets.

The International Energy Agency (iea (EIA) of the United States produces capacity and availability data. generating data for all types of power facilities in the country. Each unit is identified by the NERC area to which it belongs. For aggregating unit-level data to regional levels, we calculate NERC-region capacity factors by fuel.

A. Data Filtering and Detection of Outliers

Although the research is based on factual sources data, inaccurately reported generation or capacity might result in an exaggerated capacity factor for a specific facility. Plants may also be unable to run for a whole year owing to maintenance issues. In this sort of research, we can't foresee extended maintenance intervals, therefore we focus on projecting generation for plants that are constantly accessible throughout the year. Outliers may have an undue impact on the model, resulting in erroneous predictions. To decrease measurement error, we first exclude situations where the electricity production is greater than or equal to 1. Ways larger than one are frequently misrepresented as generating value or potential. Photovoltaic plants that buy more water than they output, such as pumps stores, can have capacity factors that are less than zero. Anomalous reports with throughput larger than three deviations from the mean (as assessed across all countries for the specified fuel type) are also removed. Each section includes information on any extra fuel-specific data cleansing methods.

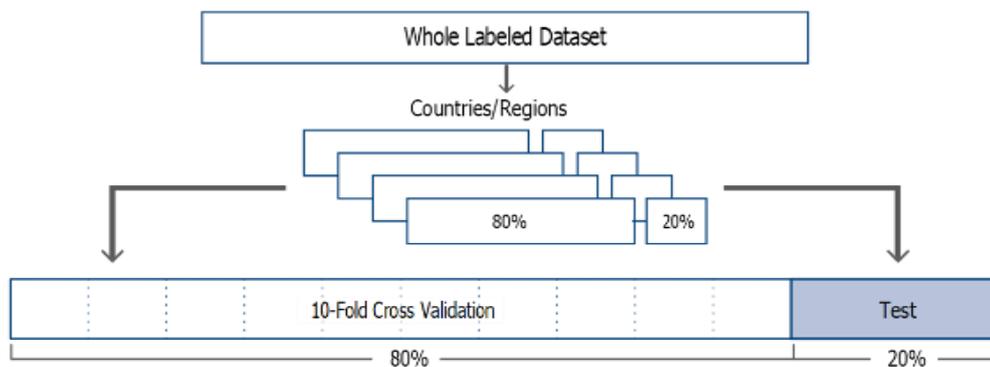
B. Evaluation and Testing of Models

Training, validity, and test sets were created from the specific dataset. The models is adjusted to match the training set., iteratively if required, depending on its

performance on the validation set. The test or unknown data is then used to assess it. The test data are 20% of the original dataset, stratified by nation to ensure that they reflect the total labeled dataset's geographical distribution. To partition We use tiered K-fold stepwise regression to divide the information into assessment and training sets. Purposeful sampling is used to solve the model using K identically sized subtests, or folds, of the training data, as shown in Figure 1 for 10 folds. The classifier is trained on (K-1) folds prior to getting analyzed, or validate, on a dataset consisting that was not used (Varian 2014). The spin of the kept-out sample population is completed. By repeating this process K times, we can achieve K confirmation scores. The identify appropriate score is the simple average of the K grades. K-fold analysis technique reduces the risk of a single training examples not being inclusive of the whole data and trying to skew the result (Shulga 2018). It also allows for constant validation for education utilizing all of the data, and that is very beneficial when data is limited, like it is in our case. To improve bend scores, we change the model. Last, we put the fine-tuned model through its paces on the test set to assess how well it works. The entire retraining, validity, and test cases is depicted in Figure 2. The process is the same for all of the fossil models.



Figure 1: Sample with K=10 Folds: Cross Validity



Source: Authors

Figure 2: Training, Validation, and Test Split

C. Model Performance Evaluation

To assess model performance, we use two measures to compare the estimated and reported capacity factors: The deviation (MAE) and mean square confidence interval (MAPE) are two types of errors (MAPE) The MAE rights in respect real and anticipated readings and analyzes the percent difference over all occurrences, providing a simple error measure, but does not measure the relative magnitude of the error. On the other token, the MAPE is unitless, making it simpler to compare accuracy across capacity factors of varying magnitudes.

V. SYSTEM ARCHITECTURE

A. Data

Global Database of Power Plants the GPPD includes 2,598 gas-fired power stations that were built before 2016, with around 70% of them reporting gen. North Asia is home to 66 percent of all coal power stations, but 93 half of those that make energy. In parallel to other fuels, gas has been used in electricity production, and many have numerous turbines with various generating technologies. We only model nuclear plants with a percent of propane capacity more than 95% (relatively to all biofuels) and a share of any operating method (subsumed as CCGT, CS, GT, ST, IC, or FC) less than 95% of complete plant for efficiency's sake. The set of data in the training set is reduced from 1,780 to 1,284 as a result of this. Gas fired power plant is shown in figure 3

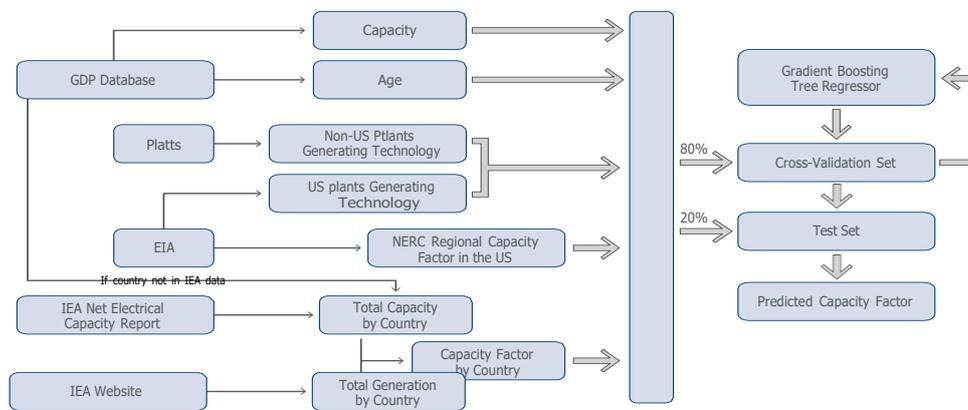


Figure 3: Gas-Fired Power Plant Generation Estimation Workflow

Table 2: Gas-Fired Plants and Plants with Reported

	NUMBER OF GAS PLANTS BY REGION	GAS PLANTS BY REGION (% OF WORLD PLANTS)	NUMBER OF GAS PLANTS WITH REPORTED GENERATION DATA (PLANTS THAT CAN BE USED FOR TRAINING)	GAS PLANTS WITH REPORTED GENERATION DATA (% OF WORLD PLANTS WITH REPORTED GENERATION)
NORTH AMERICA	1,708	65.7%	1,651	92.8%
SOUTH AMERICA	128	4.9%	0	0.0%
EUROPE	267	10.3%	73	4.1%
AFRICA	68	2.6%	1	0.0%
ASIA	425	16.4%	55	3.1%
AUSTRALIA/OCEANIA	2	0.1%	0	0.0%
TOTAL	2,598	100%	1,780	100%

Generation by Region (2016) in the GPPD A total of 186 gas plants, or 14% of the residual observational data, had thermal efficiency of less than 1%. (Equivalent to around 100 hours of high-intensity emission per year) These plants are difficult to model using our estimation methodologies. As a consequence, we were able

to delete the points, presenting us with 1,098 labeled facts. EIA-923 The EIA reports plant-level electricity and generating capacity one per unit using data obtained through "Form EIA-923." We use this material and the GPPD to designate individual station with something like

a specific generated skill or a set of generating technologies due to the structure of the facility.

Platts' World Electricity Power Plants While the EIA only collects data on generating technologies for gas plants in the United States, the WEPP database has data on gas plants all around the world. For the analysis, we retrieved pertinent data from both sources, however we do not republish the WEPP data because it is confidential.

When employing unknown datasets, The ratings on data, on the other hand, are still more accurate approximations of engine estimates quality. With an annual mean of 0.136 and an absolute range of 0.136, percentage error of more than 170 percent, the model forecasts plant generation. While still high, this is a significant improvement over the baseline, which utilizes an uncorrected average capacity factor and results in a 350% inaccuracy.

Figure 4 indicates that the suggested model exhibits more variance between plants than the baseline method, resulting in a more accurate forecast. The cluster of dots along the vertical axis represents a In 2016, there were a large variety of options available with poor performance. The theory approximates the utilization factor on low kwh plants and undervalues the utilization factor for high bandwidth frequency plants in general so the nodes are against the 45-degree line. Figure 5 shows how the defect varies with biomass production, with growers generating more commonly having a lower error. The x axis shows the size barrier over which the efficiency parameters are reassessed, as well as the capacity restriction within which installations are not included.

The typical absolute percentage error for natural gas plants never goes below 50%, implying that Our current predictors are insufficient to accurately anticipate plant-level energy. Due to the training data extraction and approach installation, we are often unable to estimate which stations will be on extended care or provide for fewer than 100 hours in a row per year. In our dataset, this relates to around 14% of plants with recorded productivity and equipment type. Tree-based models produce attempt to highlight ratings for each of the covariates shown in the classification algorithm.

These "element significance ratings," which always add up to one across all predictions, are shown in Figure 6. The throughput and duration of a plant are the comparatively more relevant determinants of its load demand in this scenario.

Solar photovoltaic

B. Solar Photovoltaic

Solar photovoltaic plants are the topic of this section. Solar thermal plants account for just They account for 1.6 percent of the entire generating capacity and thus not included in this study. Model Description (6.1) The basic pieces of a pv systems (PV) system that is linked to the grid are depicted in Figure 11. The photovoltaic effect converts incoming solar energy into electricity in PV modules. After energy is generated, components of the balance-of-system (BoS) assist in its regulation to ensure grid-quality output. The inverter, which converts power One of the significant BoS sections is the conversion of dc generator (DC) to oscillating (AC), which allows it to be provided into the

power system. When things like these happen, it's important to be prepared. cloud cover occurs, the PV system may contain a battery or other storage option, allowing it Than as a rapidly fluctuating output, provide a steady output. One of the key parts in the energy input output per MW of megawatts is the amount of irradiation intercepted more by solar panel, that may be approximated using the seasonal actual average diagonal intensity (GHI) somewhere at surface level. Temperature also plays a factor, since rising temperatures reduce solar panel efficiency (Dubey et al. 2013)..

The interpolation is two-dimensional only, with no consideration for build - up or alpine relief The annual combined cycle for each country or area incorporates fluctuations in a variety of inputs that we can't directly measure, such as the curtailment in a country or region owing to Limits in communication, transport, or organization Ways were calculated using IRENA and the EIA for each county and NERC area. in the United States, respectively. To conclude, the capacity factor is the dependent variable in the solar model, whereas the independent variables are:

- year global mean horizontal irradiance at earth's surface; yearly average operating temperature at the solar farm site; plant age;
- power of the plant; and
- by state solar rated capacity

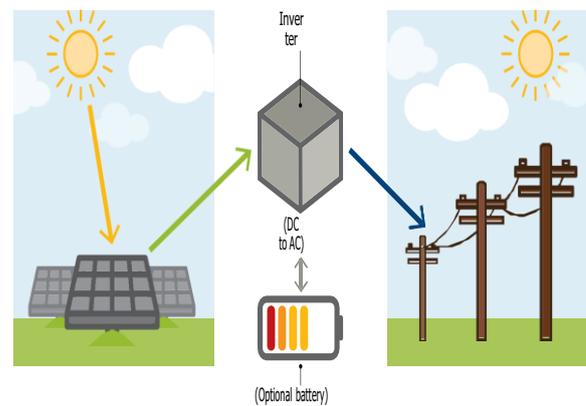


Figure 4: Simplified Grid-Connected Photovoltaic System

The computer findings are so much more appropriate than the default estimations, as seen in figure 4 . based on the holdout set of 263 observations. Based on the test set, a power plant's capacity factor is predicted to be within 15% of its genuine capacity factor on average. The observations The factors on the given in figure 5 clump and around closed interval, suggesting that they are useful in supporting within-region solar farm variation. As observed on the perfect side of the frame, the residuals to every item in the test set are equally and quick and easy way, showing that the framework is sturdy enough to produce consistent judgments.

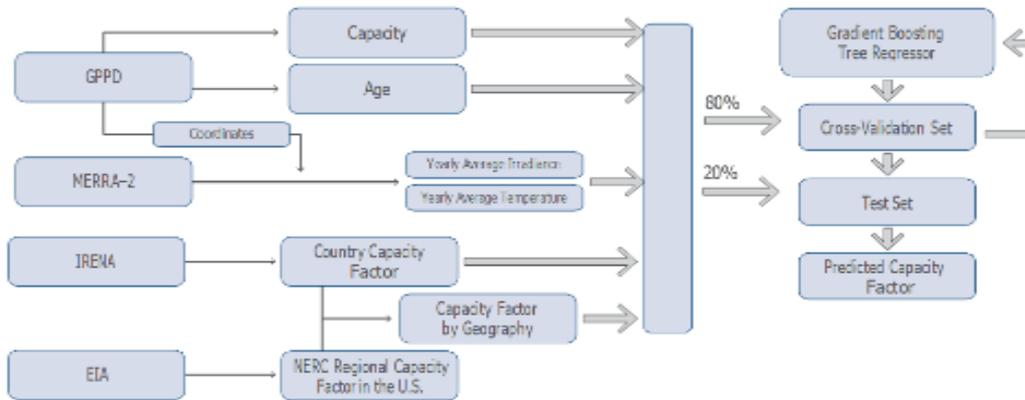


Figure 5: Workflow of solar PV

C. Hydropower Plant

Large hydro dams, tiny mill systems, and battery storage are all examples of hydroelectric plants. Pumped storage hydropower (PSH) facilities are not included in this research since they are a storage rather than a generating innovation as shown in figure 6. Drawing water to a water tank consumes more energy than the amount generated is when water returns downstream via the mill in PSH facilities, resulting in a negative yearly supply.

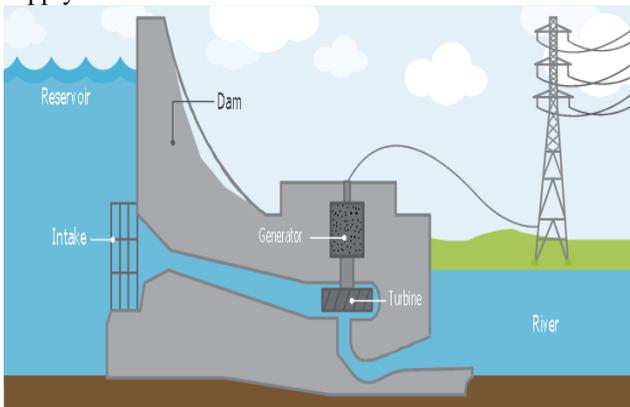


Figure 6: Simplified Hydro powerplant

The stored power can be used to fulfill Quickly meet a utility's or a region's load profile. Even though pressurized reservoirs are definitely vital parts of various energy networks, their primary function is to store power. arbitragers is difficult to be recorded with any confidence on an annual basis. Model Description Water flows through hydroelectric plants, spinning turbines that generate energy. Dams and reservoirs are used in large hydropower projects older children are mainly take

installations with a weir to catch the runoff rather than just a dam. The weir redirects a portion of the water out from principal small rivers to a turbine (Paish 2002). Rainfall and debris gathered in the watershed where it facility was located. is located are used to generate hydroelectricity. To anticipate yearly plant level generation, we use the following variables:

- Plant capacity, which defines the maximum amount of power that may be produced at any moment.
- Average runoff for the power plant site (including surface and subterranean runoff).
- River order determines the magnitude of the river that flows into the reservoir.
- Larger rivers are referred to by smaller orders. Order 1 denotes the main stem river from source to sink; order 2 denotes all tributaries flowing into a first-order river; order 3 denotes all tributaries flowing into a second-order river; and order 0 denotes conglomerates of minor coastal watersheds (Lehner and Grill 2013).
- Annual average capacity factor by nation, which includes statistics from other countries or regions that we don't directly see or measure. Hydroelectric power plant operations are often governed by full-system and regulatory requirements, which include environmental limits (Niu and Insley 2013)

responds more slowly (Kao et al. 2015). Because . Over the course of a year, daily needs become less significant, and runoff volume becomes a primary factor of generation (Kao et al. 2015). We split surface and subsurface runoff in the model because surface runoff responds quickly to precipitation events, whereas subsurface, or base flow, runoff water collects throughout a drainage region, just monitoring runoff at the reservoir location offers little information.

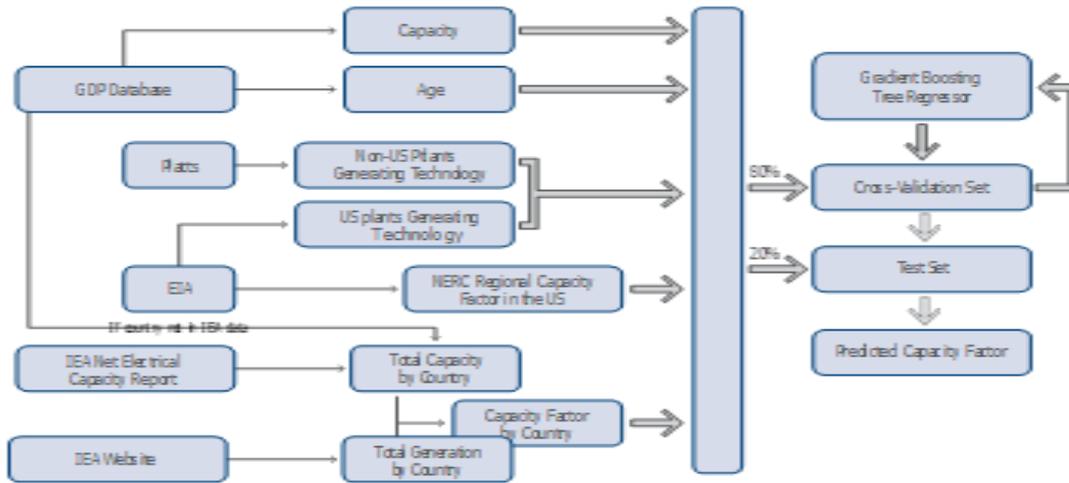


Figure 7: Workflow of Hydro powerplant

Following Kao et al., we total observations over the relevant drainage region as a predictor of generation (2015). The drainage area is measured and pertinent meteorological data within the area is aggregated using ERA5 global climate data and HydroBASINS. At a worldwide scale, Hydro BASINS is a set of polygon layers that displays watershed and sub-basin borders. Each polygon has its own ID and is polygon where a certain plant is located (figure 7) using the plant coordinates, and then retrace all upstream polygons to identify the whole drainage basin.. The arrows that cross polygons indicate the real movement and buildup of water. Hydro BASINS determines the upstream polygon for each polygon. We discover 3 by searching up the upstream polygon of 2 We'll keep searching until we find a polygon without any upstream polygons (polygon 5 in this case). The ensemble of all polygons obtained in this procedure is then used to determine the drainage area of this plant.

VI. SIMULATION AND RESULTS

Power plants have the ability to create a specific quantity of electricity over a period of time, but they are not really generating power if they are taken offline (for example, for maintenance or refueling).

A. The Capacity Factor

Energy enthusiasts can use capacity factors to assess. the reliability of various power facilities. It simply counts how many times a facility runs at maximum capacity. A plant with a capacity factor of 100 percent is always producing electricity

```
df_generation =
df.groupby('country')['estimated_generation_gwh_2020'].
sum().sort_values(ascending=False).to_frame()
px.bar(df_generation, title='Electricity Generation Per Country')
```

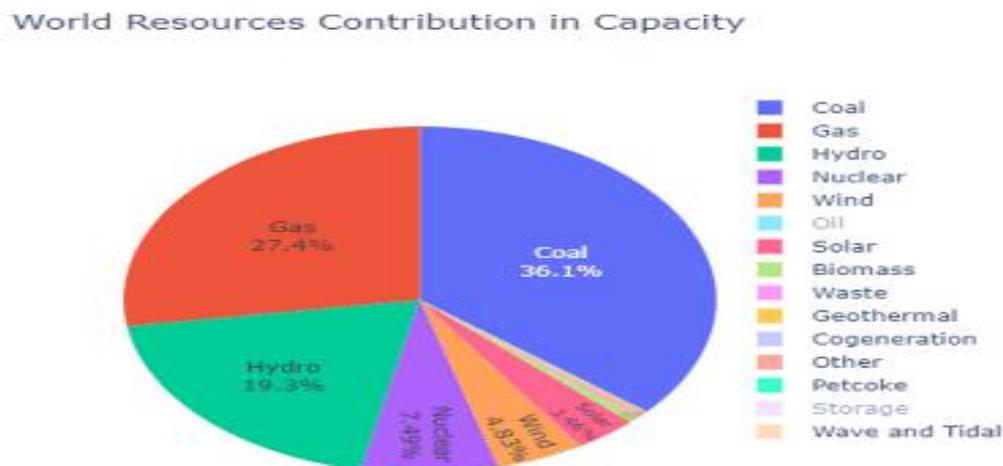


Figure 8: World Resources Contribution in capacity

World Resources Contribution in capacity is shown in figure 8

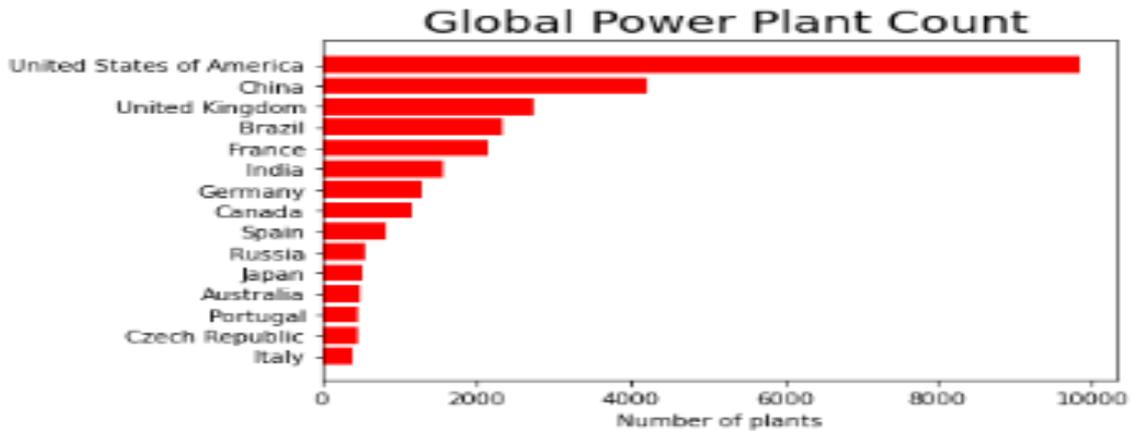


Figure 9: Global Power plant estimated count

Global Power plant estimated count is shown in figure 9

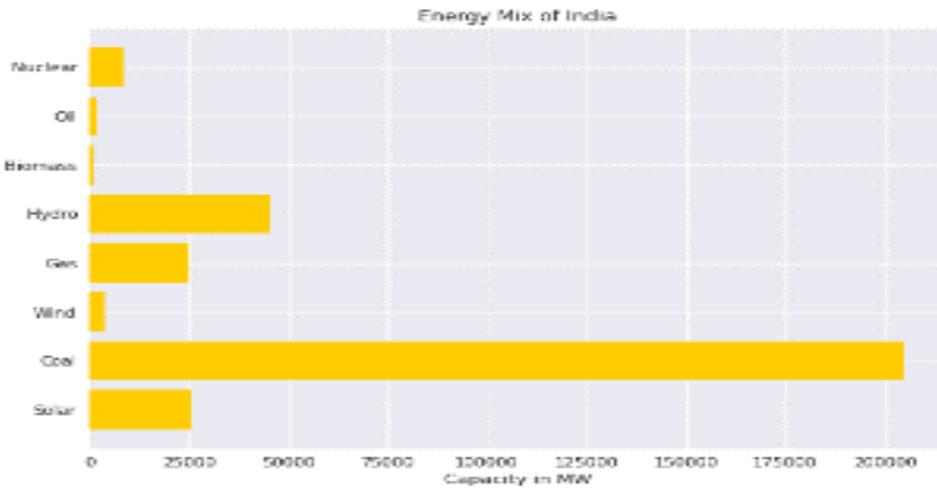


Figure 10: Energy mix of india estimated

Figure 10 shows Energy mix of india estimated.

VII. CONCLUSION

When system variables are essential, this would allow the models to reach better accuracy, but it would rely on the availability of adequate high-frequency data. Finally, when fresh training data becomes available, we continue to incorporate it. While the extra data will have no impact on our estimate methods, it will improve the prediction models. They will also assist us in the development of time series, whether by adding yearly generation estimates for various years or by producing higher frequency data. We will also update technical information for the plants as it becomes available, expanding the number of species for which the procedures outlined in this technical note may be used.

REFERENCES

[1] Dodamani, S.; Shetty, V.; Magadam, R. Short term load forecast based on time series analysis: A case study. In Proceedings of the IEEE International Conference on

Technological Advancements in Power and Energy (TAP Energy), Kollam, India, 24–26 June 2015; pp. 299–303. [Google Scholar]
 [2] Dudek, G. Short-term load forecasting using neural networks with pattern similarity-based error weights. *Energies* 2021, 14, 2334. [Google Scholar] [CrossRef]
 [3] Edigera, V.Ş.; Akarb, S. ARIMA forecasting of primary energy demand by fuel in Turkey. *Energy Policy* 2007, 35, 1701–1708. [Google Scholar] [CrossRef]
 [4] Hussain, I.; Ali, S.M.; Khan, B.; Ullah, Z.; Mehmood, C.A.; Jawad, M.; Farid, U.; Haider, A. Stochastic wind energy management model within smart grid framework: A joint bi-directional Service Level Agreement (SLA) between smart grid and wind energy district prosumers. *Renew. Energy* 2019, 134, 1017–1033. [Google Scholar] [CrossRef]
 [5] IEA India Energy Energy Outlook 2021. Available online: <https://www.iea.org/reports/india-energy-outlook-2021> (accessed on 25 August 2021).
 [6] IEA South Asia Energy Outlook 2019. Available online: <https://www.iea.org/reports/southeast-asia-energy-outlook-2019> (accessed on 25 August 2021).
 [7] Jawad, M.; Ali, S.M.; Khan, B.; Mehmood, C.A.; Farid, U.; Ullah, Z.; Usman, S.; Fayyaz, A.; Jadoon, J.; Tareen, N.; et al. Genetic algorithm-based non-linear auto-regressive with exogenous inputs neural network short-term and medium-term

- uncertainty modelling and prediction for electrical load and wind speed. *J. Eng.* 2018, 2018, 721–729. [Google Scholar] [CrossRef]
- [8] Jawad, M.; Qureshi, M.B.; Khan, M.U.S.; Ali, S.M.; Mehmood, A.; Khan, B.; Wang, X.; Khan, S.U. A robust optimization technique for energy cost minimization of cloud data centers. *IEEE Trans. Cloud Comput.* 2021, 9, 447–460. [Google Scholar] [CrossRef]
- [9] Jawad, M.; Rafique, A.; Khosa, I.; Ghous, I.; Akhtar, J.; Ali, S.M. Improving disturbance storm time index prediction using linear and nonlinear parametric models: A comprehensive analysis. *IEEE Trans. Plasma Sci.* 2019, 47, 1429–1444. [Google Scholar] [CrossRef]
- [10] Khan, K.S.; Ali, S.M.; Ullah, Z.; Sami, I.; Khan, B.; Mehmood, C.A. Statistical energy information and analysis of Pakistan economic corridor based on strengths, availabilities, and future roadmap. *IEEE Access* 2020, 8, 169701–169739. [Google Scholar] [CrossRef]
- [11] Kiprijanovska, I.; Stankoski, S.; Ilievski, I.; Jovanovski, S.; Gams, M.; Gjoreski, H. HousEEC: Day-ahead household electrical energy consumption forecasting using deep learning. *Energies* 2020, 13, 2672. [Google Scholar] [CrossRef]
- [12] Musbah, H.; El-Hawary, M. SARIMA model forecasting of short-term electrical load data augmented by fast fourier transform seasonality detection. In *Proceedings of the IEEE Canadian Conference of Electrical and Computer Engineering (CCECE)*, Edmonton, AB, Canada, 5–8 May 2019; pp. 1–4. [Google Scholar]
- [13] National Transmission and Despatch Company Limited. *Power System Statistics 45th Edition*. Available online: <https://ntdc.gov.pk/ntdc/public/uploads/services/planning/power%20system%20statistics/Power%20System%20Statistics%2045th%20Edition.pdf> (accessed on 25 August 2021).
- [14] Shah, I.; Iftikhar, H.; Ali, S.; Wang, D. Short-term electricity demand forecasting using components estimation technique. *Energies* 2019, 12, 2532. [Google Scholar] [CrossRef]
- [15] Wood, A.J.; Wollenberg, B.F.; Sheblé, G.B. *Power Generation, Operation, and Control*; John Wiley & Sons, Inc.: Hoboken, NJ, USA, 2014; Chapter 12; pp. 566–569. [Google Scholar]