

Predicting Full Power Output from a Combined Cycle Power Plant

Amin Fesharaki, Benjamin Earnest, and Jeffrey Joyner

Shiley-Marcos School of Engineering, University of San Diego

Jun 27, 2022

Abstract

Predicting full load power available from a power plant, based on assumed or forecasted environmental or ambient conditions, has strong implications for improved operational efficiencies and profit from available megawatt hours. This paper conducts a secondary data analysis, utilizing operating data from a combined cycle power plant over a six year period. There are four main parameters (predictors) that influence power plant performance. These parameters include ambient temperature, atmospheric pressure, relative humidity, and steam turbine exhaust pressure (vacuum). The data is used to build both linear and non-linear regression models that predict hourly full load electrical power from the combined cycle power plant. Model performance is measured by mean absolute error (MAE), Root Mean Squared Error (RMSE), and R squared. After analyzing the performance metrics on the test data set for each model used in this study, it was determined that Random Forest tuned by out-of-bag estimates yielded the best results. The final model used had a mtry value of 2 with an RMSE of 3.1709 and an R^2 of .965 on the training data set. However, the performance metrics on the testing data sets had a slightly higher RMSE of 3.7197, a lower R^2 of .952, and an MAE of 2.6460. The Random Forest model tuned by out-of-bag estimates slightly outperformed the most successful model in Tufekci's 2014 study while using the same feature subset. Furthermore, the prediction accuracy for the optimal machine learning model is suitable enough to replace thermodynamic approaches to model a real world system. By doing so, the time and effort to model thermodynamic systems analytics can be greatly reduced. This data was compared to a 2014 study, which found the optimal model to be a bagging algorithm with REPTree (Tufekci, 2014). The optimal model for this study slightly outperformed the 2014 study.

Problem Statement

Combined cycle power plants (CCPP) are thermodynamically complex systems. By definition, they utilize more than one means to generate power, and one of the generators will use waste heat or energy from the other generator as its energy source for electric power generation. This is advantageous due to the efficiencies gained. For our analysis, we revisit a problem first explored by Tufekci (2014), who did a predictive analysis on a combined cycle power plant made up of two gas turbines (GT) and one steam turbine (ST). Figure 1 shows a functional diagram of the combined cycle power plant under consideration. Two GT generators output approximately 160 MW of electrical power each. The exhaust from the GT generators is then run through steam generators as a heat source, creating high pressure steam that is used to power an ST generator.

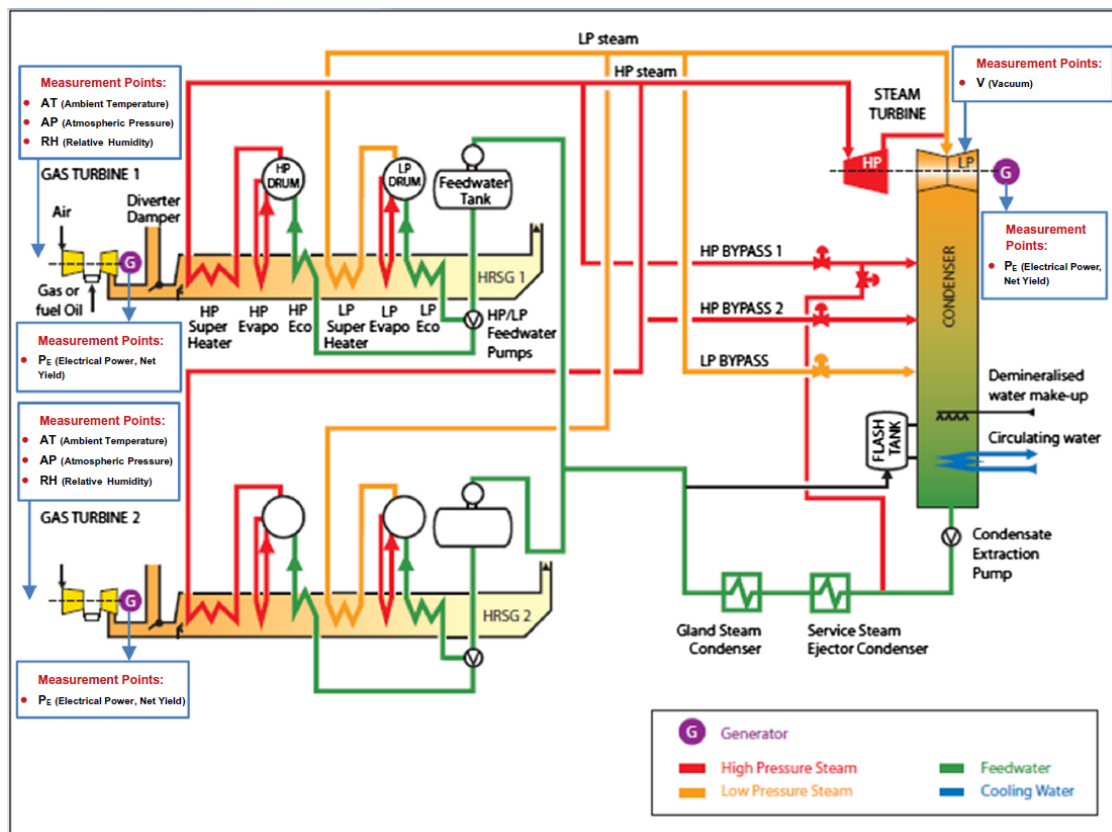


Figure 1 Combined Cycle Power Plant (Tufekci, 2014)

Predicting power plant electrical output provides benefits from both an economic and a plant efficiency perspective. The opportunity to plan for and maximize output at least one day prior, as well as income from the available mega-watt hours is a desirable tool. The problem, as Tufekci (2014) found, requires many assumptions to deal with the thermodynamic complexities that come with combined cycle power plants. Without these assumptions, the machine learning problem would be much more challenging, and the thermodynamic analysis would require thousands of non-linear equations (Tufekci, 2014). To avoid these complexities, Tufekci (2014) considered primarily, for GTs, the ambient parameters that impact performance, including ambient temperature, atmospheric pressure, and relative humidity. ST generator performance has a relationship with exhaust vacuum, which will be utilized in this analysis. Our goal for this analysis is to identify the optimal predictive model for full power output from a CCPP, then compare optimal model performance to the study done by Tufekci, and determine if the modeling tools available in 2022 provide opportunities to create more effective predictive models to solve the same problem.

Exploratory Data Analysis

This effort is a secondary data analysis. The data used was downloaded from the UCI Machine Learning Repository. It consists of 5 variables (4 predictors and 1 target) and 9568 observations. The predictor variables include ambient temperature (AT), atmospheric pressure (AP), relative humidity (RH), and steam turbine exhaust vacuum (V). The target variable is full load electrical power output (PE). The observations represent data collected from the CCPP over a six year period, each observation representing hourly measurements over the six year period. Note, each observation represents maximum output of the CCPP, based on the conditions of that day. Table 1 shows the descriptive statistics for the predictor and target variables.

Table 1 *Descriptive Statistics for Predictors and Target Variables*

Characteristic	Descriptive Statistic			
	Mean	Standard Deviation	Range	Median
Ambient Temperature (AT)	19.65	7.45	35.3	20.34
Atmospheric Pressure (AP)	1013.3	5.93	40.4	1012.9
Relative Humidity (RH)	73.31	14.60	74.6	74.97
Turbine Exhaust Pressure (V)	54.31	12.71	56.2	52.08
Electrical Power Output (PE)	454.4	17.07	75.5	451.6

Note. Units are AT(degrees C), AP (mbar), RH (%), V (cm Hg), PE (MW)

Our first step in exploring the data set was to check for any missing values that may require removal, imputation, or other actions to ensure they don't inappropriately influence or interfere with model development. We found there were no missing values. Further, we also found that the data in the predictor and target variable columns were numerical and continuous. This lends itself well to regression models.

Now that we understand the structure and quality of the data, we looked for relationships between variables, investigated the distributions, and looked for outliers that may influence model training. Figure 2 shows the pairwise scatter plots of each predictor and the target variable.

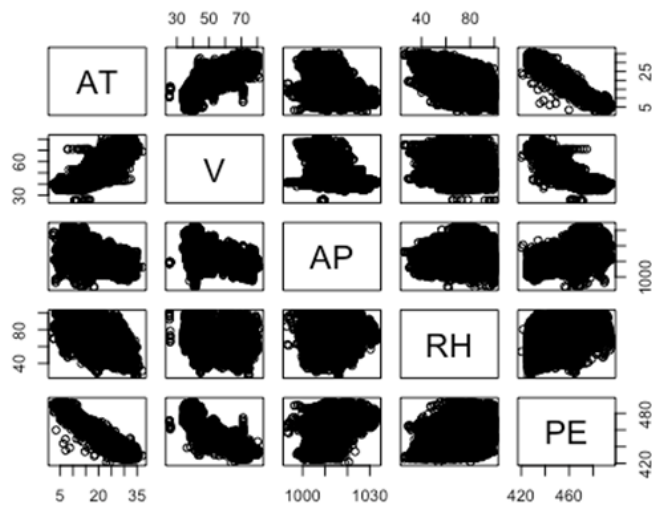


Figure 2 Pairwise scatter plots of variables

This shows linear relationships may exist between the PE and AT, as well as PE and V. We verified this by checking the correlation between all variables. The correlations were calculated to be -0.94 between PE/AT, and -0.87 between PE/V. We also checked the distribution of the variables, using histograms, and looked for outliers using box plots. Figure 3 shows the histograms of the predictor variables. They generally follow a normal distribution, and do not appear skewed. Vacuum (V) shows the least normal distribution, and will need to be addressed in preprocessing.

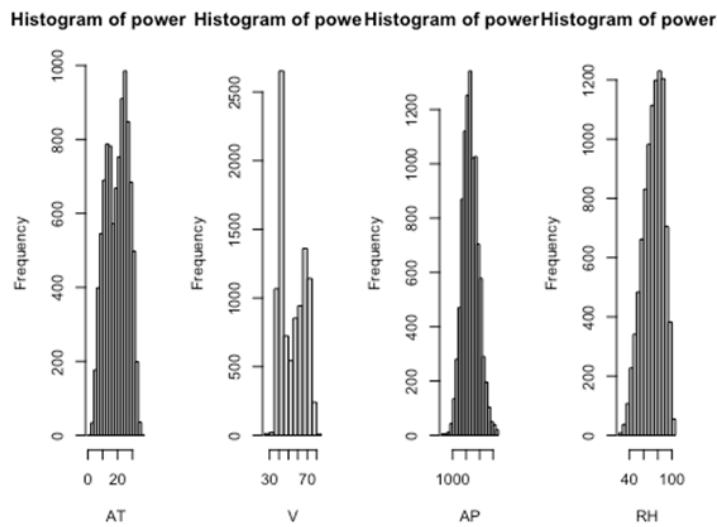


Figure 3 Distribution of predictor variables

The boxplots in figure 4 show very few outliers, but make it clear the variables are on very different scales, which will also need to be dealt with during preprocessing, before data splitting and model training.

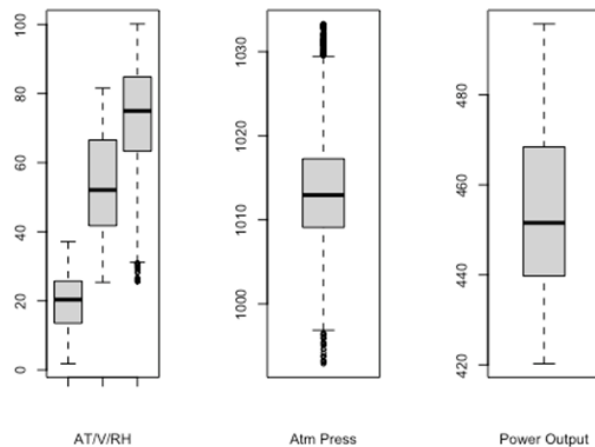


Figure 4 Box plots of all variables

Data Pre-Processing

To preprocess, we first checked for any near zero variance (degenerate) predictors that may not add value to the model training. We used the `nearZeroVar()` function in R, passed in the data set, and

found there are not any degenerate predictors. Next we checked for any highly correlated predictors, and found that AT and RH had a fairly strong correlation. So strong we could potentially remove AT as a predictor. This was confirmed by principal component (PC) analysis, which identified three PC's that account for 95% of variance in the target. However, the 2014 study included all of the predictors in the optimal model, so we also will for comparison. We will investigate variable importance later on for the optimal model to validate whether AT should stay, or recommend removing in future analysis.

The predictors were centered, scaled, and transformed using the Box Cox method, to prepare the data for model training. This is done to improve numerical stability for many of the modeling calculations. Many of the linear regression models are sensitive to the scale of predictors, while others, like tree based methods, are not. The only real downside to centering, scaling, and transforming is the loss of interpretation of the data, as it will no longer be in its original units (Kuhn & Johnson, 2013). The output can be seen in the new boxplots, which show the data each having normal distributions and on the same scale in figure 5.

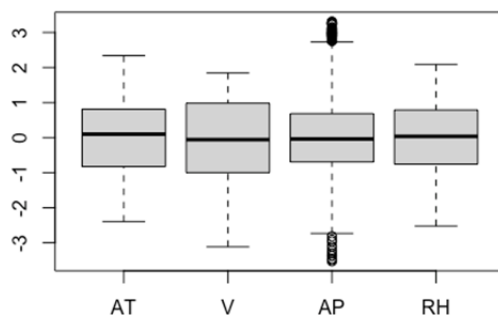


Figure 5 Box plots of centered, scaled, and transformed predictor variables

Data Splitting

The data was split into training and test data sets, first by separating the predictors into “x” or independent variables, then the target variable (PE) as the “y” or independent variable. The predictors

were centered, scaled, and transformed after partitioning into training and test data sets. We used the `createDataPartition()` function in R to create an 80/20 split of the data, where 80% of the 9568 observations were partitioned into the training set, and the remaining 20% were set aside as the test data set.

Modeling Strategies

A mix of regression modeling strategies was implemented for comparison of best model performance according to the lowest Root Mean Square Error (RMSE), Mean Absolute Error (MAE), highest R-squared values. Our goal was to produce a model that would best predict power output given our predictor variables and to see which predictor variables have the most influence on the full power output of the CCPP.

With the data now pre-processed and partitioned, the training and testing variables created within the R-language code were incorporated into the various models for computation and analysis. We estimated the level of performance of the models by analyzing how it performed on data it has not yet seen, which is our test set of data. We used 10-fold cross validation as the resampling method across all models, in which data were partitioned into a training set of data and a test set of data. A random number seed of 100 was also used across all models for reproducibility of results. The algorithms used the training set of data to train the various machine learning models and to find their parameters. Then each model's performance against the test set of data was evaluated.

We also intended to evaluate if linear or non-linear models perform better. If the data were linear then linear regression would have outperformed the non-linear models. On the contrary, if the data were non-linear, the non-linear models would outperform the linear regression model.

Validation and Testing

Several linear regression models were fitted to the data via R programming language in an attempt to match the pattern of the relationship between the predictor variables and the response variable. This category of models works best when the nature of the relationship between the predictor variables and the response variable is linear in nature, or rather, can be represented by a straight line on a graph whose algebraic equation uses the variable to the power of 1. They work by finding the straight line that produces the smallest amount of RSME, or x-axis distance between the data points and the regression line. Linear regression models tried were Ordinary Least Squares (OLS), Ridge, Lasso, ElasticNet, Partial Least Squares (PLS), Principal Component Regression (PCR), and OLS Principal Component Analysis (PCA). The linear model function was used as the training method for OLS and OLS PCA models. With the linear regression function we trained the model on the training data for the OLS model. We then used the resulting tuning variable to rank the variables in order of importance. For the Ridge model, the penalty variable of lambda was the tuning parameter. We made a grid data frame of every combination of variables to use for tuning in the Ridge, ElasticNet, PLS and PCR models.

Non-linear regression models used to model the data were Support Vector Machines (SVM) Radial, K-Nearest Neighbor (KNN), Neural Network, Multivariate Adaptive Regression Splines (MARS), and SVM Polynomial. The “earth” package was used to train the MARS model. The train function was used to tune the MARS model via external resampling. We were also able to rank the variables in terms of importance according to the MARS model. For the KNN model, the k value was used for tuning grind in the train function. For SVM Radial and SVM Polynomial, the tuneLength argument uses 14 cost values as a default.

Tree Regression models, which use multiple, long series of if-then statements for predictor variable values and corresponding response variable outcome values, were also applied to the data to determine if this style of modeling would produce results superior to those of the non-linear models.

The Tree Regression models used on the data were Random Forest, Boosting, Cubist, Classification and Regression Tree (CART), and Bagged. With the Random Forest model we defined the number of trees. Similarly in the Bagged Trees model we defined the number of bags.

Results and Final Model Selection

Prediction accuracy using root mean-squared error (RMSE) will indicate the best performance models. RMSE is used to measure the difference between values predicted by the model and the values that are observed. In addition, 10-fold cross validation was used for all models. Moreover, linear, nonlinear, and tree regression models were used to analyze the data. With respect to the lower RMSE score, tree and nonlinear regression models outperformed the linear regression models. The only exception was the bagged tree model as it was ranked lower than five linear models. Linear models include ordinary least squares, partial least squares, principal component regression, and three penalized linear models: Lasso, Ridge, and Elastic Net. Support vector machines (radial and poly), multivariate adaptive regression splines, K-Nearest neighbors, and neural network made up the nonlinear regression models. Lastly, the tree regression models used were CART, random forest, boosting, bagged, and cubist.

Table 2

Performance Metrics On Test Set For All Models

Table 2

Model	Type	RMSE	R ²	MAE
Random Forest	Tree	3.8025	0.9501	2.7439
SVM Radial	Nonlinear	4.0610	0.9430	2.9680
KNN	Nonlinear	4.0612	0.9429	2.8716
Boosting	Tree	4.0679	0.9426	2.9963
Neural Network	Nonlinear	4.1510	0.9402	3.1339
MARS	Nonlinear	4.3579	0.9342	3.3181
Cubist	Tree	4.3640	0.9342	3.0827
SVM Poly	Nonlinear	4.4659	0.9309	3.4597
CART	Tree	4.5106	0.9295	3.4101
OLS	Linear	4.6642	0.9247	3.6306
Ridge	Linear	4.6642	0.9247	3.6306
ElasticNet	Linear	4.6642	0.9247	3.6306
Lasso	Linear	4.7216	0.9230	3.6852
PLS	Linear	4.9263	0.9160	3.8140
Bagged	Tree	5.1254	0.9091	3.9202
PCR	Linear	5.5115	0.8947	4.2382
OLS PCA	Linear	32.640	0.8962	29.075

Table 2 indicates RMSE, R^2 , and MAE for each model. The range of the target variable (electrical power output) is 75.5 with a mean of 454.4. This indicates that all of the models have suitable RMSE values with the exception of the Ordinary Least Squares model using Principal Component Analysis. As shown in the table, nonlinear models outperform linear models, which indicates that the data should be considered nonlinear. If the data were linear, the linear models would outperform the other models. Moreover, Ordinary Least Squares yielded the best linear regression with an RMSE score of 4.664. The best nonlinear regression was the Support Vector Machine (SVM) with radial basis function kernel and a RMSE score of 4.061. For the SVM model, the tuning parameter 'sigma' was held constant at a value of 0.4113507 with a C of 4. However, the best performing model was Random Forest which had the lowest RMSE score of 3.803. The data concludes that Random Forest models achieved the lowest RMSE and

highest R^2 . Therefore, two additional models were built to see what variation of Random Forest was most optimal. The original model had no tuning and was set at 500 trees with 1 variable tried at each split.

Table 3

Performance Metrics On Test Set For Random Forest Models

Table 3

Random Forest Model	RMSE	R Squared	MAE	
Tuned Model Using Out-of-Bag Estimates	3.7197	0.9521	2.6460	
Tuned Model Using Cross Validation	3.7257	0.9519	2.6488	
Original Model	3.8027	0.9501	2.7439	

In addition, the Random Forest model was tuned using out-of-bag (OOB) estimates and using cross validation. Among all three Random Forest models, the one tuned using out-of-bag had the lowest RMSE with the highest R^2 . Performance metric values on the test set are shown in Table 3. The OOB model was tuned to find the optimal mtry value; the final model used 500 trees with 2 variables tried at each split. In addition, the model explained 96.56% variance using these parameters.

Table 4

Variable Importance for Optimal Model - Random Forest (OOB)

Table 4

Predictor	Overall Importance
Relative Humidity	105.554
Ambient Temperature	70.8531
Ambient Pressure	52.3890
Turbine Exhaust Pressure	40.1479

Variable importance gives insights on how much a specific variable influences the model to make accurate predictions. Table 4 indicates the predictor variable importance used in the out-of-bag tuned Random Forest mode: in descending order, relative humidity with an importance of 105; ambient temperature with 70; ambient temperature with 52; and ambient pressure with 40. Therefore, we can assume that relative humidity was the most influential predictor used in the model.

Discussion and Conclusions

Machine learning models can be implemented to replace the need to incorporate thermodynamic approaches. Thermodynamic approaches often use assumptions in real life application in order to use appropriate nonlinear equations relating to the system. Therefore, machine learning models can create easier and more effective predictions by removing the obstacles surrounding assumptions and nonlinear equations when modeling the system.

The goal of this study is to improve upon Tufekci's study on the prediction of electrical power using machine learning methods. Tufekci's study concluded that the most successful model was bagging REPTree (using all predictors) with a RMSE of 3.787 and MAE of 2.818. To achieve this goal, the study focuses on finding the best model with a secondary objective of determining the predictor impact on the model. The most successful model will be regarded as the model with the lowest RMSE value. After analyzing the performance metrics on the test data set for each model used in this study, it was determined that Random Forest tuned by out-of-bag estimates yielded the best results. The final model used had a mtry value of 2 with an RMSE of 3.1709 and an R^2 of .965 on the training data set. However, the performance metrics on the testing data sets had a slightly higher RMSE of 3.7197, a lower R^2 of .952, and an MAE of 2.6460. The range of electrical power output (75.5) and mean (454.4) indicates that the RMSE value is suitable for accurate model prediction.

Furthermore, each model included all four predictors to obtain the best model. In addition, the predictor variable importance, which gives insight on how influential each predictor is, was calculated for each of the predictors in the optimal model. Among all the predictors, relative humidity had the highest importance by a significant amount, followed by ambient temperature, ambient pressure, and turbine exhaust pressure. With this insight, plant operators can potentially focus on manipulating the relative humidity and ambient temperature to reach a desired electrical output level.

In conclusion, the Random Forest model tuned by out-of-bag estimates slightly outperformed the most successful model in Tufekci's 2014 study while using the same feature subset. Furthermore, the prediction accuracy for the optimal machine learning model is suitable enough to replace thermodynamic approaches to model a real world system. By doing so, the time and effort to model thermodynamic systems analytics can be greatly reduced. Additionally, eliminating assumptions can improve prediction accuracy as well. Future works on this study should include different tuned models in efforts to perfect the ability to predict electrical output. As time goes on, better models are introduced to the data science industry, which can lead to better results than what the model provided in this study. In addition, data should be gathered for different types of power plants with similar unit designs since different unit designs can drastically affect the predictor values.

References

Kuhn, M., & Johnson, K. (2013). *Applied Predictive Modeling*. New York: Springer.

Tufekci, P. (2014). Prediction of full load electrical power output of a base load operated combined cycle power plant using machine learning methods. *Electrical Power and Energy Systems*, 60, 126-140.

Appendix