

# Case Studies 2024

## Forecasting the Electricity Price Approach Using Machine Learning

Professors

Prof. Dr. Matei Demetrescu, Prof. Dr. Paul Navas

Author: Fedi Ghanmi

Group Members

Alicia Hemmersbach, Ketevan Gurtskaia, Lev Luskin

21<sup>st</sup> Juin 2024

# Table of Contents

<b>1</b>	<b>Introduction</b>	<b>2</b>
<b>2</b>	<b>Problem and Data Description</b>	<b>3</b>
2.1	Project Objective . . . . .	3
2.2	Data Description . . . . .	3
2.2.1	Data Sources . . . . .	3
2.2.2	Data Preprocessing . . . . .	4
<b>3</b>	<b>Methods</b>	<b>5</b>
3.1	Statistical Methods . . . . .	6
3.1.1	Winsorization . . . . .	6
3.1.2	Regression Trees . . . . .	6
3.1.3	Random Forests . . . . .	7
3.1.4	Feature Importance in Tree Methods . . . . .	8
3.2	Relevant Tools . . . . .	9
<b>4</b>	<b>Empirical Results</b>	<b>10</b>
4.1	Price and Return Analysis . . . . .	10
4.2	Prediction Using Lagged Dependant Variable . . . . .	12
4.3	Prediction Using Additional Variables . . . . .	12
4.4	Feature Importance Results . . . . .	13
<b>5</b>	<b>Conclusion</b>	<b>16</b>
	<b>Bibliography</b>	<b>17</b>
	<b>Appendix</b>	<b>18</b>

# 1 Introduction

Immediately following the outbreak of the Russian-Ukraine war in 2022, electricity prices spiked across Europe. Among other factors, this increase can primarily be attributed to the reduction in natural gas supply. In response, most European countries sought alternative energy sources to compensate for the disruption. Germany is one of these countries, where in addition to the mentioned developments, have also shut down all nuclear production sites, officially becoming nuclear-free by 2023. These circumstances affected the merit order and significantly influenced electricity prices in the German market.

This study aims to address this phenomenon by understanding price fluctuation patterns in Germany and forecasting its future returns. To do this, historical records of them must be available. Hence a data gathering step on prices was made. Only the period between 2015 until 2024 was considered. Subsequently, two tree based models were used to predict the hour-ahead returns firstly only on the lagged dependant variable and secondly by using additional lagged predictors that were also used to forecast electricity demand in a previous study. It will then be determined whether these features contribute to the accuracy of price forecasting or not. As a result of this process, it was found that the prices could actually be forecasted over time. Random forest always outperformed a single decision tree with the least error measure value reaching an MSFE  $\approx 47$ . Ultimately, even with both tree methods, the variables used to predict electricity quantities proved to be merely random noise for price prediction when used collectively. Not only did they fail to improve the model, but they actually worsened the results, as demonstrated in the case of a single decision tree regression. In that instance, the MSFE increased from approximately 64 to 69. In the end, it is concluded that considering the variables individually provides more added value to price forecasting than when they are considered collectively.

In the following sections of this report, the initial chapter will explain more why a good forecast for returns is needed and describe in more details the data used to solve it. Following that, the methods employed will be briefly explained and a theoretical background on them will be provided. Finally, the empirical findings that were briefly discussed in this introduction will be detailed in its corresponding section and conclusions from this study will be drawn.

## 2 Problem and Data Description

In this section, firstly the project objective are outlined with the respective research questions needed to be answered by this study and secondly the data is described in details.

### 2.1 Project Objective

The primary objectives of this study is to have a good forecast of the electricity prices. Any market participant would need such information in order to maximize their profit by utilizing electricity when it is most affordable. It could also help grid operators <sup>1</sup> to balance the supply with the demand and ensure stable provision across the entire day. In this case, the first research question to be asked is: Could the return rates be predicted across time? To respond to this, lagged prices across hours will first be considered as a baseline model. If the answer is yes, then working with 2 statistical models, is it likely to state that in both cases other variables like temperature or CO<sub>2</sub> prices among many that proved useful to the load prediction are also equally important to the forecast of the returns? The primary measure that will be used to be guided in the right direction is the root mean squared forecast error <sup>2</sup>. Its minimization is the key objective while also taking into consideration the risk of overfitting using an out-of-sample set.

### 2.2 Data Description

Below in table 2.1 the variables used, their contextual meaning and units are recapitulated.

The date frame of interest spans from 1<sup>st</sup> of January 2015 until 15<sup>th</sup> of March 2024. Since this sample have an hourly frequency it contains  $N = 80.688$  observation.

#### 2.2.1 Data Sources

The CO<sub>2</sub> futures data come from the European Energy Exchange AG (EEX, 2002) and the temperature data was retrieved from an open source weather data provider (Zippenfenig, 2023). It is assumed that the capital of Germany, Berlin, represents a good proxy for the climate degree across the whole country. Additionally both the actual prices and

---

<sup>1</sup>Also called Transmission System Operators or TSO

<sup>2</sup>Abbreviation: MSFE

Name	Interpretation	Units
datetime_clean	date and time index	Hourly freq.
Price	Electricity prices	Euro(€) / MW
return	Electricity returns	Percentage (%)
co2_prices	CO2 price	Euro (€) / ton
temp	Temperature information	Celsius (°C)
solar_forecast	Day ahead solar energy forecast	Megawatt (MW)
fossil_output	Energy produced by fossil gas	Megawatt (MW)
geothermal_output	Energy from geothermal sources	Megawatt (MW)
hydro_output	Energy from hydro sources	Megawatt (MW)
other_output	Energy produced by other sources	Megawatt (MW)

Table 2.1: Variables Summary

the generation output per unit data were acquired from (ENTSO-e, 2008) which stands for European Network of Transmission System Operators for Electricity. They contain information respectively of the hourly recorded cost and the output of energy per each generation unit in Germany. The price data obtained include the prices within the bidding zone comprising Germany, Luxembourg, and Austria. These figures are assumed to serve as a good approximation for prices in Germany alone and are utilized for this project. They are already adjusted by default to the daylight savings time according to the CET zone <sup>3</sup>. Fossil\_output and geothermal\_output are collected from a new data source and are much cleaner with no missing values too.

### 2.2.2 Data Preprocessing

The forecasting on this data will be conducted using an out-of-sample test set that corresponds to 30% of the total data. The remaining 70% is used to train the model. Other features used in this study that were not adjusted by default to the Daylight Savings Time like CO<sub>2</sub> prices, solar\_output, hydro\_output, and other\_output were adjusted, pre-processed and interpolated using the same techniques applied from a previous project (Ghanmi, 2024). Hence they are already clean for this research. Their selection for this project was also because they were identified as the most promising features from the same study. The endogenous variable total\_load, which represents electricity demand, will be excluded to avoid any endogeneity issues. It is true that in the realm of machine learning, using a predictor that enhances out-of-sample forecast accuracy is generally acceptable.

---

<sup>3</sup>Central European Time Zone

But that comes with the condition that the test data is a good representative of the true generating process otherwise the model might perform well initially but deteriorate in the long run. Therefore, to prevent further assumption violations in our modeling techniques, The previously mentioned variables listed in table 2.1 without the total\_load will be proceeded with.

Since prices in a time series usually exhibit trends, detrending them is a necessity. Hence, a new variable *return* is introduced which is defined as the transformation below:

$$\text{return} = \frac{\text{price}_t - \text{price}_{t-1}}{\text{price}_{t-1}} \quad (2.1)$$

Here, eq. (2.1) can be regarded as the rate of change of price between two consecutive time steps  $t$  and  $t - 1$ (which is defined by hours in this study). One type of returns observed after the price preprocessing are funny returns where they refer to unexpected result. In this case, our model will struggles with handling returns that approach infinity. This issue occurs when the denominator in eq. (2.1) is nearly zero, causing the rate of change to suddenly spike. A small constant of 0.01 is added to these prices to fix this. Additionally, extremely large returns have been identified, which can negatively impact the performance of our tree methods. Models that rely on averaging like our case are susceptible to being influenced by such outliers. A known approach used to deal with this is called winsorization. See section 3.1.1 for more details of the method. In this project, a 99% Winsorization is applied. This percentage was selected based on practical intuition to keep our data values from not deviating too far from a 3 digit figure. Other percentages might yield better or worse results, but exploring these alternatives are beyond the scope of this study.

### 3 Methods

In this part both the statistical methods as well as the relevant tools used as part of this work will be depicted.

## 3.1 Statistical Methods

In this section, all statistical methods employed in our project are explained with their respective original authors.

### 3.1.1 Winsorization

According to (Dodge, 2003) winsorization was named after Charles P. Winsor (1895-1951). It is a method that is used to handle extreme data values and reduce the effect of huge outliers. An  $x\%$  winsorization replaces values below the lower  $\frac{(1-x)}{2}$  percentile and above the upper  $x + \frac{(1-x)}{2}$  percentile with the values at these respective percentiles. For example, a 90% winsorization replaces observations below the 5<sup>th</sup> percentile with the 5<sup>th</sup> percentile value and observations above the 95<sup>th</sup> percentile with the 95<sup>th</sup> percentile value itself.

### 3.1.2 Regression Trees

Decision trees are a type of supervised machine learning technique that follow the intuition of real-life trees as seen in the abstract fig. 3.1 below.

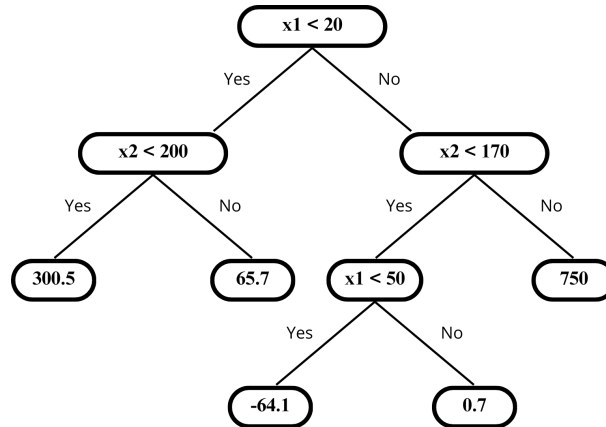


Figure 3.1: Simple abstraction of a regression tree on two random variables

Mathematically, they model the relationship between independent variables to determine the value of a dependent variable. A split is the process of dividing the data into subsets based on a specific condition related to the independent covariates which is in this example  $x_1$  and  $x_2$ . Each split within the tree results in an internal node. These nodes can further split, leading to additional internal nodes. The final splits, which contain the last sampled observations, are referred to as leaf nodes and represent the predicted values of the dependent variable. This process of splitting aims to accurately predict the

dependent variable. The topmost node in a tree is called a root node and represent the entire dataset before any splitting is made. The term regression trees was first introduced by (Breiman et al., 1984). They are an example of decision trees where each leaf node represents a continuous target and the prediction at each one of them is the average of the observations at that level. At each node, the algorithm chooses the split that minimizes some loss measure. The root mean squared forecast error can be considered in this case as it is defined as:

$$MSFE_z = \sqrt{\frac{1}{T} \sum_{t=1}^T (z_t - \hat{z}_t)^2}. \quad (3.1)$$

In eq. (3.1)  $z_t$  is the actual value,  $\hat{z}_t$  is the predicted value, and  $N$  is the number of observations. As already mentioned, the tree construction starts with the root node and splits it into subsets that are as similar as possible. The splitting for a variable  $j$  at a split point  $s$  can be formulated as below. See (Hastie et al., 2009):

$$\min_{j,s} \left[ \min_{c_1} \sum_{x_i \in R_1(j,s)} (y_i - c_1)^2 + \min_{c_2} \sum_{x_i \in R_2(j,s)} (y_i - c_2)^2 \right]. \quad (3.2)$$

The goal of eq. (3.2) is to minimize the squared difference between the observed values  $y_i$  and their means  $c_1$  and  $c_2$  after the split and to obtain  $R_1(j, s)$  and  $R_2(j, s)$  which are the two sets created by splitting on variable  $j$  at point  $s$  (Hastie et al., 2009). To avoid overfitting, the growth of the tree can be constrained by what is so called stopping criteria such as minimum node size or maximum tree depth (Breiman et al., 1984). The minimum node size is a hyperparameter chosen by the modeler that dictates the smallest number of observations a node must have before a split is attempted. This prevents the creation of nodes with small number of instances that could lead to less generalization. The maximum tree depth is the maximum number of layers a tree can grow. It serves as an upper bound to prevent the tree from becoming too complex and modeling noise in the data. Many other criteria also exist and are covered in (Hastie et al., 2009).

### 3.1.3 Random Forests

The Random forests method is an ensemble learning method that extends the simple case of a single tree as seen in section 3.1.2. It operates by constructing multiple decision trees



where each one is no longer fitted to the entire dataset but only to a specific bootstrap sample (Breiman, 2001). A bootstrap sample is a sampling technique with replacement. This means that each instance is randomly selected and can appear multiple times within each sample (Efron, 1979). These samples are then used to train each tree, and predictions are generated based on them. This can also be referred to as the bagging technique. For the case of regression, the algorithm averages the predictions from all individual  $B$  trees as eq. (3.3) shows:

$$\hat{f} = \frac{1}{B} \sum_{b=1}^B f_b(x') \quad (3.3)$$

What characterizes random forests as an ensemble method is its ability to reduce the variance through both the bagging process and the final averaging of the trees. For the tunable parameters, the number of trees to be selected is unlike the maximum tree depth parameter where when it is excessively large would lead to overfitting. On the contrary, the larger the number of trees the more accurate the prediction will be without any risk of overfitting. This is because the model will eventually average all outputs as eq. (3.3) depicts. This concept is mathematically summarized in eq. (3.4) according to (Breiman, 2001).

$$E_{X,Y} (Y - av_k h(X, \Theta_k))^2 \rightarrow E_{X,Y} (Y - E_{\Theta} h(X, \Theta))^2 \quad (3.4)$$

More clearly eq. (3.4) states that as the number of trees in a random forest increases to infinity, the expected residuals of the averaged prediction  $av_k$  from a finite number of  $k$  trees obtained from the function  $h(X, \Theta)$  converges to the expected residual of the ideal average of the prediction over all possible tree parameters  $\Theta$ . One constraint that should be mentioned is that this comes with an increased computational cost. As the number of trees grows, the computational expense increases, and the improvement in accuracy diminishes, eventually becoming negligible compared to the extra computational time required.

### 3.1.4 Feature Importance in Tree Methods

What tree methods have special that other methods do not is their ability to identify and return the most important features considered by the model during the fitting process.

The measure used in this project, which is also used in the case of general regression tree models is the mean decrease in impurity and is formulated for random forests according to (Breiman et al., 2017) as the eq. (3.5) below:

$$\text{unnormalized avg importance}(x) = \frac{1}{n_T} \sum_{i=1}^{n_T} \sum_{\text{node } j \in T_i | \text{split variable}(j)=x} p_{T_i}(j) \Delta i_{T_i}(j), \quad (3.5)$$

For a feature  $x$ , the sum is taken over all trees  $T_i$  considering all the nodes where feature  $x$  was used to make a split ( $\text{split variable}(j) = x$ ) and for each of these nodes, the importance is calculated by multiplying the proportion of samples in the node  $p_{T_i}(j)$  with the impurity decrease  $\Delta i_{T_i}(j)$  of that node. What is meant by impurity decrease is to actually have final nodes with the least variance possible, such that they are homogeneous and reflects similar properties. The feature importance is also of a similar concept in simple regression trees except that the importance is now no longer calculated over all the trees in the forest but for one single tree. However, one of the many limitation in the impurity-based importance is that they are derived from the training set. As a result, they may in some cases not necessarily indicate how effectively the variables will perform in making predictions that generalize well on the test set (Breiman, 2001). But as an approximation, when identifying these crucial variables, one can reduce the dimensionality of the data which can even lead in some cases to more accurate models due to the noise elimination. In this project, the feature importance scores generated using Python packages are presented as percentages to facilitate comparison. They are in more sense the normalized reduction in impurity from each feature over the total reduction.

## 3.2 Relevant Tools

As part of the work, Python version 3.11 (Guido van Rossum, 1989) with its open source development environment Jupyter Notebook (Jupyter, 2014) were used to conduct our study. External packages used are summarized in table A1 in the appendix.

## 4 Empirical Results

The results obtained from the practical implementation of this study are detailed in this section and are divided into different parts.

### 4.1 Price and Return Analysis

The price variation in fig. 4.1 shows stable fluctuations from 2015 until the end of 2020. Starting in 2021, there is a gradual increase until the end of the year, followed by a spike, taking prices to nearly 500 euros per megawatt at around July 2022, therefore reflecting an average of a 1900% increase. This change can be attributed to the Ukraine war, which led to restrictive natural gas supply, and the closure of all nuclear facilities in Germany, further contributing to a merit order shift. Later on, this spike gradually decreases to eventually approach the old range by early 2024.

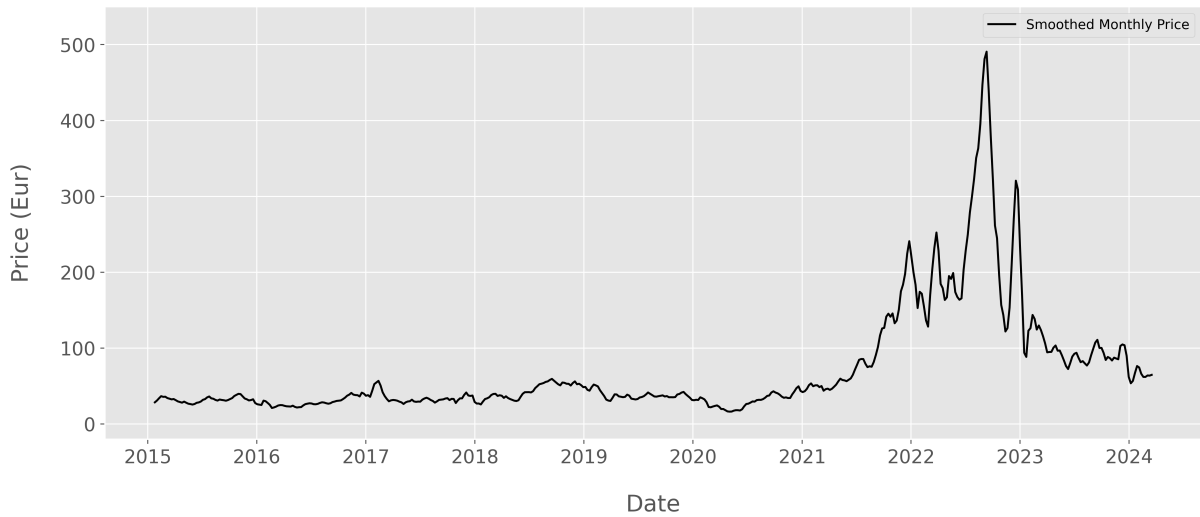


Figure 4.1: Running monthly average variation of electricity prices from 2015 to 2024 in Germany

The sudden jump in prices has made the variations amplitude between 2015 and 2021 seem insignificant. Therefore, the same smoothed curve at a lower interval is re-plotted in fig. 4.2. There, the electricity prices exhibit frequent rises and falls that eventually revolved around an average of 34 euros per megawatt.

These oscillations do not appear to be regular or follow a consistent pattern. For instance, between 2017 and 2018, the price did not vary as much as in the years before or after. It is evident that, if not for the political change mentioned earlier, we would

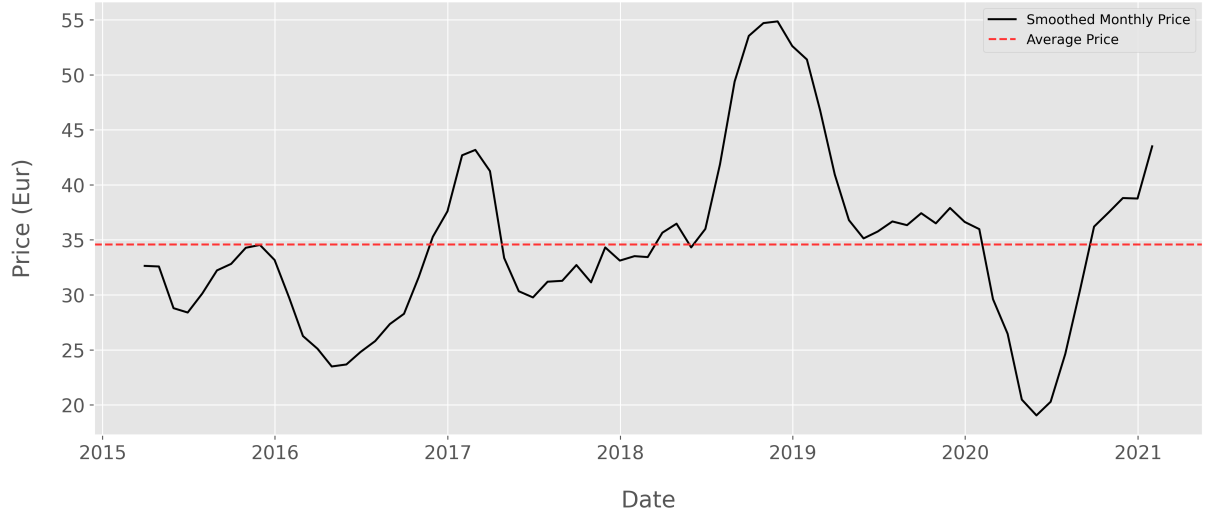


Figure 4.2: Running monthly average variation of electricity prices from 2015 to 2021 in Germany

have expected the prices to remain within the same interval of variation as in fig. 4.2. While other events may have also played a role, they will not be further addressed in this project. The analysis of prices has served merely as a basis to justify the distribution of returns, as the returns are derived from the relative one-hour change in price.

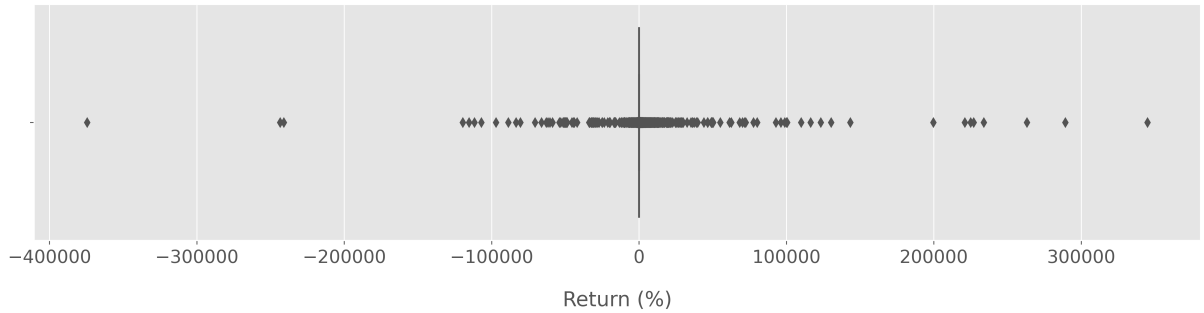


Figure 4.3: Distribution of the electricity returns in Germany from 2015 to 2024

As shown in fig. 4.3, the returns exhibit extreme values, causing the range to span from approximately -400.000% to 400.000%. This difference can well clearly be caused by the previously analyzed price spikes especially after the political shift starting 2022. Experiencing such large outliers can also occur during the night when prices are very low, and grid operators are required to sell electricity at negative fees to manage excess production. In this case, the difference in hourly prices is substantial, leading to relatively huge returns across hours. Using values in such range to later train tree models may not be optimal. Therefore, these outliers were limited using winsorization, as justified in section 2.2.2. This transformation is presumed to reduce the error rate in this case, not

necessarily because the method itself is "a great tool" (which still needs to be in a general sense verified), but because narrowing the range of returns will automatically reduce the average spread of errors.

In addition, it is assumed that any inconsistencies in the data due to political and environmental conditions or any transformation made to reduce them will be handled by the tree methods and constitutes one of the many ways it could have been addressed. Feeding our tree methods the returns we just analyzed and transformed to see if they can be efficiently predicted will now be the aim, as demonstrated in the following section.

## 4.2 Prediction Using Lagged Dependant Variable

Using only the return lagged three times, table 4.1 below is produced.

Model	MSFE
Tree Regression	64.14
Random Forests	49.66

Table 4.1: Mean squared forecast error of each model using only the lagged returns

It was decided to maintain the same lagged order setting as in the study of (Ghanmi, 2024). It is acknowledged that other lags could potentially yield better results, but exploring these alternatives is beyond this scope and may be pursued in future research. The numerical findings show enough reasonable results to state that returns do have some patterns that could be forecasted. The MSFE of the random forests was lower than the regression tree going from 64.14 to 49.66. This dissimilarity is believed to arise from the bagging strategy of the random forest and is expected to always show better error rates by comparison to the tree method. These values will serve as our baseline results and will be compared with the RMSE obtained after introducing additional variables in the next section.

## 4.3 Prediction Using Additional Variables

Adding further variables already mentioned in table 2.1 to both of our models, the following table 4.2 is obtained.

First of all, the tree regression has resulted in a worst and not a better error rate, increasing to 69.62. Secondly the random forest, as expected, had better RMSE. However,

Model	MSFE
Tree Regression	69.62
Random Forests	47.16

Table 4.2: Mean squared forecast error of each model using only the lagged returns with additional variables

it showed a very slight improvement, with the error rate decreasing from 49 to 47.16. This slight decrease relative to the number of variables added suggests that many of them may likely be just noise. It may be possible to explain the difference in MSFE between models by claiming that the bagging strategy of the random forests has minimized the impact of such noise, preventing the error from increasing as it did with a single decision tree. Independently from that, it is now evident that if these features are included altogether they may not necessarily lead to promising prices forecasts just because they have proven useful in predicting its demand (Ghanmi, 2024). To shed more light on this change, the feature importance results from both methods will be examined.

#### 4.4 Feature Importance Results

The MDI <sup>4</sup> output for both methods is plotted in the table 4.3 below.

Feature	MDI %	Feature	MDI %
return_lag_1	25.24	return_lag_1	28.03
return_lag_2	13.04	return_lag_2	11.82
other_output_lag_1	11.03	return_lag_3	9.37
return_lag_3	9.20	other_output_lag_1	7.62
solar_forecast_lag_1	7.86	fossil_output_lag_1	7.62
fossil_output_lag_1	7.78	geo_output_lag_1	7.56
geo_output_lag_1	7.23	solar_forecast_lag_1	7.49
co2_prices_lag_1	6.59	hydro_output_lag_1	7.36
temp_lag_1	6.12	temp_lag_1	7.15
hydro_output_lag_1	5.92	co2_prices_lag_1	5.97

(a) The normalized mean decrease in impurity of features for the single regression tree method

(b) The normalized mean decrease in impurity of features for the random forests method

Table 4.3: Comparative table of the difference in the normalized mean decrease in impurity for each of the regression tree and random forests method

The most relevant agreement that occurs between both methods is when they identify

<sup>4</sup>Abbreviation for Mean Decrease in Impurity

the lagged returns (lag 1 and lag 2) as having the highest importance, with a combined effect of around 38%. Temperature also have the same rank for both method with around a 7% of contribution. Other\_output, solar\_forecast, and co2\_prices show a drop in rank from the regression tree method to the random forest, accompanied by a slight decrease in their importance percentage. In contrast, return\_lag\_3, fossil, geo and hydro output all moved up in rank leading to a better mean decrease in impurity according to the Random Forest. The lowest contribution is attributed to hydro\_output in the decision trees and co2\_prices in the random forests with a value almost reaching 6%.

It was relatively expected that returns would significantly contribute to the decrease in impurity. However, other variables also played a role in reducing variance but their collective addition to the model did not lead to improved error measures according to table 4.2. As already mentioned, perhaps constructing a tree with all variables together generated poor error rates due to the noise created by their interactions. Therefore, an isolation of each variable individually is attempted to study their predictive power.

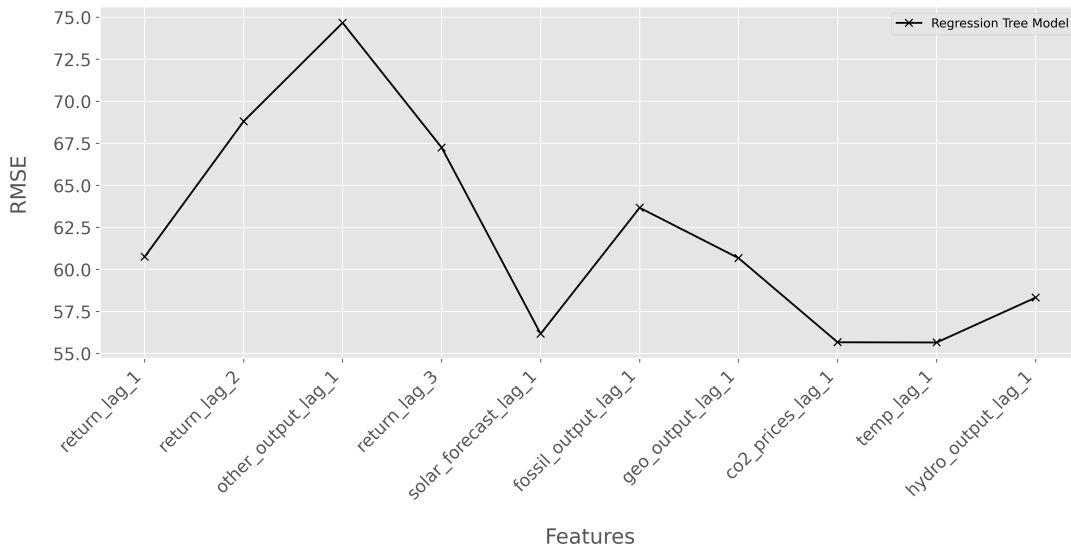


Figure 4.4: Performance of individual variables on the regression tree ranked by their decrease in impurity.

A pattern is observed in fig. 4.4 where, for the first three variables, a higher decrease in impurity corresponds to a lower RMSE. However, this pattern breaks around lower order variables, where the RMSE drops again.

The same variation explained is also seen in fig. 4.5 above. For the lagged temperature, both models produced an RMSE that was equal to or even better than that of the lagged returns while possessing lower importance. This proves that the evaluation of the decrease

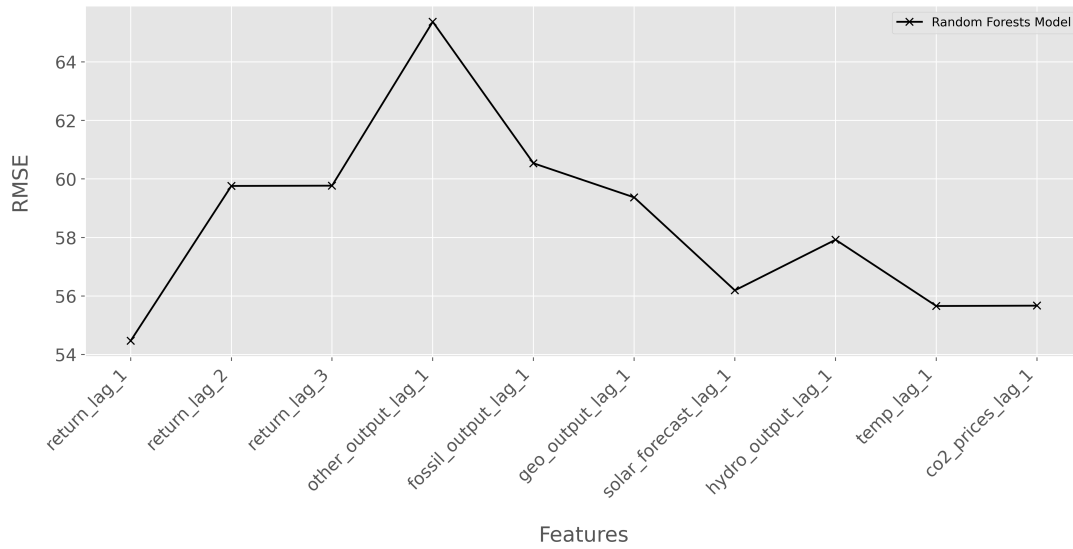


Figure 4.5: Performance of individual variables on the random forest, ranked by their decrease in impurity.

in impurity did not necessarily identify the variables that would be the best predictors. Instead, it served as an approximation based on the fitted data, showing how well each variable contributed to the tree formation relative to others. This discrepancy could be due to the biased evaluation of the MDI method already mentioned in section 3.1.4. It is at last reasonable to deduce that when these variables interact collectively, they may introduce noise and hinder the model's performance. Other specific combinations of these covariates will not be considered in this study. However, it is likely that certain interactions would result in lower error rates compared to others and even better than when all variables are used simultaneously. What matters up to this point is the demonstrated difference between the mean decrease in impurity as a measure of variable importance and its effectiveness as an indicator of the predictive power of that variable. This discovery allowed us to conclude that the same variables, which proved useful for the load, actually have the potential to decrease our baseline model errors if a careful combination is chosen.



## 5 Conclusion

Understanding demand patterns through return fluctuations and making profit by forecasting it was the main objective that motivated this research. This study takes place in the German market and seeks to reduce electricity usage costs for market consumers and ensure stable provision of it. Hence, having a model that describes how prices fluctuate would allow to make informed decisions about electricity usage.

The principal question that was started with was whether such returns could be forecasted over time or not and whether they could be supplemented by the same variables that proved promising when predicting electricity demand. Eventually, it was figured that electricity returns can actually be predicted over time and by using tree models it was possible to forecast one hour ahead values. From a loss of 64 using tree regression to a final loss of 47 using random forests, the models performance had almost a 36% improvement. This reduction in error is believed to be attributed to the use of the random forest technique and its bagging process that allows the model to learn more data patterns. Extending these models to include all new predictors simultaneously was not a critical reinforcement and did not lead to an important increase in predictions. However this proved to be contrary, where selecting a robust subset of them led to a more reliable forecasting. So it is eventually possible to say that some variables can also prove to be useful for price prediction. Relying solely on the mean decrease in impurity to identify these features was not enough as alternative strategies may perhaps be more useful for this task in this specific research.

Throughout the study, fixed settings and assumptions were maintained like the choice of the winsorization among many other. By holding these constant, reasonable results were achieved that not only aided in price forecasting but also demonstrated that using different techniques and parameters could potentially yield lower errors. This study can pave the way for further future research that could for instance explore the integration of specific tree methods designed for time series and the inclusion of more comprehensive data processing techniques to further enhance predictive capabilities. This study followed one of many possible approaches, and while it had its limitations, it still proved to be effective.

# Bibliography

Leo Breiman. Random forests. *Machine Learning*, 2001.

Leo Breiman, Jerome Friedman, Richard Olshen, and Charles Stone. *Classification and Regression Trees*. Wadsworth International Group, 1984.

Leo Breiman, Jerome Friedman, Richard Olshen, and Charles Stone. *Classification and Regression Trees*. Routledge, 2017.

Yadolah Dodge. The oxford dictionary of statistical terms. *Oxford University Press*, 2003.

EEX. European energy exchange. <https://www.eex.com/en/>, 2002. Accessed: May 16, 2024.

Bradley Efron. Bootstrap methods: another look at the jackknife. *The Annals of Statistics*, 1979.

ENTSO-e. Entso-e. <https://www.entsoe.eu/>, 2008. Accessed: May 14, 2024.

Fedi Ghanmi. Forecasting the electricity load: Approach using linear models and seasonal filtering. 2024.

Guido van Rossum. Python software foundation. <https://www.python.org/>, 1989. Accessed: May 12, 2024.

Trevor Hastie, Robert Tibshirani, and Jerome Friedman. *The Elements of Statistical Learning: Data Mining, Inference, and Prediction*. Springer, 2009.

Project Jupyter. Jupyter notebook. <https://jupyter.org>, 2014. Version 6.5.2.

Patrick Zippenfenig. Open-meteo.com weather api, 2023. URL <https://open-meteo.com/>.

# Appendix

## Appendix A

Package	Version
pip3	23.0.1
Jupyter Notebook	6.5.2
Pandas	1.5.1
matplotlib	3.7.1
seaborn	0.12.2
statsmodels	0.13.5
numpy	1.23.4
scikit-learn	1.5
mlxtend	0.23.1
datetime	* 5
warnings	* 5

Table A1: List of Python packages used across the project

---

<sup>5</sup>The datetime and warnigns package are a part of python core distribution and hence do not have a specific version.