

Technical report 21-011

Forecasting day-ahead electricity prices: A review of state-of-the-art algorithms, best practices and an open-access benchmark*

J. Lago, G. Marcjasz, B. De Schutter, and R. Weron

If you want to cite this report, please use the following reference instead:

J. Lago, G. Marcjasz, B. De Schutter, and R. Weron, "Forecasting day-ahead electricity prices: A review of state-of-the-art algorithms, best practices and an open-access benchmark," *Applied Energy*, vol. 293, July 2021. Article 116983. doi:[10.1016/j.apenergy.2021.116983](https://doi.org/10.1016/j.apenergy.2021.116983)

Delft Center for Systems and Control
Delft University of Technology
Mekelweg 2, 2628 CD Delft
The Netherlands
phone: +31-15-278.24.73 (secretary)
URL: <https://www.dcsc.tudelft.nl>

* This report can also be downloaded via https://pub.bartdeschutter.org/abs/21_011.html

Forecasting day-ahead electricity prices: A review of state-of-the-art algorithms, best practices and an open-access benchmark

Jesus Lago^{a,b,c,*}, Grzegorz Marcjasz^d, Bart De Schutter^a, Rafał Weron^d

^a*Delft Center for Systems and Control, Delft University of Technology, Delft, The Netherlands*

^b*Algorithms, Modeling, and Optimization, Energyville, Genk, Belgium*

^c*Energy Technology, Flemish Institute for Technological Research (VITO), Mol, Belgium*

^d*Department of Operations Research and Business Intelligence, Wrocław University of Science and Technology, Wrocław, Poland*

Abstract

While the field of electricity price forecasting has benefited from plenty of contributions in the last two decades, it arguably lacks a rigorous approach to evaluating new predictive algorithms. The latter are often compared using unique, not publicly available datasets and across too short and limited to one market test samples. The proposed new methods are rarely benchmarked against well established and well performing simpler models, the accuracy metrics are sometimes inadequate and testing the significance of differences in predictive performance is seldom conducted. Consequently, it is not clear which methods perform well nor what are the best practices when forecasting electricity prices. In this paper, we tackle these issues by performing a literature survey of state-of-the-art models, comparing state-of-the-art statistical and deep learning methods across multiple years and markets, and by putting forward a set of best practices. In addition, we make available the considered datasets, forecasts of the state-of-the-art models, and a specifically designed `python` toolbox, so that new algorithms can be rigorously evaluated in future studies.

Keywords: Electricity price forecasting, Deep learning, Open-access benchmark, Forecast evaluation, Best practices for price forecasting

1. Introduction

The increasing penetration of *renewable energy sources* (RES) in today's power systems makes electricity generation more volatile and the resulting electricity prices harder to predict [1–4]. On the other hand, advances in *electricity price forecasting* (EPF) constantly provide new tools with the ultimate objective of narrowing the gap between predictions and actual prices. The progress in this field, however, is not steady and easy to follow. In particular, as concluded by all major review publications, comparisons between EPF methods are very difficult since studies use different datasets, different software implementations, and different error measures; the lack of statistical rigor complicates these analyses even further [5–8]. In particular:

- There are several studies comparing *machine learning* (ML) and statistical methods but the conclusions of these studies are contradictory. Typically, studies considering advanced statistical techniques only compare them with simple ML methods [9–11] and show that statistical methods are obviously better. Conversely, studies proposing new ML methods only compare them with simple statistical methods [12–16] and show that ML models are more accurate.

*Corresponding author

Email address: j.lagogarcia@tudelft.nl, jesus.lagogarcia@vito.be (Jesus Lago)

- In many of the existing studies [17–23] the testing periods are usually too short to yield conclusive results. In some cases, the test datasets are limited to one-week periods [22, 24–30]; this ignores the problem of special days, e.g. holidays, and is not representative for the performance of the proposed algorithms across a whole year. As argued in [5], to have meaningful conclusions, the test dataset should span at least a year.
- Some of the existing papers do not provide enough details to reproduce the research. The three most common issues are: (i) not specifying the exact split between the training and test dataset [31–37], (ii) not indicating the inputs used for the prediction model [35, 36, 38–40], and (iii) not specifying the dataset employed [21, 33, 41, 42]. This obviously prevents other researchers from validating the research results.

These three problems have aggravated over the last years with the increase in popularity of *deep learning* (DL). While new published papers on DL for EPF appear almost every month, and most claim to develop models that obtain state-of-the-art accuracy, the comparisons performed in those papers are very limited. Particularly, the new DL methods are usually compared with simpler ML methods [28, 30, 43–47]. This is obviously problematic as such comparisons are not fair. Moreover, as the proposed methods are not compared with other DL algorithms, new DL methods are continuously being proposed but it is unclear how the different models perform relatively to each other.

Similar problems arise in the context of *hybrid methods*. In recent years, very complex hybrid methods have been proposed. Typically, these hybrid models are based on combining a decomposition technique, a feature selection method, an ML regression model, and sometimes a type of genetic algorithm for optimization purposes. As with DL algorithms, these studies usually avoid comparisons with well-established methods [21, 25, 34, 42, 48–50] or resort to comparisons using outdated methodologies [22, 24, 26, 37, 41, 51, 52]. In addition, while a specific genetic algorithm or decomposition technique is considered, most of the studies do not analyze the effect of selecting a variant of these techniques [21, 24, 50–52]. Thus, the relative importance of each of the different components of the hybrid methods it is not even clear.

1.1. Motivation and contributions

The above mentioned problems call for three actions. Firstly, implementing in a popular programming environment (e.g. `python`), thoroughly testing and making available *a set of simple but powerful open-source forecasting methods*, which can potentially obtain state-of-the-art results, and that researchers can easily use to evaluate any new forecasting model.

Secondly, collecting and making freely available to the EPF community *a set of representative benchmark datasets* that researchers can use to evaluate and compare their methods using long testing periods. Although, some datasets are available for download without restrictions, e.g. as supplements to published articles [53] or sample transaction data [54], they are typically limited in scope (one market, a 2-3 year timespan or price series only). Hence, conclusions from such datasets are limited, results can hardly be extrapolated to other markets, and the relevance of the studies using such data are not entirely clear.

Thirdly, *putting forward a set of best practices* so that the conclusions of EPF studies become more meaningful and fair comparisons can be made.

In this paper, we try to tackle the above via three distinct contributions:

1. We analyze the existing literature and select what could arguably be considered as state-of-the-art among statistical and machine learning methods: the *Lasso Estimated AutoRegressive* (LEAR) model¹ [55] and the *Deep Neural Network* (DNN) [59], a relatively simple and automated DL method that optimizes hyperparameters and features using Bayesian optimization. Then, we make our models open-source and available to other researchers as part of an open-source `python` library <https://>

¹Originally introduced in [55] under the name *LassoX* and based on the *fARX* model, a parameter-rich autoregressive specification with exogenous variables. The name refers to the *least absolute shrinkage and selection operator* (LASSO) [56] used to jointly select features and estimate their parameters. Very similar models were used in [57] under the name *24lasso_{DOW,nl}* and in [58] under the name *24Lasso₁*.

github.com/jeslago/epftoolbox specially designed for this study to provide a common research framework for EPF research [60]. Besides the models, we also provide extensive documentation [61] for the library.

2. We propose a set of five open-access benchmark datasets spanning six years each, that represent a range of well-established day-ahead, auction type power markets from around the globe. The datasets contain day-ahead electricity prices at an hourly frequency and two relevant exogenous variables each. They can be accessed from the mentioned `python` library [60] that is specially designed for this study. Together with the datasets, the library also includes the forecasts of the open-access methods across the five benchmark datasets so that researchers can quickly make further comparisons without having to re-train or re-estimate the models.
3. We provide a set of best practice guidelines to conduct research in EPF so that new studies are more sound, reproducible, and the obtained conclusions are stronger. In addition, we include some of the guidelines, e.g. adequate evaluation metrics or statistical tests, in the mentioned `python` library [60] that is specially designed for this study to provide a common research framework for EPF research

1.2. Paper structure

The remainder of the paper is organized as follows. Section 2 performs a literature review of the current state of EPF. Sections 3 and 4 respectively present the open-access benchmark datasets and the open-source benchmark models. Section 5 describes the set of guidelines and best practices when performing research in EPF. Section 6 discusses the forecasting results for all five datasets. Finally, Section 7 provides a summary and a checklist of the requirements for meaningful EPF research.

2. Literature review

The field of EPF aims at predicting the spot and forward prices in wholesale markets, either in a point or probabilistic setting. However, given the diversity of trading regulations available across the globe, EPF always has to be tailored to the specific market. For instance, the workhorse of European short-term power trading is the *day-ahead* market with its once-per-day uniform-price auction, see Fig. 1. On the other hand, the Australian National Electricity Market operates as a real-time power pool, where a dispatch price is determined every five minutes and six dispatch prices are averaged every half hour as pool prices [62], while electricity forward markets share many aspects with those of other energy commodities (oil, gas, coal), and quite often are only financially settled [63].

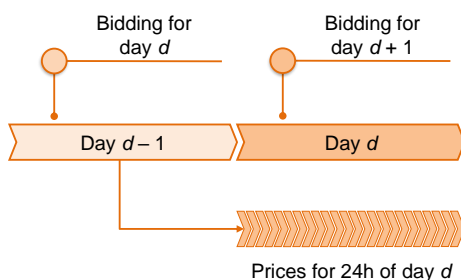


Figure 1: Illustration of the *day-ahead* auction market, where wholesale sellers and buyers submit their bids before gate closure on day $d - 1$ for the delivery of electricity during day d ; the 24 hourly prices for day d are set simultaneously, typically around midday.

As the field of EPF is very diverse, a complete literature review is out of the scope of this paper. Instead, this section is intended to provide an overview of the three families of methods, i.e. statistical, ML, and hybrid methods, proposed for point forecasting in day-ahead markets since 2014, i.e. since the last comprehensive literature review of Weron [5]. The more recent reviews either focused on short-term [6] and medium-/long-term [7] probabilistic EPF, were not that comprehensive in scope [64, 65], or concerned

electricity derivatives [63]. Furthermore, our survey puts a special emphasis on DL and hybrid methods as this is the area of EPF characterized by the most rapid development and, at the same time, troubled by non-rigorous empirical studies which motivated us to write this paper in the first place.

2.1. Statistical methods

Most models in this class rely on linear regression and represent the dependent (or output) variable, i.e. the price $p_{d,h}$ for day d and hour h , by a linear combination of independent (or predictor, explanatory) variables, also called regressors, inputs, or features:

$$p_{d,h} = \boldsymbol{\theta}_h \mathbf{X}_{d,h} + \varepsilon_{d,h}, \quad (1)$$

where $\boldsymbol{\theta}_h = [\theta_{h,0}, \theta_{h,1}, \dots, \theta_{h,n}]$ is a row vector of coefficients specific to hour h , $\mathbf{X}_{d,h} = [1, X_{d,h}^1, \dots, X_{d,h}^n]^\top$ is a column vector of inputs and $\varepsilon_{d,h}$ is an error term; the intercept $\theta_{h,0}$ can be set to zero if the data is demeaned beforehand. Note that here we are using a notation common in day-ahead forecasting, which emphasizes the vector structure of these price series, see Fig. 1. Alternatively we could use single indexing: p_t with $t = 24d + h$. Although the multivariate modeling framework has been shown to be marginally more accurate than the univariate counterpart, both approaches have their pros and cons [57, 66].

In the last few years, there have been several key contributions in the field of statistical methods for EPF. Arguably, the most relevant of them has been the appearance of linear regression models with a large number of input features that utilize regularization techniques [56, 67]. Classically, the regression model in (1) is estimated using ordinary least squares by minimizing the *residual sum of squares* (RSS), i.e. squared differences between the predicted and actual values. However, if the number of regressors is large, using the *least absolute shrinkage and selection operator* (LASSO) [56] or its generalization the *elastic net* [67] as implicit feature selection methods have been shown to improved the forecasting results [55, 57–59, 68, 69]. In particular, by jointly minimizing the RSS and a penalty factor of the model parameters (see Section 4.2 for details), these two implicit regularization techniques set some of the parameters to zero and thus effectively eliminate redundant regressors. As shown in the cited studies [55, 57–59, 68, 69], these parameter-rich² regularized regression models exhibit superior performance. It is important to note that such an approach, called here *Lasso Estimated AutoRegressive* (LEAR), is in fact hybrid since LASSO (and elastic nets) are considered ML techniques by some authors. However, we classify it as statistical because the underlying model is autoregressive.

Aside from proposing parameter-rich models and advanced estimators, researchers have also improved the field by considering a variety of additional preprocessing techniques. Most notably, models using so-called variance stabilizing transformations [9, 57, 70, 71] and long-term seasonal components [72–75] have been proposed and shown to result in statistically significant improvements. However, the applicability of these two techniques varies greatly: due to very common occurrence of price spikes, variance stabilizing transformations have become a standard and replaced the commonly used logarithmic transformation (no longer applicable due to zeros and negative values³) to normalize electricity prices. By contrast, the applicability of long-term seasonal components has been more limited and it is unknown whether their beneficial effect is limited to relatively parsimonious regression models or also holds for parameter-rich models.

A third innovation in the field is an ensemble, i.e. a method that combines individual forecasting models to enhance the accuracy, that combines multiple forecasts of the same model calibrated on different windows. In this context, two different studies [76, 77] showed that the best results are obtained with a combination of a few short (spanning 1-4 months) and a few long calibration windows (of approximately two years). Said ensembles were able to significantly outperform predictions obtained for the best ex-post selected calibration window [76, 77]. But again, it has not been shown to date whether this effect is limited to relatively parsimonious regression models or also holds for LEAR models.

²We define a parameter-rich linear model as a model with multiple regressors.

³The logarithmic of 0 or a negative value is undefined.

2.2. Deep learning

In the last five years, a total of 28 deep learning papers in the context of EPF have been published⁴. Moreover, this number has been steadily increasing: while in 2016 there was only one paper and in 2017 none, in 2018 there were 11, and in 2019 there were 16. Despite this trend, most of the published studies are very limited: the comparisons are too simplistic, e.g. avoid state-of-the-art statistical methods, and their results cannot be generalized.

The first published DL paper [12] proposes a deep learning network using stacked denoising autoencoders. The paper, despite being the first, provides a better evaluation than most studies: the new method is compared not only against machine learning techniques but also against two statistical methods. Yet, the evaluation is limited as it is done considering three months of test data and employing simple models for comparison. In the second published DL paper [43], a DNN for modeling market integration is proposed. While the method is evaluated over a year of data, the study is also limited as the proposed model is not compared against other machine learning or statistical methods.

In the third published paper [59], four DL models (a DNNs, two *recurrent neural networks (RNNs)*, and a *convolutional network (CNN)*) are proposed. This study is, to the best of our knowledge, the most complete study up to date. In particular, the proposed DL models are compared using a whole year of data against a benchmark of 23 different models, including 7 machine learning models, 15 statistical methods, a commercial software. Moreover, among the statistical methods, the comparison includes the fARX-Lasso and fARX-EN, i.e. the state-of-the-art statistical methods. While the study shows the superiority of the DL algorithms, very strong conclusions are not possible as the study only considers a single market.

The studies that followed in 2018 focused on one of three topics: 1) evaluating the performance of different deep recurrent networks [13, 23, 37, 78]; 2) proposing new hybrid methods based on CNNs and LSTMs [14, 44, 79, 80]; or 3) employing regular DNN models [23]. Independently of the focus, they were all more limited than the first and the third studies [12, 59] as they failed to compare the new DL models with state-of-the-art statistical methods and/or to employ long enough datasets to derive strong conclusions.

In detail, [13] studies the use of RNNs for forecasting electricity prices but the comparison is done in a single market and against simple statistical methods (a seasonal *auto regressive integrated moving average (ARIMA)* model, a Markov regime-switching model, and a self exciting threshold model). Moreover, while the comparison includes other DL methods, it avoids comparison with simpler ML techniques. Ref. [44] proposes a hybrid DL methods composed of a CNN and a *long short-term memory (LSTM)* (a type of recurrent network) for forecasting balancing prices. However, the new model is only compared against simple ML benchmarks and the evaluation is done using different periods comprising three months for training and 1 month for testing. Similarly, [14] proposes another hybrid model combining a CNN and an LSTM, but the model is only compared against two naive statistical methods: an *auto regressive moving average (ARMA)* and a *generalized autoregressive conditional heteroscedasticity (GARCH)* model.

In [23] a regular DNN model is proposed but the model is only evaluated on a test dataset comprising a single day and compared against a simple MLP. In [29], the use of an LSTM model for EPF is evaluated, but the method is only compared with three neural networks and a simple statistical method, and the evaluation is done using only 4 weeks of data. Likewise, [78] proposes a model based on an LSTM but a comparison against other methods is not done and the test dataset only comprises 2 weeks of data. In [37], another LSTM model is proposed but, as other studies, the test dataset comprises some months of data and the method is only compared against a simple decision tree and a support vector regressor; moreover, the exact split between the training and test dataset is not specified and it is unclear what is exactly the performance of the model. An exception to these studies is [81] which proposes a series of DL models and compares them

⁴This data is based on two searches in Scopus looking for keywords in the title, abstract, and keywords. The first search is based on the following query TITLE-ABS-KEY((((("forecasting electricity") OR ("predicting electricity")) AND (("electricity spot") OR ("electricity day-ahead") OR ("electricity price")))) OR (((("price forecasting") OR ("price prediction") OR ("forecasting price") OR ("predicting price") OR ("forecasting spikes") OR ("forecasting VAR")) AND (("electricity spot price") OR ("electricity price") OR ("electricity market") OR ("day-ahead market") OR ("power market")))) AND ("deep") AND ("learning")). The second search is very similar but replacing ("deep") AND ("learning") by ("neural") AND ("network").

for a year of data against several advanced statistical methods such as LASSO and a simpler ML method. The main drawbacks of the study are that it is based on a single market and that it only considers a simple ML method as a benchmark. In addition, the study focuses on intraday electricity prices, while most of the literature (including the current paper) considers forecasting day-ahead electricity prices.

In 2019, the main focus of the papers was the same as in 2018: 1) evaluating the performance of different deep recurrent networks (mostly LSTMs) [16, 30, 45, 47, 82–84], 2) proposing new hybrid deep learning methods usually based on LSTMs and CNNs [17, 28, 36, 82, 85–87], or 3) employing regular DNN models [15, 46, 88]. Similarly, as with most studies in 2018, the new studies were more limited than [12, 59] as no comparisons with state-of-the-art statistical methods were made and long test datasets were seldom used. In this context, even though some studies [16, 88] tried to compare the proposed methods with existing DL models [59], they either failed to re-estimated the benchmark models for the new case study [16] or they overfitted the DL benchmark models [88].

In detail, [30] proposes different LSTM models but the new models are only compared against 5 other ML techniques and using a test period of 4 weeks. In [28], a CNN model is proposed but the new model is just compared against three simple ML methods and using a test dataset that comprises a week. In [45], a model based on an LSTM is proposed but it is only compared against three simple ML methods and for a period of 12 weeks. In [46], the performance of a DNN is compared to that of an SVR model and, as the comparison only includes these two models, it is obviously very limited. In [15], a DNN is used as part of a two-step forecasting method; as in many other studies, the comparison is performed for one month of data and limited to two simple ML models (a SVR and an MLP) and a standard linear model. In [47], two DL models are proposed but the models are only compared to very simple ML methods (extreme learning machines and standard MLPs) and using a test dataset spanning eight months. In [16], a bidirectional LSTM to forecast prices in the French market is proposed; however, the study only considers historical prices as input features and the proposed method is only compared against DL models and a simple autoregressive model. In addition, the benchmark DL models are copied from [59] (a completely different case study that considers exogenous inputs and a different market) without re-tuning the hyperparameters to the new case study.

In [88], a neural network that uses data from order books is proposed and compared against DL methods from the literature, e.g. the ones proposed in [59]. While the new model outperforms existing DL methods, the DL methods from the literature are trained to overfit the training dataset⁵. Therefore, the comparison is not meaningful (the DL benchmark models will necessarily perform poorly in the test dataset) and it cannot be assessed how the new model performs. In [85], a hybrid DL forecasting method is proposed based on stacked denoising autoencoders for pre-training, regular autoencoders for feature selection, and a rough DNN as a forecasting method. As other studies, the method is only compared against other simpler ML models. Moreover, the importance of each of the four modules of the hybrid method is not studied and the study does not re-calibrate the models with new data: the models are trained once and evaluated during a whole year. Similarly, [86] proposes a CNN hybrid model that uses mutual information, random forests, gray correlation analysis, and recursive feature elimination for feature selection. Unlike most models, the algorithm is trained to classify prices instead of predicting their scalar values; however, details of how this process is done are not provided. In addition, the method is only compared against simpler ML methods and evaluated for less than a year of data (the study uses 1 year for testing and training but the split is not specified). Likewise, [36] proposes a hybrid model based on CNNs and RNNs in the context of microgrids; as other studies, the method is evaluated in a small dataset, it is not compared against state-of-the-art statistical methods, and the exact split between training and test datasets is not specified.

⁵In the training dataset, the proposed model and some naive ML benchmark models yield a root mean square error in the order of 6. For the test dataset, for the same models, that error is between 9 and 12. By contrast, the training error of the benchmark DL model is 2, and the test error is 20. Having a training error that is 1/3 of the error of other models but a test error that is 10 times larger than the training error is a clear sign for overfitting (especially when for the rest of the models the test error is just 1.5 larger than the training error).

2.3. Hybrid methods

Within the field of EPF, the research area that has received most attention in the last 5 years has been hybrid forecasting methods. In this time frame, more than 100 articles proposing new hybrid methods have been published⁶, i.e. approximately 5 times more than articles based on DL. Hybrid models are very complex forecasting frameworks that are composed of several algorithms. Usually, they comprise at least two of the following five modules:

- An algorithm for decomposing data.
- An algorithm for feature selection.
- An algorithm to cluster data.
- One or more forecasting models whose predictions are combined.
- Some type of heuristic optimization algorithm to either estimate the models or their hyperparameters.

In terms of decomposition methods, the most widely used technique is the wavelet transform [17, 19, 22, 24, 34, 41, 49, 51, 52, 89]. Alternatives methods include empirical mode decomposition [32, 90], variational mode decomposition [27, 48], and singular spectrum analysis [91, 92].

For feature selection, the most commonly used algorithms are correlation analysis [32, 41, 42, 93, 94] and the mutual information technique [18, 42, 52, 95–97]. Other algorithms include classification and regression trees with recursive feature elimination [50] or Relief-F [50].

For clustering data, the algorithms are usually based on one of the following four: k-means [26, 98], self-organizing maps [19, 26, 99], enhanced game theoretic clustering [26], or fuzzy clustering [52, 100]

The selection of forecasting models is much more diverse. The most widely used method is the standard MLP [19, 20, 32, 41, 42, 51, 91, 92, 94, 96, 97], followed by the *adaptive network-based fuzzy inference system* (ANFIS) [19, 90, 95], radial basis function network [20, 24, 100], and autoregressive models like ARMA or ARIMA [20, 22, 24, 90]. Other models include LSTM [17], linear regression [50], extreme learning machine [22, 50], CNN [50], Bayesian neural network [26, 99], exponential GARCH [90], echo state neural network [27], Elman neural networks [18], and support vector regressors [20]. It is important to note that in many of the approaches, the hybrid method does not consider a single forecasting model but combines several of them [19, 20, 24, 50, 90, 97].

Just as for the forecasting model, the diversity of the heuristic optimization algorithms is also large. While the most often utilized algorithm is particle swarm optimization [22, 48, 51, 95, 96, 100], many other approaches are also used: differential evolution [27], genetic algorithm [95], backtracking search [95], deterministic annealing [100], bat algorithm [41], vaporization precipitation-based water cycle algorithm [93], cuckoo search [92, 94], or honey bee mating optimization [24].

In spite of the large number of published works, the research in hybrid methods suffers from the same problems as discussed earlier. First, most of the studies either avoid comparison with well-established methods [18–21, 25, 27, 34, 42, 48–50, 90, 93, 95, 100] or resort to comparisons using outdated methodologies [22, 24, 26, 41, 51, 52, 91, 92]. Hence, the accuracy of the new proposed methods cannot be accurately established.

Second, the considered studies usually employ very small datasets consisting either of a few days [17–22] or a few weeks [18, 19, 22, 24–27, 41, 42, 49, 51, 91–93, 95, 100]. Thus, drawing conclusions is nearly impossible and it is unclear whether the accuracy results are just the outcome of selecting a convenient test period.

Besides these two problems, for many hybrid methods the effect of selecting variants of the different hybrid components is not analyzed [20, 21, 24, 25, 27, 41, 42, 50–52, 91, 92]. Thus, it is not clear how relevant or useful the individual components are.

⁶This data is based on two searches in Scopus looking for keywords in the title, abstract, and keywords. The first search is based on the following query `TITLE-ABS-KEY(((forecast*) OR (predict*)) AND (electricity) AND (price*) AND (hybrid))`. The second search is very similar but replacing the keyword `hybrid` by `neural AND network`. Note that, while this search is not as complete as the one for DL, it provides enough material for building an overview of the state of the field.

2.4. State-of-the-art models

Because of the described problems when comparing EPF models, it is very hard to establish what are the state-of-the-art methods. Nevertheless, considering the studies performed in the last years, it can be argued that the LEAR is a very accurate (if not the most accurate) linear model. Moreover, it can also be argued that the accuracy of this model can be further improved by transforming the prices using variance stabilizing transformations, combining forecasts obtained for different calibration windows, and/or using long-term seasonal decomposition.

For the case of ML models, the selection is harder as the existing comparisons are of worse quality. Considering the most complete benchmark study in terms of forecasting models [59], it seems that a simple DNN with two layers is one of the best ML models. In particular, while more complex models, e.g. LSTMs, could potentially be more accurate, there is at the moment no sound evidence to validate this claim.

In the case of hybrid models, establishing what is the best model is an impossible task. Firstly, while many hybrid methods have been proposed, they have not been compared with each other nor with the LEAR or DNN models. Secondly, as most studies do not evaluate the individual influence of each hybrid component, it is also impossible to establish the best algorithms for each hybrid component, e.g. it is unclear what are the best clustering, feature selection method, or data decomposition methods.

With that in mind, we will consider the LEAR and the DNN for the proposed open-access benchmark. In particular, not only are these two methods highly accurate, but they are also relatively simple. As such, we think that they are the best benchmarks to compare new complex EPF forecasting methods with.

3. Open-access benchmark dataset

The first contribution of the paper is to provide a large open-access benchmark dataset on which new methods can be tested, together with the day-ahead forecasts of the proposed open-access methods. In this section, we introduce this dataset, which can be accessed⁷ using the `python` library built for this study.

3.1. General characteristics

For a benchmark dataset in EPF to be fair it should satisfy three conditions: 1) comprise several electricity markets so that the capabilities of new models can be tested under different conditions, 2) be long enough so that algorithms can be analyzed using out-of-sample datasets that span 1-2 years, and 3) be recent enough to include price effects due to the integration of RES.

Based on these conditions, we propose five datasets representing five different day-ahead electricity markets, each of them comprising 6 years of data. The prices of each market have very distinct dynamics, i.e. they all have differences in terms of the frequency and existence of negative prices, zeros prices, and price spikes. In addition, as electricity prices depend on exogenous variables, each dataset comprises two additional time series: day-ahead forecasts of two influential exogenous factors that differ from each market. The length of each dataset equals 2184 days, which translates to six years of 364 days or $6 \times 52 = 312$ weeks⁸. All available time series are saved using the local time, and the daylight savings are treated by either arithmetically averaging two values from the extra hour or interpolating the neighboring values for the missing observation.

3.2. Nord Pool

The first dataset represents the Nord Pool (NP), i.e. the European power market of the Nordic countries, and spans from 01.01.2013 to 24.12.2018. The dataset contains hourly observations of day-ahead prices, the day-ahead load forecast, and the day-ahead wind generation forecast. The dataset was constructed using the

⁷Note that we do not own the data in the dataset. However, it can be freely accessed from different websites, e.g. the ENTSO-E transparency platform [101]. In this context, the proposed `python` library [60, 61] provides an interface to easily access the data.

⁸Electricity prices have weekly seasonality. Thus, by approximating a year by 52 weeks because we ensure that the metrics are not offset because a certain day, e.g. Monday, is harder to predict than the others.

data freely available on the webpage of the Nordic power exchange Nord Pool [54]. Figure 2 (top) displays the electricity price time series of the dataset; as can be seen, the prices are always positives, zero prices are rare, and prices spikes seldom occur.

3.3. PJM

The second dataset is obtained from the *Pennsylvania-New Jersey-Maryland* (PJM) market in the United States. It covers the same data points as Nord Pool, i.e. from 01.01.2013 to 24.12.2018. The three time series are: the zonal prices in the *Commonwealth Edison* (COMED) (a zone located in the state of Illinois) and two day-ahead load forecast series, one describing the system load and the second one the COMED zonal load. The data is freely available on the PJM’s website [102]. Figure 2 (bottom) represents the electricity price time series of the dataset; as with the NP market, the prices are always positive and zero prices are rare; however, unlike with the prices in the NP market, spikes appear frequently.

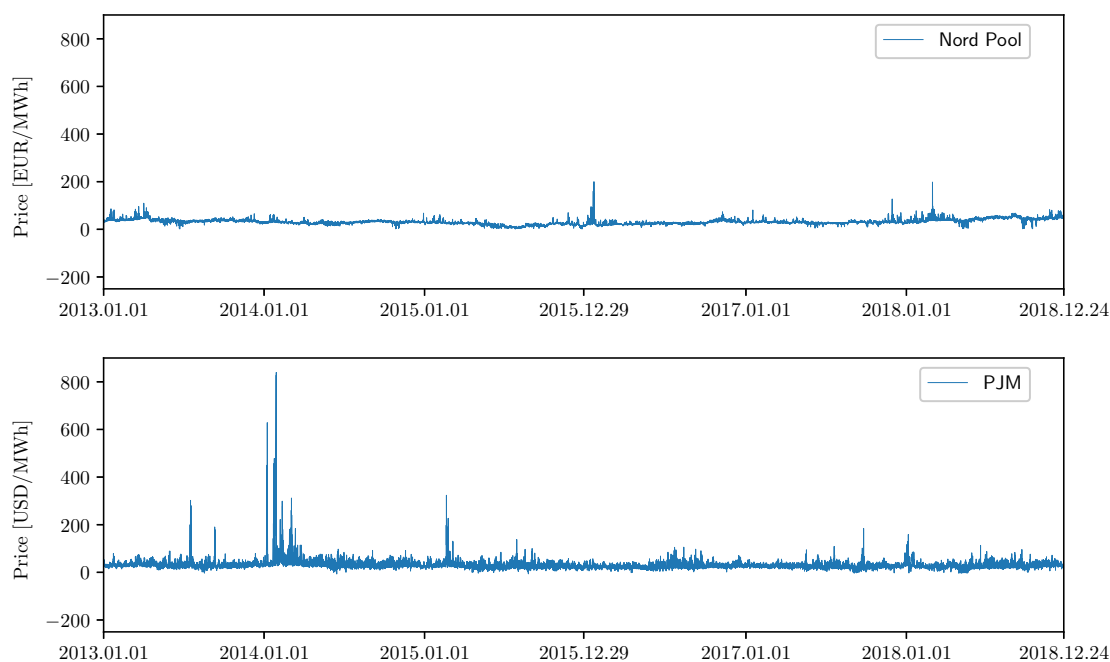


Figure 2: Electricity prices time series for two of the five datasets, i.e. Nord Pool and PJM, considered in the open-access benchmark dataset (Note that each dataset also includes two exogenous time series that are not plotted here).

3.4. EPEX-BE

The third dataset represents the EPEX-BE market, the day-ahead electricity market in Belgium, which is operated by EPEX SPOT. The dataset spans from 09.01.2011 to 31.12.2016. The two exogenous data series represent the day-ahead load forecast and the day-ahead generation forecast in France. While this selection might be surprising, it has been shown [43] that these two are the best predictors of Belgian prices. The price data is freely available in the ENTSO-E transparency platform [101] and the ELIA website [103], and the load and generation day-ahead forecasts are freely available in [104]. It is important to note that this dataset is particularly interesting because it is harder to predict. Figure 3 (top) shows the electricity price time series of the dataset; unlike the prices in the PJM and NP markets, negative prices and zero prices appear more frequently, and price spikes are very common.

3.5. EPEX-FR

The fourth dataset represents the EPEX-FR market, the day-ahead electricity market in France, which is also operated by EPEX SPOT. The dataset spans the same period as the EPEX-BE dataset, i.e. from 09.01.2011 to 31.12.2016. Besides the electricity prices, the dataset comprises the day-ahead load forecast and the day-ahead generation forecast. As before, the price data is freely obtained from the ENTSO-E transparency platform [101], and the load and generation day-ahead forecasts are freely available on the webpage of RTE [104], i.e. the *transmission system operator* (TSO) in France. Figure 3 (middle) displays the electricity price time series of the dataset; as in the EPEX-BE market, negative prices, zero prices, and spikes are very common.

3.6. EPEX-DE

The last dataset describes the EPEX-DE market, the German electricity market, which is also operated by EPEX SPOT. The dataset spans from 09.01.2012 to 31.12.2017. Besides the prices, the dataset comprises the day-ahead zonal load forecast in the TSO Amprion zone and the day-ahead wind generation forecast in the whole country. The price data is freely obtained from the ENTSO-E transparency platform [101], the zonal load day-ahead forecasts is freely available in the website of Amprion [105], and the wind forecast in the information platform of the German transmission system operators [106]. Figure 3 (bottom) displays the electricity price time series of the dataset; as can be seen, while negative and zero prices occur more often than in the other four markets, price spikes are more rare.

3.7. Training and testing periods

For each dataset, the testing period is defined as the last 104 weeks, i.e. the last two years, of the dataset. The exact dates of the testing datasets are defined in Table 1.

Table 1: Start and end dates of the testing (out-of-sample) datasets for each electricity market.

Market	Test period
Nord pool	27.12.2016 – 24.12.2018
PJM	27.12.2016 – 24.12.2018
EPEX-FR	04.01.2015 – 31.12.2016
EPEX-BE	04.01.2015 – 31.12.2016
EPEX-DE	04.01.2016 – 31.12.2017

It is important to note that, as we will argue in Section 5, selecting two years as the testing period is paramount to ensure good research practices in EPF.

Unlike the testing dataset, the training dataset cannot be defined as it will vary between different models. In, general, the training dataset will comprise any data that is prior to the data under study. However, the exact data will change depending on two concepts, i.e. calibration window and recalibration:

- While there are four years of data available for estimating the model, it might be desirable to employ only recent data, e.g. to avoid estimating effects that no longer play a role. The amount of past data employed for estimation defines the calibration window.
- The model can be estimated once and then evaluated for the full test dataset, or it can be continuously recalibrated on daily basis to incorporate the input of recent data.

For example, let us consider predicting the prices in the NP on 15.02.2017. A model using a calibration window of 52 weeks and no recalibration would employ a training dataset comprising the data between 29.12.2016 and 26.12.2016, i.e. one year prior to the start of the test period. By contrast, a model using a calibration window of 104 weeks and daily recalibration would employ the data between 18.02.2015 and 14.02.2017.



Figure 3: Electricity prices time series for three of the five datasets, i.e. EPEX-BE, EPEX-FR, and EPEX-DE, considered in the open-access benchmark dataset (Note that each dataset also includes two exogenous time series that are not plotted here). The EPEX-BE and EPEX-FR time series are similar because the EPEX-FR and EPEX-BE are highly coupled markets [43]. To keep the plots readable, the upper limit of the y-axis is below the maximum price; this only affects one spike in EPEX-FR and another one in EPEX-BE.

4. Open-access benchmark models

The second contribution of the paper is to provide a set of state-of-the-art open-source forecasting methods as an open-source `python` toolbox. As explained in Section 2.4, the LEAR [55] and the DNN [59] models are not only highly accurate but also relatively simple. Therefore, we implement these two methods and provide their code freely available as part of the proposed toolbox [60, 61]. It is important to note that the use of the proposed open-access methods is fully documented and automated so researchers can test and use them without expert knowledge.

For the sake of simplicity, the description provided here is limited to the bare minimum. For further details on the two models we refer to the original papers [55, 59].

4.1. Input features

Before describing each model, let us define the input features that are considered. Independently of the model, the available input features to forecast the 24 day-ahead prices of day d , i.e. $\mathbf{p}_d = [p_{d,1}, \dots, p_{d,24}]^\top$,

are the same:

- Historical day-ahead prices of the previous three days and one week ago, i.e. \mathbf{p}_{d-1} , \mathbf{p}_{d-2} , \mathbf{p}_{d-3} , \mathbf{p}_{d-7} .
- The day-ahead forecasts of the two variables of interest (see Section 3 for details) for day d available on day $d-1$, i.e. $\mathbf{x}_d^1 = [x_{d,1}^1, \dots, x_{d,24}^1]^\top$ and $\mathbf{x}_d^2 = [x_{d,1}^2, \dots, x_{d,24}^2]^\top$; note that the variables of interest are different for each market.
- Historical day-ahead forecasts of the variables of interest the previous day and one week ago, i.e. \mathbf{x}_{d-1}^1 , \mathbf{x}_{d-7}^1 , \mathbf{x}_{d-1}^2 , \mathbf{x}_{d-7}^2 .
- A dummy variable \mathbf{z}_d that represents the day of the week. In the case of the linear model, following the standard practice in the literature [55, 58, 77], this is a binary vector $\mathbf{z} = [z_{d,1}, \dots, z_{d,7}]^\top$ that encodes every day of the week by setting all elements to zero except the element that identifies the day of the week, e.g. $[1, 0, 0, 0, 0, 0, 0]$ represents Monday and $[0, 1, 0, 0, 0, 0, 0]$ Tuesday. In the case of the neural network, for the sake of simplicity, the day of the week is modeled with a multi-valued input $z_d \in \{1, \dots, 7\}$.

In total, we consider a total of 247 available input features for each LEAR model and 241 input features for each DNN model. It is important to note that, while the available input features are the same, each method performs a different feature selection procedure:

- Each of the LEAR models finds the optimal set of features using LASSO as an embedded feature selection, i.e. each of the models uses L1-regularization to select among the 247 features.
- For the DNN, as in the original study [59], the input features are optimized together with the hyper-parameters using the tree Parzen estimator [107] (see Section 4.3 for details).

In both cases, the feature selection is fully automated and does not require expert intervention.

4.2. The LEAR model

The first model is the LEAR model [55], a parameter-rich ARX model estimated using LASSO as an implicit feature selection approach. To enhance the model as shown by [9], the data is preprocessed with the *arc hyperbolic sine* (asinh) variance stabilizing transformation. Long-term seasonal decomposition is not considered for the sake of simplicity; particularly, while it has been shown to further improve the performance of the LEAR, we leave it out for future research.

As in [77], to further enhance the model, we recalibrated daily over different calibration window lengths: 8 weeks, 12 weeks, 3 years, and 4 years. We consider short windows (8-12 weeks) in combination with long windows (3-4 years) because it has been empirically shown to lead to better results [77]. In this context, we consider a minimum of 8 weeks as lower windows might not have enough information to correctly estimate parameter-rich models [77].

The LEAR model to predict price $p_{d,h}$ on day d and hour h is defined by:

$$\begin{aligned}
p_{d,h} &= f(\mathbf{p}_{d-1}, \mathbf{p}_{d-2}, \mathbf{p}_{d-3}, \mathbf{p}_{d-7}, \mathbf{x}_d^i, \mathbf{x}_{d-1}^i, \mathbf{x}_{d-7}^i, \boldsymbol{\theta}_h) + \varepsilon_{d,h} \\
&= \sum_{i=1}^{24} \theta_{h,i} \cdot p_{d-1,i} + \sum_{i=1}^{24} \theta_{h,24+i} \cdot p_{d-2,i} + \sum_{i=1}^{24} \theta_{h,48+i} \cdot p_{d-3,i} + \sum_{i=1}^{24} \theta_{h,72+i} \cdot p_{d-7,i} \\
&+ \sum_{i=1}^{24} \theta_{h,96+i} \cdot x_{d,i}^1 + \sum_{i=1}^{24} \theta_{h,120+i} \cdot x_{d,i}^2 + \sum_{i=1}^{24} \theta_{h,144+i} \cdot x_{d-1,i}^1 + \sum_{i=1}^{24} \theta_{h,168+i} \cdot x_{d-1,i}^2 \\
&+ \sum_{i=1}^{24} \theta_{h,192+i} \cdot x_{d-7,i}^1 + \sum_{i=1}^{24} \theta_{h,216+i} \cdot x_{d-7,i}^2 + \sum_{i=1}^7 \theta_{h,240+i} \cdot z_{d,i} + \varepsilon_{d,h} \tag{2}
\end{aligned}$$

where $\boldsymbol{\theta}_h = [\theta_{h,1}, \dots, \theta_{h,247}]^\top$ are the 247 parameters of the LEAR model for hour h . Many of these parameters become zero when (2) is estimated using LASSO:

$$\hat{\boldsymbol{\theta}}_h = \underset{\boldsymbol{\theta}_h}{\operatorname{argmin}} \{ \text{RSS} + \lambda \|\boldsymbol{\theta}_h\|_1 \} = \underset{\boldsymbol{\theta}_h}{\operatorname{argmin}} \left\{ \text{RSS} + \lambda \sum_{i=1}^{247} |\theta_{h,i}| \right\}, \quad (3)$$

where $\text{RSS} = \sum_{d=8}^{N_d} (p_{d,h} - \hat{p}_{d,h})^2$ is the sum of squares residuals, $\hat{p}_{d,h}$ the price forecast, N_d is the number of days in the training dataset, and $\lambda \geq 0$ is the *tuning* (or *regularization*) hyperparameter of LASSO. Due to the computational speed of estimating with LASSO, during every daily recalibration, the hyperparameter λ that regulates the L_1 penalty is optimized. This can be done using an *ex-ante* cross-validation procedure [108]. In this study, to further reduce the computational cost, we propose an efficient hybrid approach to perform the optimal selection of λ . See Section 4.2.2 for details.

4.2.1. Regularization hyperparameter

The hyperparameter λ of LASSO can be optimized in multiple ways, each one of them with different merits and disadvantages. A first approach is to optimize λ once and then keep it fixed for the whole test period. Although it requires very low computation costs, the limitation of this approach is that it assumes that the hyperparameter λ does not change over time. This assumption might hinder the performance of the estimator as the regularization parameter does not change even when the market might do.

A second approach is to recalibrate the hyperparameter on a periodic basis using a validation dataset. Although this method yields good results, tuning the recalibration frequency and calibration window is complicated, the computational cost is large, and the results may vary between datasets [58].

A third option is to recalibrate the hyperparameter periodically, but using *cross-validation* (CV): splitting the data into disjoint partitions, using each possible partition once as a test dataset with the remaining data as the training dataset, and selecting the hyperparameter that performs the best across all partitions [108]. Although this approach is highly accurate, its computation costs are very large.

A fourth option is to periodically update the hyperparameter but using information criteria, e.g. the *Akaike information criterion* (AIC) or the Bayesian information criterion [57, 69, 109]. As before, this involves training multiple LASSO models to compute the information criteria for each possible hyperparameter value, which in turn leads to a high computational cost.

Lastly, one can use the *least angle regression* (LARS) LASSO [110] for estimating the model instead of the coordinate descent implementation. This estimation procedure has the advantage of computing the whole LASSO solution path, which in turn allows to compute the information criteria or perform CV much faster.

4.2.2. Selecting the regularization hyperparameter

To select λ we propose a hybrid approach. On a daily basis, we estimate the hyperparameter using the LARS method with the in-sample AIC. Then, using the optimal λ obtained from the LARS method, we recalibrate the LEAR using the traditional coordinate descent implementation.

The reason for proposing this hybrid approach is that it provides a good trade-off between computational complexity and accuracy. In particular, it leverages the computational efficiency of LARS for ex-ante λ selection with the predictive performance on short calibration windows of the coordinate descent LASSO.

It is important to note that we have studied multiple approaches to select λ : (i) daily recalibration, CV, with coordinate descent; (ii) daily recalibration, CV, with LARS; (iii) daily recalibration with LARS and AIC. However, the computational cost of the first method was too high (in the same order of magnitude as the cost of the DNN model), and the accuracy of the other two was not good. By contrast, the proposed approach had a performance on par with coordinate descent LASSO using CV, but with a computational cost that was an order of magnitude lower.

4.3. The DNN model

The second model is the DNN [59], one of the most simple DL models whose input features and hyperparameters can be optimized and tailored for each case study without the need of expert knowledge.

4.3.1. Structure

The DNN is a deep feedforward neural network that contains 4 layers, employs the multivariate framework (single model with 24 outputs), is estimated using Adam, and its hyperparameters and input features are optimized using the tree Parzen estimator [107], i.e. a Bayesian optimization algorithm. The DNN model is visualized in Figure 4.

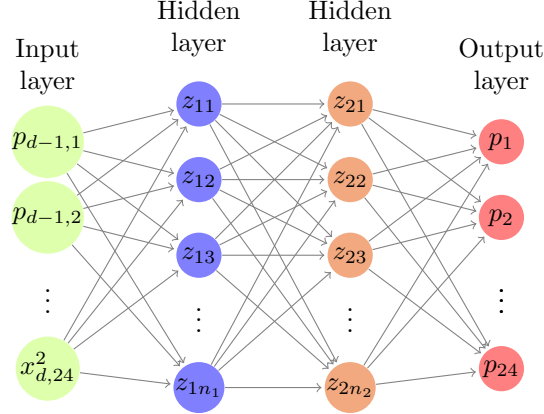


Figure 4: Visualization of a sample DNN model

4.3.2. Training dataset

For estimating the hyperparameters, the training dataset is fixed and comprises the four years prior to the testing period. For evaluating the testing dataset, the DNN is recalibrated on a daily basis using a calibration window of four years.

In all cases, the training dataset is split into a training and a validation dataset, with the latter being used for two purposes: performing early stopping [111] to avoid overfitting and optimizing hyperparameters/features. While the validation dataset always comprises 42 weeks, the split between the training and validation datasets depends on whether the validation dataset is used for hyperparameter/feature selection or for the recalibration step:

- For estimating the hyperparameters, as the validation dataset is used to guide the optimization process, the validation dataset is selected as the last 42 weeks of the training dataset. This is done to keep the training and validation datasets completely independent and to avoid overfitting⁹.
- For the testing phase, as the validation dataset is only used for early stopping, it is defined by randomly selecting 42 weeks out of the total 208 weeks employed for training. This is done to ensure that the dataset used for optimizing the DNN parameters includes up-to-date data¹⁰.

As example, let us consider the training and evaluation of a DNN in the Nord Pool market. Before evaluating the DNN, the hyperparameter and features of the DNN are optimized. For that, the employed dataset comprises the data between 01.01.2013 and 26.12.2016, of which the training dataset represents the first 166 weeks, i.e. 01.01.2013 to 07.03.2016, and the validation dataset the last 42, i.e. 08.03.2016 to 26.12.2016. During the evaluation of the model, i.e. after the hyperparameter and feature selection, the training and validation datasets comprise the last four years of data but are randomly shuffled. For

⁹Similar as it is done when splitting the dataset between the training and test dataset.

¹⁰For hyperparameter optimization, as the validation dataset represents the most recent weeks of data, the neural network is trained with data that is almost one year old. While this is not a big problem when deciding on the DNN structure, it should be avoided during testing to ensure that the DNN captures new market effects.

example, to evaluate the DNN during 15.02.2017, the training and validation datasets would represent the data between 20.02.2013 and 14.02.2017, of which 166 randomly selected weeks would define the training dataset and the remaining 42 the validation dataset.

4.3.3. Hyperparameter and feature selection

As in the original DNN paper [59], the hyperparameters and input features are optimized together using the tree-structured Parzen estimator [107], a Bayesian optimization algorithm based on sequential model-based optimization. To do so, the features are modeled as hyperparameters, with each hyperparameter representing a binary variable that selects whether a specific feature is included in the model (as explained in [43]). In more detail, to select which of the 241 available input features are relevant, the method employs 11 decision variables, i.e. 11 hyperparameters:

- Four binary hyperparameters (1-4) that indicate whether to include the historical day ahead prices \mathbf{p}_{d-1} , \mathbf{p}_{d-2} , \mathbf{p}_{d-3} , \mathbf{p}_{d-7} . The selection is done per day¹¹, e.g. the algorithm either selects all the prices \mathbf{p}_{d-j} of j days ago or it cannot select any price from day $d - j$, hence the four hyperparameters.
- Two binary hyperparameters (5-6) that indicate whether to include each of the day-ahead forecasts \mathbf{x}_d^1 and \mathbf{x}_d^2 . As with the past prices, this is done for the whole day, i.e. a hyperparameter either selects all the elements in \mathbf{x}_d^j or none.
- Four binary hyperparameters (7-10) that indicate whether to include the historical day-ahead forecasts \mathbf{x}_{d-1}^1 , \mathbf{x}_{d-1}^2 , \mathbf{x}_{d-7}^1 , and \mathbf{x}_{d-7}^2 . This selection is also done per day.
- One binary hyperparameter (11) that indicates whether to include the variable z_d representing the day of the week.

In short, 10 binary hyperparameters indicating whether to include 24 inputs each and another binary hyperparameter indicating whether to include a dummy variable.

Besides selecting the features, the algorithm also optimizes eight additional hyperparameters: 1) the number of neurons per layer, 2) the activation function, 3) the dropout rate, 4) the learning rate, 5) whether to use batch normalization, 6) the type of data preprocessing technique, 7) the initialization of the DNN weights, and 8) the coefficient for L1 regularization that is applied to each layer’s kernel.

Unlike the weights of the DNN that are recalibrated on a daily basis, the hyperparameter and features are optimized only once using the four years of data prior to the testing period. It is important to note that the algorithm runs for a number T of iterations, where at every iteration the algorithm infers a potential optimal subset of hyperparameters/features and evaluates this subset in the validation dataset. For the proposed open-access benchmark models, T is selected as 1500 iterations to obtain a trade-off between accuracy and computational requirements¹².

4.4. Ensembles

For the open-access benchmark, in order to have benchmark predictions when evaluating ensemble techniques, we also propose ensembles of LEAR and DNNs as open-access benchmarks of ensembles methods. For the LEAR, the ensemble is built as the arithmetic average of forecasts across four calibration window lengths: 8 weeks, 12 weeks, 3 years, and 4 years. For the DNN, the ensemble is built as the arithmetic average of four different DNNs that are estimated by running the hyperparameter/feature selection procedure four times. In particular, the hyperparameter optimization is asymptotically deterministic, i.e. the global

¹¹This is done for the sake of simplicity to speed up the optimization procedure of the feature selection. In particular, an alternative could be to use a binary hyperparameter for each individual historical prices; however, in most markets, that would mean using 24 as many hyperparameters as there are 24 different prices per day.

¹²It can be empirically observed that the performance of the models barely improves after 1000 iterations. Moreover, performing 1500 iteration takes approximately just one day on a regular quadcore laptop like the i7-6920HQ, a computation cost very acceptable when the algorithm has to run only once.

optimum is found for an infinite number of iterations. However, for a finite number of iterations and using a different initial random seed, the algorithm is non-deterministic and every run provides a different set of hyperparameters and features. Although each of these hyperparameter/feature subsets represent a local minimum, it is impossible to establish which of the subsets is better since their relative performance on the validation dataset is nearly identical. This effect can be explained due to the DNN being a very flexible model and thus different network architectures being able to obtain equally good results.

4.5. Software implementation

The proposed open-access models are developed in `python`: the LEAR is implemented using the `scikit-learn` library [112] and the DNN model using the `Keras` library [113]. The reason for selecting `python` is that it is one of the most widely used programming languages, especially in the context of ML and statistical inference.

5. Guidelines and best practices in EPF

As motivated in the introduction, the field of EPF suffers from several problems that prevent having reproducible research and establishing strong conclusions on what methods work best. In this section, we outline some of these issues and provide some guidelines on how to address them.

5.1. Length of the test period

A common practice in EPF is to evaluate new methods on very short test periods. The typical approach is to evaluate the method on 4 weeks of data [18, 19, 22, 24–26, 29, 30, 41, 42, 49, 51, 87, 91–96, 99], with each week representing one of the four seasons in the year. This is problematic for three reasons:

- Selecting four weeks can lead to cherry-picking the weeks where a given method excels, e.g. a method that performs bad with spikes could be evaluated in a week with fewer spikes, leading in turn to biased estimations of the forecasting accuracy. While this is an ethical issue that most researchers would avoid, establishing four weeks testing periods as the standard does facilitate the malpractice and it should be avoided.
- Assuming that the four weeks are randomly selected and no bias is introduced in the selection, it is still not possible to guarantee that these four weeks are representative of the price behavior on a whole year. Particularly, even within a given season, the price dynamics can change dramatically, e.g. during winter there are weeks with a lot of sun and wind but there are also weeks without them. Therefore, picking only a week per season rarely represents the average performance of a forecaster in a give dataset.
- There are situations in the electrical grid that do not occur very often but that can have a very large effect on electricity prices, e.g. when several power plants are under maintenance at the same time. Forecasting methods need to be evaluated under those conditions to ensure that they are also accurate under extreme events. By selecting four weeks most of these effects are neglected.

To avoid this problem, we recommend using a minimum of one year as a testing period. This ensures that forecasting methods are evaluated considering the complete set of effects that take place during the year. To guarantee that all researchers have access to this type of data, the open-access benchmark dataset that we propose contains data from several markets and employs a testing period of two years. In addition, the open-access benchmark can be directly accessed using the proposed `epftoolbox` library [60, 61].

5.2. Benchmark models

A second issue with many EPF publications is that new methods are not compared with well-established methods [14, 16, 18–21, 23, 25, 27, 34, 36, 42, 46, 48–50, 78, 90, 93, 95, 100] or resort to comparisons using either outdated methodologies or simplified methods [13, 15, 22, 24, 26, 28–30, 37, 41, 44, 45, 47, 51, 52, 85, 86, 91, 92].

This poses a problem since it becomes very hard to establish which algorithms work best and which ones do not. To address this issue, we recommend using well-established state-of-the-art open-source methods and a common benchmark dataset. With that in mind, we have provided and make freely available an open-access benchmark dataset comprising 5 markets (as described in Section 3), and we have implemented, thoroughly tested, and made freely available two state-of-the-art forecasting methods (as described in Section 4) and their day-ahead predictions for all 5 datasets over a period of two years (as described in Section 6). Additionally, we have implemented all these resources in an easy-to-use toolbox [60] and built an adequate documentation [61].

5.3. Open-access

A third issue in the field of EPF is that datasets are usually not made publicly available and the code of the proposed methods is not shared. This poses four obvious problems:

- Research cannot be reproduced as data is not available. This goes against one of the main principles of science as all research should be reproducible.
- The progress of EPF research is hindered since it is hard to establish which methodologies work well. Consequently, researchers spend unnecessary time re-evaluating methodologies that have been evaluated already.
- Comparing new methods with published ones becomes very challenging because researchers have to re-implement methods from the literature. As a result, comparisons with state-of-the-art methods are often avoided, and new methods are usually compared with simple and easy-to-implement methods.
- When new methods are proposed, they cannot be compared with published methods under the same circumstances. This leads to comparisons under different conditions and opens up the door to wrong implementations of the original methods, which in turn leads to results that are not correct.

As these problems are critical, we directly try to address them by providing an open-access benchmark/toolbox comprising five datasets, two state-of-the-art methods, and a set of day-ahead forecasts of the latter two methods. In addition, we encourage researchers in EPF to share the developed codes and to either share their datasets or use an open-access benchmark dataset.

5.4. Evaluation metrics for point forecasts

In the field of EPF, the most widely used metrics to measure the accuracy of point forecasts are the *mean absolute error* (MAE), the *root mean square error* (RMSE), and the *mean absolute percentage error* (MAPE):

$$\text{MAE} = \frac{1}{24 N_d} \sum_{d=1}^{N_d} \sum_{h=1}^{24} |p_{d,h} - \hat{p}_{d,h}|, \quad (4)$$

$$\text{RMSE} = \sqrt{\frac{1}{24 N_d} \sum_{d=1}^{N_d} \sum_{h=1}^{24} (p_{d,h} - \hat{p}_{d,h})^2}, \quad (5)$$

$$\text{MAPE} = \frac{1}{24 N_d} \sum_{d=1}^{N_d} \sum_{h=1}^{24} \frac{|p_{d,h} - \hat{p}_{d,h}|}{|p_{d,h}|}, \quad (6)$$

where $p_{d,h}$ and $\hat{p}_{d,h}$ respectively represent the real and forecasted price on day d and hour h , and N_d is the number of days in the out-of-sample test period, i.e. in the test dataset.

Since absolute errors are hard to compare between different datasets, the MAE and RMSE are not always very informative. Moreover, since electricity costs and profits are often linearly dependent on the electricity prices, metrics based on quadratic errors, e.g. RMSE, are hard to interpret and do not accurately represent the underlying problem of most forecasting users. In particular, in most electricity trade applications, the underlying risk, profits, and costs depend linearly on the price and on the forecasting errors. Hence, linear metrics represent better than quadratic metrics the underlying risks of forecasting errors.

Similarly, since MAPE values become very large with prices close to zero (regardless of the actual absolute errors), the MAPE is usually dominated by the periods of low prices and is also not very informative. While the *symmetric mean absolute percentage error* (sMAPE) defined¹³ as:

$$\text{sMAPE} = \frac{1}{24 N_d} \sum_{d=1}^{N_d} \sum_{h=1}^{24} 2 \frac{|p_{d,h} - \hat{p}_{d,h}|}{|p_{d,h}| + |\hat{p}_{d,h}|} \quad (7)$$

solves some of these issues, it has (as any metric based on percentage errors) a statistical distribution with undefined mean and infinite variance [115].

5.4.1. Scaled errors

In this context, several studies advocate for the use of scaled errors [5, 115, 116], where a scaled error is simply the MAE scaled by the in-sample MAE of a naive forecast. A scaled error has the nice interpretation of being lower/larger than one if it is better/worse than the average naive forecast evaluated in-sample.

A metric based on this concept is the *mean absolute scaled error* (MASE), and in the context of one-step ahead forecasting is defined as [115]:

$$\text{MASE} = \frac{1}{N} \sum_{k=1}^N \frac{|p_k - \hat{p}_k|}{\frac{1}{n-1} \sum_{i=2}^n |p_i^{\text{in}} - p_{i-1}^{\text{in}}|}, \quad (8)$$

where p_i^{in} is the i^{th} price in the in-sample, i.e. training, dataset (note that in EPF $i = 24d + h$), p_{i-1}^{in} is the one-step ahead naive forecast of p_i^{in} , i.e. \hat{p}_i^{in} , N is the number of out-of-sample (test) datapoints, and n the number of in-sample (training) datapoints. For seasonal time series, the MASE may be defined using the MAE of a seasonal naive model in the denominator [5, 116].

5.4.2. Relative measures

While scaled errors do indeed solve the issues of more traditional metrics, they have other associated problems that make them unsuitable in the context of EPF:

1. As MASE depends on the in-sample dataset, forecasting methods with different calibration windows will naturally have to consider different in-sample datasets. As a result, the MASE of each model will be based on a different scaling factor and comparisons between models cannot be drawn.
2. The same argument applies to models with and without rolling windows. The latter will use a different in-sample dataset at every time point while the former will keep the in-sample dataset constant.
3. In ensembles of models with different calibration windows, the MASE cannot be defined as the calibration window of the ensemble is undefined.
4. Drawing comparisons across different time series is problematic as electricity prices are not stationary. For example, an in-sample dataset with spikes and an out-of-sample dataset without spikes will lead to a smaller MASE than if we consider the same market but with the in-sample/out-sample datasets reversed.

¹³Note, that there are multiple versions of sMAPE, here we consider the most sensible one according to [114].

To solve these issues, we argue that a better metric is the *relative* MAE (rMAE) [115]. Similar to MASE, rMAE normalizes the MAE by the MAE of a naive forecast. However, instead of considering the in-sample dataset, the naive forecast is built based on the out-of-sample dataset. In the context of EPF, rMAE is defined as:

$$\text{rMAE} = \frac{\frac{1}{24N_d} \sum_{d=1}^{N_d} \sum_{h=1}^{24} |p_{d,h} - \hat{p}_{d,h}|}{\frac{1}{24N_d} \sum_{d=1}^{N_d} \sum_{h=1}^{24} |p_{d,h} - \hat{p}_{d,h}^{\text{naive}}|}, \quad (9)$$

where the $\frac{1}{24N_d}$ factor cancels out in the numerator and the denominator. There are three natural choices for the naive forecasts:

- $\hat{p}_{d,h}^{\text{naive},1} = p_{d-1,h}$,
- $\hat{p}_{d,h}^{\text{naive},2} = p_{d-7,h}$,
- $\hat{p}_{d,h}^{\text{naive},3} = \begin{cases} p_{d-1,h}, & \text{if } d \text{ is Tue, Wed, Thu, or Fri,} \\ p_{d-7,h}, & \text{if } d \text{ is Sat, Sun, or Mon.} \end{cases}$

In the context of EPF, rMAE using $\hat{p}_{d,h}^{\text{naive},2} = p_{d-7,h}$ is arguably the best choice for two reasons: (i) it is easier to compute than the one based on $\hat{p}_{d,h}^{\text{naive},3}$ and, unlike the rMAE based on $\hat{p}_{d,h}^{\text{naive},1}$, it captures weekly effects; (ii) given a set of forecasting models, the relative ranking of the accuracy of the models is independent from the naive benchmark used (see last paragraph of this subsection for an explanation). Hence, for the remainder of the article we will use rMAE to explicitly refer to the rMAE based on $\hat{p}_{d,h}^{\text{naive},2}$. It is important to note that, similar to rMAE, one could also define the *relative* RMSE (rRMSE) by dividing the RMSE of each forecast by the RMSE of a naive forecast.

Since the dependence on the in-sample dataset is removed, using a rolling window is no longer a problem as the out-of-sample dataset stays the same. Similarly, models with different calibration windows can be compared and the rMAE of ensembles is properly defined. Moreover, as the metric is normalized by the MAE of a naive forecast for the same sample, the problem with drawing conclusions in non-stationary time series is mitigated. As before, we can also define the rMAE for seasonal time series:

Due to its better properties, rMAE should always be used to evaluate new methods in EPF. In particular, while it can be used in conjunction with other metrics, it is important to include and employ rMAE to obtain more fair evaluations and comparisons.

With that in mind, the accuracy of the open-access models in the open-access benchmark dataset is computed considering rMAE, sMAPE, MAPE, MAE, and RMSE. Then, an analysis of the different metrics is provided (see Section 6.4.2). Finally, the forecasts themselves are provided as csv files so that the accuracy results can be updated in case more adequate metrics are developed in the future.

As a final remark, let us note that, given a set of forecasting models, the relative ranking of the accuracy of the models is independent from the naive benchmark used for the rMAE or MASE. Changing it simply changes the denominator but preserves the numerator, and since the change in the denominator is the same across all methods, the relative ranking is preserved. Furthermore, as the numerator is the MAE, it follows that the ranking based on the rMAE or MASE will be the same as that based on the MAE.

5.5. Statistical testing

While using adequate metrics to compare the accuracy of the forecasts is important, it is also necessary to analyze whether any difference in accuracy is statistically significant. This is paramount to conclude whether the difference in accuracy does really exist and is not simply due to random differences between the forecasts. Despite its importance, the use of statistical testing has been downplayed in the EPF literature [5]. In particular, most publications only compare the accuracy in terms of an error metric and do not

analyze the statistical significance of the accuracy differences. This trend needs to be corrected in order to compare forecasting approaches with statistical rigor. Particularly, new studies need to ensure that:

- Any new method is compared against well-established methods using a statistical test.
- The forecasts of the proposed methods are provided as open-access datasets. This ensures that, when new models are proposed, the difference in accuracy with the published methods can be analyzed in terms of statistical testing.

To facilitate statistical testing, we include in the proposed open-source `epftoolbox` library [60, 61] the two most widely used statistical tests in EPF, i.e. the Diebold-Mariano and the Giacomini-White tests.

5.5.1. The Diebold-Mariano test

The Diebold-Mariano (DM) test [117] is probably the most commonly used tool to evaluate the significance of differences in forecasting accuracy. It is an asymptotic z -test of the hypothesis that the mean of the *loss differential* series:

$$\Delta_{d,h}^{A,B} = L(\varepsilon_{d,h}^A) - L(\varepsilon_{d,h}^B) \quad (10)$$

is zero, where $\varepsilon_{d,h}^Z = p_{d,h} - \hat{p}_{d,h}$ is the prediction error of model Z for day d and hour h , and $L(\cdot)$ is the loss function. For point forecasts, we usually take $L(\varepsilon_{d,h}^Z) = |\varepsilon_{d,h}^Z|^p$ with $p = 1$ or 2 , which corresponds to the absolute and squared losses, respectively; for probabilistic forecasts, $L(\cdot)$ may be any strictly proper scoring rule, in particular the pinball loss, the continuous ranked probability score, or the energy score [6, 65, 66]. Given the loss differential series, we compute the statistic:

$$\text{DM} = \sqrt{N} \frac{\hat{\mu}}{\hat{\sigma}}, \quad (11)$$

where $\hat{\mu}$ and $\hat{\sigma}$ are the sample mean and standard deviation of $\Delta_{d,h}^{A,B}$, respectively, and N is the length of the out-of-sample test period. Under the assumption of covariance stationarity of $\Delta_{d,h}^{A,B}$, the DM statistic is asymptotically standard normal, and one- or two-sided asymptotic tail probabilities can be easily computed.

It is important to note three things. Firstly, the DM test is model-free, i.e. it compares forecasts (of models), not models themselves. Secondly, although in the standard formulation [117] the DM test compares forecasts via the null hypothesis of the expected loss differential being zero, it is more informative to compute the p -values of two one-sided tests:

1. with the null hypothesis $H_0 : E(\Delta_{d,h}^{A,B}) \leq 0$,
2. with the alternative hypothesis null $H_1 : E(\Delta_{d,h}^{A,B}) \geq 0$.

The lower the p -value¹⁴, i.e. the closer it is to zero, the more the observed data is inconsistent with the null hypothesis. If the p -value is less than the commonly accepted level of 5%, the null hypothesis is typically rejected. In the DM test, this means that the forecasts of model B are significantly more accurate than those of model A.

Thirdly, the DM test requires (only) that the loss differential be covariance stationary.¹⁵ This may not be satisfied by forecasts in day-ahead markets, since the predictions for all 24 hours of the next day are computed at the same time, using the same information set. Hence, following [65], we recommend two variants of the DM test in the context of day-ahead EPF:

- a *univariate* variant with 24 independent tests performed¹⁶, one for each hour of the day, and comparisons based on the number of hours for which the predictions of one model are significantly better than those of another, i.e. the number of hours for which the null hypothesis is rejected,

¹⁴Recall, that the p -value is the probability of obtaining results (in our case – loss differentials) at least as large as the ones actually observed, assuming that the null hypothesis is correct.

¹⁵Actually covariance stationarity is sufficient but may not be strictly necessary [118].

¹⁶We assume that a day-ahead market has 24 prices. For markets with prices every half hour, the univariate variant comprises 48 independent tests.

- a *multivariate* variant with the test performed jointly for all hours using the ‘daily’ or multivariate loss differential series:

$$\Delta_d^{A,B} = \|\varepsilon_d^A\|_p - \|\varepsilon_d^B\|_p, \quad (12)$$

where ε_d^Z is the 24-dimensional vector of prediction errors of model Z for day d , $\|\varepsilon_d^Z\|_p = (\sum_{h=1}^{24} |\varepsilon_{d,h}^Z|^p)^{1/p}$ is the p -th norm of that vector with $p = 1$ or 2 .

The univariate version of the test has the advantage of providing a deeper analysis as it indicates which forecast is significantly better for which hour of the day [6, 55, 59, 66, 119, 120]. The multivariate version, introduced in [57], enables a better representation of the results as it summarizes the comparison in a single p -value, which can be conveniently visualized using heat maps arranged as chessboards [9, 10, 58, 76], see Figure 5.

5.5.2. The Giacomini-White test

In some of the more recent EPF studies [77, 121, 122], the DM test has been replaced by the Giacomini-White (GW) test [123] for *conditional predictive ability*. The latter is preferred because it can be regarded as a generalization of the DM test for *unconditional predictive ability*: while both tests can be used for nested and non-nested models¹⁷, only the GW test accounts for parameter estimation uncertainty through ‘conditioning’ [65].

Like the DM test, also the GW test has two variants in day-ahead EPF – the univariate and the multivariate. Without loss of generality, let us focus on the latter. It starts by building a multivariate loss differential series, see (12), for a pair of forecasts (of models A and B). Next, the test considers the following regression:

$$\Delta_d^{A,B} = \phi' X_{d-1} + \epsilon_d, \quad (13)$$

where X_{d-1} contains elements from the information set on day $d-1$, i.e. a constant and lags of $\Delta_d^{A,B}$. Note that $\epsilon_d \neq \varepsilon_d^Z$, i.e. ϵ_d is not the 24-dimensional vector of prediction errors for day d and model Z but simply an error term in the regression. Also note that using this notation the DM test can be written as [124]:

$$\Delta_d^{A,B} = \mu + \epsilon_d, \quad (14)$$

i.e. with X_{d-1} containing just a constant. Finally, like for the DM test, to check the significance of differences in forecasting accuracy, the p -values of two one-sided tests can be computed. The interpretation and possible visualization (see Figure 5) are analogous to that of the DM test.

5.6. Recalibration

An issue with many EPF studies is that forecasting models are not recalibrated. Instead, they are often estimated once using the training dataset and directly evaluated in the whole test dataset. This is problematic as it does not represent real-life conditions where forecasting models are retrained (often on a daily basis) to account for the latest market information.

To have models that are evaluated in realistic conditions, they need to be retrained considering the new incoming flow of market information. As an example, for the day-ahead market, a forecasting model should be retrained on a daily basis as new information is available. Considering a testing period of a year, this means that a realistic evaluation requires estimating the forecasting model 365 times.

¹⁷This holds as long as the calibration window does not grow with the sample size [124]. This is satisfied for rolling windows, but not for extended calibration windows.

5.7. *Ex-ante hyperparameter optimization*

A common issue in the current EPF literature is that the hyperparameter selection is often either done ex-post [49, 51, 125–128] or its details are not sufficiently explained [13, 21, 37, 48, 79, 82, 91–93, 96, 99]. As an example, when models based on neural networks are proposed, the details on how the number of neurons are selected are usually not provided. In other cases, while the approach is provided, it is often based on analyzing different configurations of neurons using the test dataset and selecting the one that works best, i.e. ex-post hyperparameter selection.

Not providing enough details on how hyperparameters are selected is an obvious problem as it prevents reproducing research. Similarly, performing hyperparameter optimization ex-post leads to overfitting the test dataset, i.e. the model is partially optimized using the same dataset used for evaluating the model, and it grants the model an unfair and non-existent advantage over other models.

To prevent this, the selection of hyperparameters should be explicitly explained and always performed ex-ante using a validation dataset. With that motivation, for the open-access methods proposed, not only do we explain how the hyperparameters are obtained, but we also provide within the toolbox [60, 61] a module for hyperparameter selection and the files containing the results of the hyperparameter optimization of the current study.

5.8. *Computation time*

An even more common problem is the fact that new models are very rarely compared in terms of their computational requirements [19, 20, 22, 24, 32, 37, 41, 42, 51, 90–92, 94–97, 100]. Although a model might be marginally better than another, it might not be worthwhile to deploy it in a practical application if its computational requirements are much larger. Particularly, higher computational requirements might pose two problems:

1. As mentioned before, forecasting models should ideally be recalibrated on a daily basis. Hence, a forecasting method is only suitable if its computational time allows this recalibration to take place. In this context, the maximum available time for estimating a model will depend on each electricity market but, as a rule of thumb, it can be argued that any model that requires more than 30 min or 1 h will unlikely be suitable for forecasting prices in the spot markets.
2. Besides recalibration, the second issue with computation time is its cost. If the computational capabilities are too large, the benefits of using a marginally better forecast might be lower than the cost of running the forecasting model on a much more expensive computer.

Hence, when new forecasting models are proposed, we argue that it is very important to provide their computation times. Moreover, we also argue that for a model to be better than the existing methods, it does not necessarily have to be the most accurate one. Instead:

1. If its computational time is large, i.e. in the order of minutes, the model should indeed be more accurate than all state-of-the-art models, e.g. DNNs.
2. If its computational time is small, i.e. in the order of seconds, the model should be more accurate than the state-of-the-art models with low computational requirements, e.g. LEAR.

In this article, we provide an analysis of the computational requirements of the proposed open-access models so other researchers can easily make such comparisons.

5.9. *Reproducibility*

Another related issue is that some studies lack enough details to replicate the research. Missing details vary from study to study but the four most common are:

1. the dataset used for testing and evaluation is not defined [31–37];
2. the dataset used for training is not defined [21, 33, 35, 41, 42];
3. the inputs of the model are unclear [35, 36, 38–40];
4. the selection of hyperparameters is unclear [13, 21, 37, 48, 79, 82, 91–93, 96, 99].

To correct this, future EPF papers should provide enough details to allow replication and reviewers should verify that all necessary details of the employed datasets are always provided.

5.10. Data contamination

Another recurrent issue in the EPF literature is data contamination, which appears when part of the training dataset is used for testing. Particularly, when working with time series data the test dataset should always comprise the last part of the dataset to avoid data contamination. If this is not done, the models can overfit the testing dataset and their accuracy can be overestimated.

Despite the importance of correctly separating the training/validation dataset from the testing dataset, some studies in EPF:

1. Do not specify the split between the training, validation, and test datasets [21, 31–37, 41, 42]. If the datasets are not specified, it is not possible to know whether data contamination occurs.
2. Randomly sample the test dataset from the full dataset [129–132], e.g. in a dataset comprising a year of data randomly selecting 4 weeks for testing and the remaining data for training.
3. Have a partial or total overlap between the training/validation dataset and the testing dataset [51, 125, 126, 133], e.g. by performing hyperparameter optimization ex-post.

To correct this issue, it is important that any future research in EPF ensures that: 1) the split between the datasets is correctly described; 2) the test dataset does never overlap with the training or validation datasets; 3) the test dataset is always selected as the last segment of the full dataset.

5.11. Software toolboxes

A less pressing yet relevant issue is the use of state-of-the-art software toolboxes. When comparing new methods with methods from the literature, the latter should be modeled using adequate toolboxes. Particularly, it is important to use toolboxes that are continuously updated as implementing methods using outdated libraries leads to unfair evaluations.

For example, in the context of neural networks, there are several open-source state-of-the-art toolboxes [113] that are continuously updated and that grant access to the latest development in the field of DL. Yet, in the context of EPF, new methods are often compared with neural networks that are modeled using the `MATLAB` toolbox [32, 38, 41, 42, 49, 91, 92, 94, 130], a toolbox that for many years was outdated and did not include many of the neural network developments that are critical in EPF, e.g. state-of-the-art activation functions or stochastic gradient descent algorithms [59]. As a result, many of the existing comparisons in EPF are based on evaluations where the accuracy of neural networks might be underestimated.

Besides using state-of-the-art software toolboxes, e.g. the `python` library `keras` for deep learning, it is also important to employ (whenever possible) free-to-access libraries so that research can be replicated by anyone.

5.12. Combining forecasts

As a final guideline, it is important to indicate the importance of ensembles in the context of EPF. In general, although exceptions exist [134], combining different models leads to a higher accuracy [77, 120] and it is thus a good idea to build forecasts based on multiple models. However, as even the arithmetic average improves the accuracy of individual models, new ensemble techniques should be studied in comparison with other ensemble techniques, i.e. as done in [120], and not simply w.r.t. the individual models.

To maximize the forecasting accuracy, it is important to employ diverse forecasts [134], e.g. forecasts generated using different data or different models. For EPF, the former can be achieved by considering models trained using different calibration window lengths [76, 121] and the latter using different modeling techniques or different sets of hyperparameters. To further maximize the performance, the number of models used in the ensemble should be limited [134], e.g. 4–10, especially in the case of heavy-tailed data for which large ensembles tend to contain outliers more often, resulting in less accurate forecasts.

With that in mind, as part of the open-access benchmark and toolbox [60, 61], we also propose a series of simple ensemble techniques. Particularly, as explained in Section 4, we provide an ensemble of four `LEAR` models that are estimated over different calibration windows and combined using a simple arithmetic average and another ensemble using four `DNNs` that are estimated for different hyperparameters and combined using the arithmetic average.

6. Evaluation of state-of-the-art methods

In this section, we present the results of the open-source benchmark methods for all five datasets. For the sake of clarity, we divide the section into two parts respectively comprising the results for the error metrics and the results for statistical testing.

6.1. Accuracy metrics

We first start by presenting the results of the open-access benchmark models in terms of accuracy metrics.

6.1.1. Individual models

Table 2 compares the performance of the two individual models and their 4 variations in terms of rMAE, MAE, MAPE, SMAPE, and RMSE. The LEAR model is displayed for 4 different calibration windows representing 56, 84, 1092, and 1456 days, i.e. 8 weeks, 12 weeks, 3 years, and 4 years. The four DNNs are obtained by performing the hyperparameter/feature optimization process four times and using the best hyperparameter/feature selection of every run (see Sections 4.3.3 and 4.4 for further details)¹⁸. Several observations can be made:

- The MAPE seems an unreliable metric as it completely disagrees with the other three linear metrics and the quadratic metric. In particular, while the rMAE, MAE, and SMAPE agree on what the best model is in all the cases, the MAPE almost never does so. This unreliability can be further seen in the German market: while the MAPE and SMAPE metrics usually have similar orders of magnitude, in the case of the German market the MAPE is approximately 10 times larger. This effect is due to prices in Germany being negative and very close to 0, leading in turn to very large MAPE values that bias the average MAPE.
- The DNN models seem to be more accurate than the LEAR models. Particularly, in terms of linear metrics, for the Nord Pool, PJM, and Belgian markets, the four DNN models perform better than all four LEAR models, and in the case of the German and French markets the best model is a DNN.
- Although the RMSE displays different results, this is expected as the metric is based on quadratic errors and not linear ones. Nonetheless, while the RMSE does display slightly different results, it still shows the superiority of the DNN model: even though the DNN is estimated to minimized absolute errors (unlike LEAR), the DNN is better in 3 of the 5 datasets. Moreover, even though the DNN seems to be worse in two markets, the RMSE metric does not correctly represent the underlying problem (see Sections 5.4 and 6.4.1) and it can be argued that it is not the best metric to assess the performance of EPF models.

6.1.2. Ensembles

The results for the ensemble methods are listed in Table 3, which compares the performance of the two ensemble models and the best DNN and LEAR models in terms of the rMAE metric, i.e. arguably the most reliable metric. From the table, several observations can be made:

- As already argued in Section 5.12, combining models usually improves the accuracy. Particularly, the ensemble of DNNs is better than the best individual DNN model for all four markets and for all reliable metrics. Similarly, the ensemble of LEAR models is better than the best individual LEAR model for all markets and reliable metrics. The exception to this observation are the MAPE and RMSE metrics but, as already noted, MAPE is an unreliable metric and RMSE does not correctly represent the underlying problem of EPF.
- In terms of rMAE, the ensemble of DNNs is the most accurate model across all markets and metrics, with the exception of the German market where the ensemble of LEAR performs slightly better.

¹⁸Note that, for the sake of simplicity, the features and hyperparameter selection for each model are not provided. However, they can be obtained from the website [60] accompanying this study

Table 2: Comparison between the two individual state-of-the-art open-source methods in terms of rMAE, MAE, MAPE, sMAPE, and RMSE. Each of the two methods is listed for four different configurations. The gray cells represent the best model for a given metric.

		DNN ₁	DNN ₂	DNN ₃	DNN ₄	LEAR ₅₆	LEAR ₈₄	LEAR ₁₀₉₂	LEAR ₁₄₅₆
NP	rMAE	0.471	0.415	0.437	0.438	0.475	0.472	0.482	0.481
	MAE	1.946	1.717	1.808	1.812	1.964	1.952	1.993	1.990
	MAPE [%]	6.04	5.46	5.93	5.85	6.34	6.36	6.10	6.14
	sMAPE [%]	5.59	5.00	5.22	5.26	5.66	5.62	5.64	5.66
	RMSE	3.579	3.341	3.502	3.596	3.671	3.664	3.605	3.604
PJM	rMAE	0.475	0.475	0.473	0.467	0.550	0.548	0.490	0.489
	MAE	3.005	3.008	2.995	2.956	3.477	3.467	3.098	3.095
	MAPE [%]	28.87	29.74	29.87	29.10	32.52	32.34	30.28	30.24
	sMAPE [%]	11.99	11.93	11.89	11.81	13.68	13.58	12.33	12.54
	RMSE	5.121	5.333	5.023	4.820	5.718	5.709	5.264	5.142
EPEX BE	rMAE	0.608	0.600	0.597	0.608	0.682	0.669	0.649	0.653
	MAE	6.181	6.094	6.066	6.173	6.924	6.798	6.594	6.634
	MAPE [%]	24.83	28.69	24.08	30.46	32.88	32.34	26.26	22.64
	sMAPE [%]	14.40	14.35	13.87	14.25	16.20	15.95	16.87	17.29
	RMSE	16.577	15.879	16.304	16.488	16.371	16.291	16.458	16.420
EPEX FR	rMAE	0.576	0.572	0.562	0.585	0.638	0.624	0.580	0.597
	MAE	4.223	4.193	4.118	4.292	4.681	4.575	4.250	4.378
	MAPE [%]	15.75	16.52	15.13	15.55	19.03	18.09	14.95	14.90
	sMAPE [%]	12.06	12.03	11.65	11.96	13.43	13.28	13.25	14.05
	RMSE	12.036	11.850	11.414	12.455	11.732	10.759	11.337	11.462
EPEX DE	rMAE	0.446	0.463	0.463	0.454	0.506	0.499	0.450	0.451
	MAE	4.071	4.222	4.223	4.148	4.619	4.555	4.108	4.118
	MAPE [%]	103.45	118.91	107.36	116.43	129.76	133.58	128.30	124.19
	sMAPE [%]	15.99	16.36	16.45	16.20	17.60	17.49	16.98	17.05
	RMSE	7.225	7.540	7.547	7.427	8.122	7.923	6.996	6.987

6.2. Statistical Testing

In this section, we present the results of the open-access benchmark models in terms of the statistical tests. For the sake of simplicity, we present together the results for individual methods and ensembles. The results are based on the multivariate GW test using the L_1 norm in (12), i.e. with the following loss differential series:

$$\Delta_d^{A,B} = \sum_{h=1}^{24} |\varepsilon_{d,h}^A| - \sum_{h=1}^{24} |\varepsilon_{d,h}^B|. \quad (15)$$

While squared losses could also be used, we do not consider them here because absolute errors better represent the underlying problem in EPF, see Section 6.4.1 for a discussion.

In Figure 5 we display the results for the five markets. More precisely, we use heat maps arranged as chessboards to indicate the range of the obtained p -values. The closer they are to zero (dark green) the more significant is the difference between the forecasts of a model on the X-axis (better) and the forecasts of a model on the Y-axis (worse). For instance, for the EPEX-DE market the first row is green indicating

Table 3: Comparison between the ensembles of the state-of-the-art open-source methods in terms of rMAE, MAE, MAPE, and sMAPE. The comparison also includes, for each market, the best individual performing DNN and LEAR model in terms of rMAE and MAE, i.e. the two most reliable metrics. The gray cells represent the best model for a given metric.

		DNN Ensemble	LEAR Ensemble	Best ¹⁹ DNN	Best LEAR
NP	rMAE	0.403	0.420	0.415	0.472
	MAE	1.667	1.738	1.717	1.952
	MAPE [%]	5.38	5.53	5.46	6.36
	sMAPE [%]	4.85	5.01	5.00	5.62
	RMSE	3.333	3.362	3.341	3.604
PJM	rMAE	0.439	0.476	0.467	0.489
	MAE	2.779	3.013	2.956	3.095
	MAPE [%]	28.66	30.13	29.10	30.24
	sMAPE [%]	11.22	11.98	11.81	12.54
	RMSE	4.637	5.127	4.820	5.142
EPEX BE	rMAE	0.573	0.604	0.597	0.649
	MAE	5.821	6.140	6.066	6.594
	MAPE [%]	26.11	20.72	24.08	26.26
	sMAPE [%]	13.33	14.55	13.87	16.87
	RMSE	16.127	15.974	15.879	16.371
EPEX FR	rMAE	0.533	0.543	0.562	0.58
	MAE	3.910	3.980	4.118	4.250
	MAPE [%]	14.77	14.68	15.13	14.95
	sMAPE [%]	10.98	11.57	11.65	13.25
	RMSE	11.738	10.676	11.414	10.759
EPEX DE	rMAE	0.438	0.433	0.446	0.450
	MAE	3.998	3.955	4.071	4.108
	MAPE [%]	106.67	122.41	103.45	128.30
	sMAPE [%]	15.68	15.75	15.99	16.98
	RMSE	7.278	7.079	7.225	6.987

that the forecasts of LEAR₂₈ are significantly outperformed by those of all other models. We can observe that:

- For all markets except the EPEX-DE the last column is green indicating that the forecasts of the DNN ensemble significantly outperform those of all other models. The only exception is the LEAR ensemble in the German market: in this case, the forecasts of the two models are not statistically different.
- The forecasts of LEAR_{ens} are statistically significantly better than those of all individual LEAR models. Together with the previous observation, i.e. the superiority of the DNN ensemble, this shows that the predictions of ensemble models usually improve upon the forecasting accuracy of individual methods.
- In one dataset (EPEX-BE), the forecasts of all the individual DNN methods are statistically significantly better than those of the individual LEAR models. In the remaining four datasets, the forecasts of the individual DNN models are significantly better than those of 2 to 3 individual LEAR models.

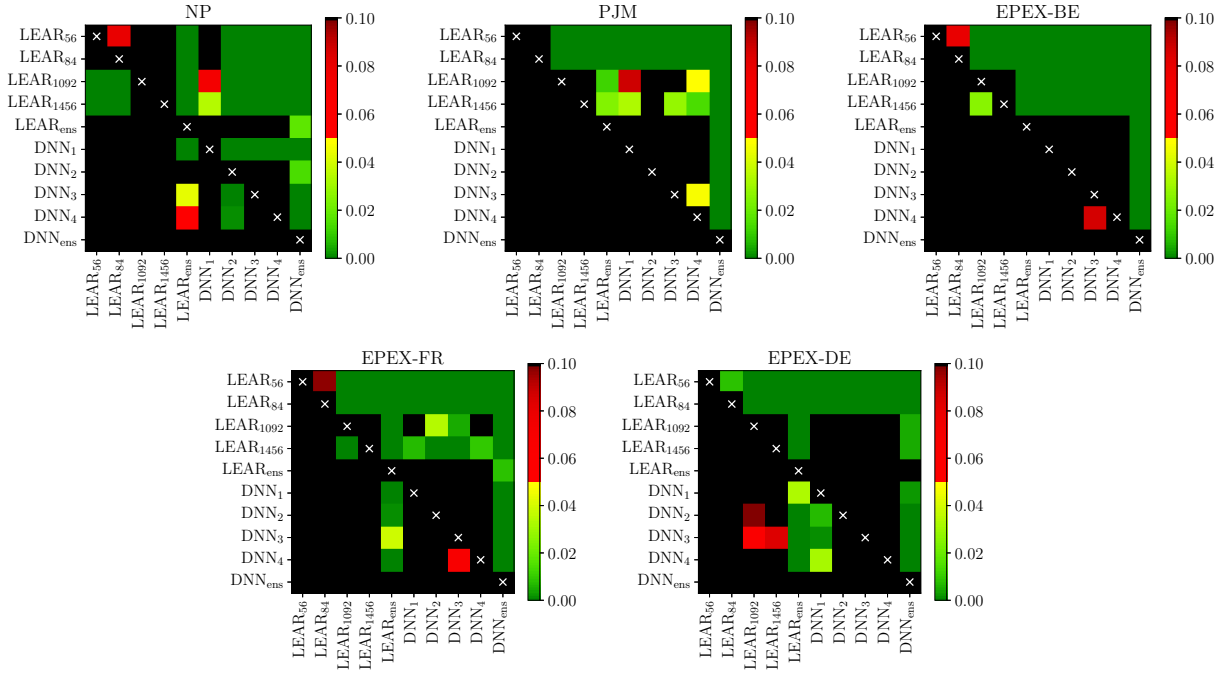


Figure 5: Results of the GW test with the multivariate loss differential series (15) for the eight individual models and the two ensembles. A heat map is used to indicate the range of the obtained p -values for each of the five markets. The closer the p -values are to zero (dark green), the more significant the difference is between the forecasts of a model on the X-axis (better) and the forecasts of a model on the Y-axis (worse). Black color indicates p -values above the color map limit, i.e. p -values larger or equal than 0.10.

- The forecasts of the individual LEAR models are never significantly better than those of the individual DNN models. Overall, it seems that forecasts based on DNNs are more likely to obtain significantly better results; this is particularly true for the DNN ensemble.

6.3. Computation time

As described in Section 5.8, besides comparing the predictive accuracy, it is also necessary to analyze the computation time of the forecasting methods. Table 4 lists a comparison of the computation time required for estimating the models considered, i.e. the time required to recalibrate each model on a daily basis. As the computation time is non-deterministic, its value is given as a range. These data were obtained using a regular laptop quad core CPU, i.e. the i7-6920HQ.

Table 4: Computation time that each benchmark model requires to perform a daily recalibration.

	Time
LEAR	1–10 seconds
LEAR Ensemble	20–25 seconds
DNN	2–5 minutes
DNN Ensemble	8–20 minutes

As can be observed, although the LEAR model performs slightly worse than the DNN model, its computation time is 30 to 100 times lower; particularly, when considering the maximum computation time of both methods, the LEAR model is 50 times faster.

6.4. Discussion and remarks

In this section, we provide some final remarks behind the motivation of the metrics employed, we briefly analyze the influence of the different metrics considered, and provide a discussion on comparing new models.

6.4.1. Absolute vs. squared errors

Throughout the text, we have mostly considered accuracy metrics based on absolute/linear errors, i.e. metrics that evaluate the accuracy of predicting the median of the distribution. Since the LEAR model is estimated by minimizing squared errors, thus leading to forecasts of the mean [116], one could argue that a metric/test based on squared errors should be preferred. While the argument has some merits, we focused on absolute metrics for three reasons:

- The metric used to evaluate the accuracy should be the one that better represents the underlying problem. In the case of EPF, since the cost of purchasing electricity is linear, linear metrics are arguably the best to quantify the risk associated with forecasting errors.
- While we provided the RMSE results, they are qualitatively the same as for MAE/rMAE. Hence, as absolute errors better represent the underlying problem of EPF and the results are similar, the RMSE results are not analyzed here in detail due to space limitations.
- While the LEAR model is indeed estimated using squared errors, this is partly done because the techniques to efficiently estimate the LASSO, e.g. coordinate descent, are based on square errors. This gives the LEAR model a computational advantage over the DNN. An alternative would be to use regularized quantile regression [135] leading, however, to an increased computational burden with little benefits on the accuracy in terms of MAE/rMAE.

6.4.2. Metrics

The obtained results validate the general guideline proposed in Section 5.4 regarding accuracy metrics: research in EPF should avoid MAPE and only use metrics like sMAPE or RMSE in conjunction with any version of rMAE. Particularly, the results validate the following three claims:

- MAE is as reliable as rMAE. However, as the errors are not relative, comparison between datasets is not possible and rMAE is preferred.
- sMAPE is more reliable than MAPE and it mostly agrees with MAE/rMAE. Yet, it disagrees with rMAE and MAE in one of the four datasets and it has the problem of an undefined mean and an infinite variance. Thus, it is less reliable than rMAE.
- MAPE is not a reliable metric as it gives more importance to datapoints close to zero. As such, using MAPE can lead to misleading results and wrong conclusions.
- RMSE is more reliable than MAPE but it does not represent correctly the underlying risks of EPF. Hence, it should not be used alone to evaluate forecasting models.

6.4.3. Performance of open-access models

Based on the extensive comparison of Sections 6.1–6.3, it can be concluded that the models based on DL are more likely to outperform those based on statistical methods. This is especially true in the context of DL ensemble models as the ensemble of DNNs obtains results that are statistically significantly better than any other model.

However, while DNNs generally outperformed the LEAR models, the latter are still the state-of-the-art in terms of low complexity and computational cost. In particular, their performance is very close to that of DNNs, but with the advantage of having computational costs that are up to 100 times lower. As such, they are the best available option when decision making has to be done within seconds.

In short, new models for EPF should either be compared against LEAR models or DNNs depending on the decision time that is available. For a method to be considered more accurate than state-of-the-art methods, it should either be more accurate than the DNN model, or more accurate than LEAR but with similar or lower computational requirements.

7. Checklist to ensure adequate EPF research

As a final contribution, and with the goal of facilitating the work of reviewers of future EPF publications, we provide a short checklist to evaluate whether any new research in EPF satisfies the requirements to be reproducible and to lead to meaningful conclusions:

- The test dataset comprises at least a year of data.
- Any new model is tested against state-of-the-art open-access models, e.g. the ones provided here.
- The computational cost of new methods is evaluated and compared against the computational cost of existing methods.
- The employed datasets are open-access.
- The study is based on multiple markets.
- rMAE is employed as one of the accuracy metrics to evaluate forecasting accuracy.
- Statistical testing is used to assess whether differences in performance are significant.
- Forecasting models are recalibrated on a daily basis and not simply estimated once and evaluated in the full out-of-sample dataset.
- Hyperparameters are estimated using a validation dataset that is different from the test dataset.
- The split and dates of the dataset are explicitly stated.
- All the inputs of the model are explicitly defined.
- The test dataset is selected as the last section of the full dataset and does not contain any overlapping data with the training or validation datasets.
- State-of-the-art and free toolboxes are used for modeling the benchmark models.

While this is just a very short summary of the guidelines described in Section 5, we think it is very useful to have them summarized together for quick evaluations of new research.

8. Conclusion

In this paper, we have derived a set of best practices for performing research in *electricity price forecasting (EPF)*. Particularly, as the field of EPF lacks a rigorous approach to compare and to evaluate new forecasting models, we have analyzed different factors affecting the quality of the research, e.g. dataset size or accuracy metrics, and we have proposed solutions to ensure that new research is adequate, reproducible, and useful.

In addition, as comparisons in EPF are often done using datasets that no other researches has ever tested, we have proposed an extensive open-access benchmark dataset comprising 6 years of recent data in 4 different markets. The aim of the benchmark dataset is to provide a common framework for future research so that new methods can be validated under the same conditions and meaningful comparisons can be obtained. To facilitate future research, we have developed an open-source `python` library named `epftoolbox` [60, 61] that provides easy access to these datasets.

Similarly, as new methods in EPF are often not compared with well-established methods, we have proposed several state-of-the-art open-source models based on statistical methods and deep learning. The methods are tuned automatically and require no expert knowledge in order to be used. These methods are provided as open-source within the proposed `epftoolbox` library [60, 61] so that other researches can employ them as benchmarks in their own studies. Although the proposed methods are currently developed in `python`, we would like to extend the support to other languages; in that spirit, we encourage other researchers to help us do so.

Finally, to have a complete open-access benchmark, we have evaluated the two proposed open-access methods in the open-access dataset and we have provided the results in terms of accuracy metrics and statistical testing. Using these results, we have shown that deep neural networks are more likely to outperform LEAR methods but that the latter are the best model for applications with short decision timeframes. Moreover, we have also shown that ensemble methods often obtain significantly better results than their individual counterparts. Based on the same results, we have also showed the importance of the guidelines as to what constitutes good practices. The most notable guidelines were that MAPE is an unreliable metric that should be avoided, that statistical testing is mandatory to obtain meaningful conclusions, and that the length of the test dataset should be at least one year.

Acknowledgment

This research has received funding from the European Union’s Horizon 2020 research and innovation program under the Marie Skłodowska-Curie grant agreement No. 675318 (INCITE), the Ministry of Science and Higher Education (MNiSW, Poland) through grant No. 0219/DIA/2019/48 and the National Science Center (NCN, Poland) through grant No. 2018/30/A/HS4/00444.

References

- [1] C. Brancucci Martinez-Anido, G. Brinkman, B.-M. Hodge, The impact of wind power on electricity prices, *Renewable Energy* 94 (2016) 474–487. doi:10.1016/j.renene.2016.03.053.
- [2] A. Gianfreda, L. Parisio, M. Pelagatti, The impact of RES in the Italian day-ahead and balancing markets, *Energy Journal* 37 (2016) 161–184. doi:10.5547/01956574.37.si2.agia.
- [3] L. Grossi, F. Nan, Robust forecasting of electricity prices: Simulations, models and the impact of renewable sources, *Technological Forecasting and Social Change* 141 (2019) 305–318. doi:10.1016/j.techfore.2019.01.006.
- [4] K. Maciejowska, Assessing the impact of renewable energy sources on the electricity price level and variability – a quantile regression approach, *Energy Economics* 85 (2020) 104532. doi:10.1016/j.eneco.2019.104532.
- [5] R. Weron, Electricity price forecasting: A review of the state-of-the-art with a look into the future, *International Journal of Forecasting* 30 (4) (2014) 1030–1081. doi:10.1016/j.ijforecast.2014.08.008.
- [6] J. Nowotarski, R. Weron, Recent advances in electricity price forecasting: A review of probabilistic forecasting, *Renewable and Sustainable Energy Reviews* 81 (1) (2018) 1548–1568. doi:10.1016/j.rser.2017.05.234.
- [7] F. Ziel, R. Steinert, Probabilistic mid- and long-term electricity price forecasting, *Renewable and Sustainable Energy Reviews* 94 (2018) 251–266. doi:10.1016/j.rser.2018.05.038.
- [8] T. Hong, P. Pinson, Y. Wang, R. Weron, D. Yang, H. Zareipour, Energy forecasting: A review and outlook, *IEEE Open Access Journal of Power and Energy* (2020) submitted. Working paper version available from RePEc: <https://ideas.repec.org/p/ahh/wpaper/worms2008.html>.
- [9] B. Uniejewski, R. Weron, F. Ziel, Variance stabilizing transformations for electricity spot price forecasting, *IEEE Transactions on Power Systems* 33 (2) (2018) 2219–2229. doi:10.1109/tpwrs.2017.2734563.
- [10] G. Marcjasz, B. Uniejewski, R. Weron, On the importance of the long-term seasonal component in day-ahead electricity price forecasting with NARX neural networks, *International Journal of Forecasting* 35 (4) (2019) 1520–1532. doi:10.1016/j.ijforecast.2017.11.009.
- [11] A. Cruz, A. Muñoz, J. Zamora, R. Espínola, The effect of wind generation and weekday on Spanish electricity spot price forecasting, *Electric Power Systems Research* 81 (10) (2011) 1924–1935. doi:10.1016/j.epsr.2011.06.002.
- [12] L. Wang, Z. Zhang, J. Chen, Short-term electricity price forecasting with stacked denoising autoencoders, *IEEE Transactions on Power Systems* 32 (4) (2016) 2673–2681. doi:10.1109/TPWRS.2016.2628873.
- [13] U. Ugurlu, I. Oksuz, O. Tas, Electricity price forecasting using recurrent neural networks, *Energies* 11 (5) (2018) 1255. doi:10.3390/en11051255.
- [14] W. Zhang, F. Cheema, D. Srinivasan, Forecasting of electricity prices using deep learning networks, in: *Proceedings of the 2018 IEEE PES Asia-Pacific Power and Energy Engineering Conference, 2018*, pp. 451–456. doi:10.1109/APPEEC.2018.8566313.
- [15] S. Luo, Y. Weng, A two-stage supervised learning approach for electricity price forecasting by leveraging different data sources, *Applied Energy* 242 (2019) 1497 – 1512. doi:10.1016/j.apenergy.2019.03.129.
- [16] Y. Chen, Y. Wang, J. Ma, Q. Jin, Brim: An accurate electricity spot price prediction scheme-based bidirectional recurrent neural network and integrated market, *Energies* 12 (12) (2019) 2241. doi:10.3390/en12122241.
- [17] Z. Chang, Y. Zhang, W. Chen, Electricity price prediction based on hybrid model of Adam optimized LSTM neural network and wavelet transform, *Energy* 187 (2019) 115804. doi:10.1016/j.energy.2019.07.134.
- [18] W. Gao, A. Darvishan, M. Toghiani, M. Mohammadi, O. Abedinia, N. Ghadimi, Different states of multi-block based forecast engine for price and load prediction, *International Journal of Electrical Power & Energy Systems* 104 (2019) 423–435. doi:10.1016/j.ijepes.2018.07.014.

- [19] M. S. Nazar, A. E. Fard, A. Heidari, M. Shafie-khah, J. P. Catalão, Hybrid model using three-stage algorithm for simultaneous load and price forecasting, *Electric Power Systems Research* 165 (2018) 214–228. doi:10.1016/j.epsr.2018.09.004.
- [20] L. Zhou, B. Wang, Z. Wang, F. Wang, M. Yang, Seasonal classification and RBF adaptive weight based parallel combined method for day-ahead electricity price forecasting, in: *Proceedings of the 2018 IEEE Power & Energy Society Innovative Smart Grid Technologies Conference*, 2018, pp. 1–5. doi:10.1109/isgt.2018.8403372.
- [21] N. Singh, S. Hussain, S. Tiwari, A PSO-based ANN model for short-term electricity price forecasting, in: *Advances in Intelligent Systems and Computing*, Springer, 2018, pp. 553–563. doi:10.1007/978-981-10-7386-1_47.
- [22] Z. Yang, L. Ce, L. Lian, Electricity price forecasting by a hybrid model, combining wavelet transform, ARMA and kernel-based extreme learning machine methods, *Applied Energy* 190 (2017) 291–305. doi:10.1016/j.apenergy.2016.12.130.
- [23] R. A. Chinnathambi, S. J. Plathottam, T. Hossen, A. S. Nair, P. Ranganathan, Deep neural networks (DNN) for day-ahead electricity price markets, in: *Proceedings of the 2018 IEEE Electrical Power and Energy Conference*, IEEE, 2018, pp. 1–6. doi:10.1109/epec.2018.8598327.
- [24] J. Olamaee, M. Mohammadi, A. Noruzi, S. M. H. Hosseini, Day-ahead price forecasting based on hybrid prediction model, *Complexity* 21 (S2) (2016) 156–164. doi:10.1002/cplx.21792.
- [25] A. Darudi, M. H. Javidi, M. Bashari, Electricity price forecasting using a new data fusion algorithm, *IET Generation, Transmission & Distribution* 9 (12) (2015) 1382–1390. doi:10.1049/iet-gtd.2014.0653.
- [26] M. Ghayekhloo, R. Azimi, M. Ghofrani, M. Menhaj, E. Shekari, A combination approach based on a novel data clustering method and Bayesian recurrent neural network for day-ahead price forecasting of electricity markets, *Electric Power Systems Research* 168 (2019) 184–199. doi:10.1016/j.epsr.2018.11.021.
- [27] A. A. Victoire, B. Gobu, S. Jaikumar, N. Arulmozhi, P. Kanimozhi, A. Victoire, Two-stage machine learning framework for simultaneous forecasting of price-load in the smart grid, in: *Proceedings of the 2018 IEEE International Conference on Machine Learning and Applications*, 2018, pp. 1081–1086. doi:10.1109/icmla.2018.00176.
- [28] M. Zahid, F. Ahmed, N. Javaid, R. Abbasi, H. Zainab Kazmi, A. Javaid, M. Bilal, M. Akbar, M. Ilahi, Electricity price and load forecasting using enhanced convolutional neural network and enhanced support vector regression in smart grids, *Electronics* 8 (2) (2019) 122. doi:10.3390/electronics8020122.
- [29] L. Jiang, G. Hu, Day-ahead price forecasting for electricity market using long-short term memory recurrent neural network, in: *Proceedings of the 2018 International Conference on Control, Automation, Robotics and Vision*, 2018, pp. 949–954. doi:10.1109/icarcv.2018.8581235.
- [30] S. Zhou, L. Zhou, M. Mao, H. Tai, Y. Wan, An optimized heterogeneous structure LSTM network for electricity price forecasting, *IEEE Access* 7 (2019) 108161–108173. doi:10.1109/ACCESS.2019.2932999.
- [31] A. Aggarwal, M. M. Tripathi, A novel hybrid approach using wavelet transform, time series time delay neural network, and error predicting algorithm for day-ahead electricity price forecasting, in: *Proceedings of the 201 International Conference on Computer Applications In Electrical Engineering-Recent Advances*, 2017, pp. 199–204. doi:10.1109/cera.2017.8343326.
- [32] Y.-Y. Hong, C.-Y. Liu, S.-J. Chen, W.-C. Huang, T.-H. Yu, Short-term LMP forecasting using an artificial neural network incorporating empirical mode decomposition, *International Transactions on Electrical Energy Systems* 25 (9) (2014) 1952–1964. doi:10.1002/etep.1949.
- [33] S. Talari, M. Shafie-khah, G. Osório, F. Wang, A. Heidari, J. Catalão, Price forecasting of electricity markets in the presence of a high penetration of wind power generators, *Sustainability* 9 (11) (2017) 2065. doi:10.3390/su9112065.
- [34] N. Singh, S. R. Mohanty, R. D. Shukla, Short term electricity price forecast based on environmentally adapted generalized neuron, *Energy* 125 (2017) 127–139. doi:10.1016/j.energy.2017.02.094.
- [35] G. M. Khan, R. Arshad, N. M. Khan, Efficient prediction of dynamic tariff in smart grid using CGP evolved artificial neural networks, in: *Proceedings of the 2017 IEEE International Conference on Machine Learning and Applications*, 2017, pp. 493–498. doi:10.1109/icmla.2017.0-113.
- [36] M. Afrasiabi, M. Mohammadi, M. Rastegar, A. Kargarian, Multi-agent microgrid energy management based on deep learning forecaster, *Energy* 186 (2019) 115873. doi:10.1016/j.energy.2019.115873.
- [37] Y. Zhu, R. Dai, G. Liu, Z. Wang, S. Lu, Power market price forecasting via deep learning, in: *Proceedings of the 44th Annual Conference of the IEEE Industrial Electronics Society*, 2018, pp. 4935–4939. doi:10.1109/iecon.2018.8591581.
- [38] D. Wang, H. Luo, O. Grunder, Y. Lin, H. Guo, Multi-step ahead electricity price forecasting using a hybrid model based on two-layer decomposition technique and BP neural network optimized by firefly algorithm, *Applied Energy* 190 (2017) 390–407. doi:10.1016/j.apenergy.2016.12.134.
- [39] N. A. Shrivastava, B. K. Panigrahi, M.-H. Lim, Electricity price classification using extreme learning machines, *Neural Computing and Applications* 27 (1) (2014) 9–18. doi:10.1007/s00521-013-1537-1.
- [40] P. Jiang, X. Ma, F. Liu, A new hybrid model based on data preprocessing and an intelligent optimization algorithm for electrical power system forecasting, *Mathematical Problems in Engineering* 2015 (2015) 1–17. doi:10.1155/2015/815253.
- [41] P. Bento, J. Pombo, M. Calado, S. Mariano, A bat optimized neural network and wavelet transform approach for short-term price forecasting, *Applied Energy* 210 (2018) 88–97. doi:10.1016/j.apenergy.2017.10.058.
- [42] M. G. Khajeh, A. Maleki, M. A. Rosen, M. H. Ahmadi, Electricity price forecasting using neural networks with an improved iterative training algorithm, *International Journal of Ambient Energy* 39 (2) (2017) 147–158. doi:10.1080/01430750.2016.1269674.
- [43] J. Lago, F. De Ridder, P. Vrancx, B. De Schutter, Forecasting day-ahead electricity prices in Europe: The importance of considering market integration, *Applied Energy* 211 (2018) 890–903. doi:10.1016/j.apenergy.2017.11.098.
- [44] P.-H. Kuo, C.-J. Huang, An electricity price forecasting model by hybrid structured deep neural networks, *Sustainability* 10 (4) (2018) 1280. doi:10.3390/su10041280.

- [45] S. Mujeeb, N. Javaid, M. Ilahi, Z. Wadud, F. Ishmanov, M. Afzal, Deep long short-term memory: A new price and load forecasting scheme for big data in smart cities, *Sustainability* 11 (4) (2019) 987. doi:10.3390/su11040987.
- [46] S. Atef, A. B. Eltawil, A comparative study using deep learning and support vector regression for electricity price forecasting in smart grids, in: *Proceedings of the 2019 IEEE International Conference on Industrial Engineering and Applications*, 2019, pp. 603–607. doi:10.1109/IEA.2019.8715213.
- [47] S. Mujeeb, N. Javaid, ESAENARX and DE-RELM: Novel schemes for big data predictive analytics of electricity load and price, *Sustainable Cities and Society* 51 (2019) 101642. doi:10.1016/j.scs.2019.101642.
- [48] S. Lahmiri, Comparing variational and empirical mode decomposition in forecasting day-ahead energy prices, *IEEE Systems Journal* 11 (3) (2017) 1907–1910. doi:10.1109/jsyst.2015.2487339.
- [49] S. E. Peter, I. J. Raglend, Sequential wavelet-ANN with embedded ANN-PSO hybrid electricity price forecasting model for Indian energy exchange, *Neural Computing and Applications* 28 (8) (2016) 2277–2292. doi:10.1007/s00521-015-2141-3.
- [50] A. Naz, M. Javed, N. Javaid, T. Saba, M. Alhussain, K. Aurangzeb, Short-term electric load and price forecasting using enhanced extreme learning machine optimization in smart grids, *Energies* 12 (5) (2019) 866. doi:10.3390/en12050866.
- [51] Anamika, R. Peesapati, N. Kumar, Electricity price forecasting and classification through wavelet–dynamic weighted PSO–FFNN approach, *IEEE Systems Journal* 12 (4) (2018) 3075–3084. doi:10.1109/jsyst.2017.2717446.
- [52] W. Gao, V. Sarlak, M. R. Parsaei, M. Ferdosi, Combination of fuzzy based on a meta-heuristic algorithm to predict electricity price in an electricity markets, *Chemical Engineering Research and Design* 131 (2018) 333–345. doi:10.1016/j.cherd.2017.09.021.
- [53] T. Hong, P. Pinson, S. Fan, H. Zareipour, A. Troccoli, R. J. Hyndman, Probabilistic energy forecasting: Global energy forecasting competition 2014 and beyond, *International Journal of Forecasting* 32 (3) (2016) 896–913. doi:10.1016/j.ijforecast.2016.02.001.
- [54] Nord pool website, www.nordpoolspot.com.
- [55] B. Uniejewski, J. Nowotarski, R. Weron, Automated variable selection and shrinkage for day-ahead electricity price forecasting, *Energies* 9 (8) (2016) 621. doi:10.3390/en9080621.
- [56] R. Tibshirani, Regression shrinkage and selection via the lasso, *Journal of the Royal Statistical Society. Series B (Methodological)* (1996) 267–288. doi:10.1111/j.2517-6161.1996.tb02080.x.
- [57] F. Ziel, R. Weron, Day-ahead electricity price forecasting with high-dimensional structures: Univariate vs. multivariate modeling frameworks, *Energy Economics* 70 (2018) 396–420. doi:10.1016/j.eneco.2017.12.016.
- [58] B. Uniejewski, R. Weron, Efficient forecasting of electricity spot prices with expert and lasso models, *Energies* 11 (8) (2018) 2039. doi:10.3390/en11082039.
- [59] J. Lago, F. De Ridder, B. De Schutter, Forecasting spot electricity prices: deep learning approaches and empirical comparison of traditional algorithms, *Applied Energy* 221 (2018) 386–405. doi:10.1016/j.apenergy.2018.02.069.
- [60] Epftoolbox library, <https://github.com/jeslago/epftoolbox>.
- [61] Epftoolbox documentation, <https://epftoolbox.readthedocs.io>.
- [62] K. Mayer, S. Trück, Electricity markets around the world, *Journal of Commodity Markets* 9 (2018) 77–100. doi:10.1016/j.jcomm.2018.02.001.
- [63] R. Aid, *Electricity Derivatives*, Springer, 2015. doi:10.1007/978-3-319-08395-7.
- [64] K. Maciejowska, R. Weron, Electricity price forecasting, in: *Wiley StatsRef: Statistics Reference Online*, Wiley, 2019, pp. 1–9. doi:10.1002/9781118445112.stat08215.
- [65] R. Weron, F. Ziel, Electricity price forecasting, in: U. Soytaş, R. Sari (Eds.), *Routledge Handbook of Energy Economics*, Routledge, 2018, pp. 506–521. doi:10.4324/9781315459653-36.
- [66] A. Gianfreda, F. Ravazzolo, L. Rossini, Comparing the forecasting performances of linear models for electricity prices with high RES penetration, *International Journal of Forecasting* 36 (2020) 974–986. doi:10.1016/j.ijforecast.2019.11.002.
- [67] H. Zou, T. Hastie, Regularization and variable selection via the elastic net, *Journal of the Royal Statistical Society: Series B (Statistical Methodology)* 67 (2) (2005) 301–320. doi:10.1111/j.1467-9868.2005.00503.x.
- [68] F. Ziel, R. Steinert, S. Husmann, Forecasting day ahead electricity spot prices: The impact of the EXAA to other European electricity markets, *Energy Economics* 51 (2015) 430–444. doi:10.1016/j.eneco.2015.08.005.
- [69] F. Ziel, Forecasting electricity spot prices using lasso: On capturing the autoregressive intraday structure, *IEEE Transactions on Power Systems* 31 (6) (2016) 4977–4987. doi:10.1109/tpwrs.2016.2521545.
- [70] S. Schneider, Power spot price models with negative prices, *Journal of Energy Markets* 4 (4) (2011) 77–102. doi:10.21314/jem.2011.079.
- [71] G. Diaz, E. Planas, A note on the normalization of Spanish electricity spot prices, *IEEE Transactions on Power Systems* 31 (3) (2016) 2499–2500. doi:10.1109/tpwrs.2015.2449757.
- [72] J. Nowotarski, J. Tomczyk, R. Weron, Robust estimation and forecasting of the long-term seasonal component of electricity spot prices, *Energy Economics* 39 (2013) 13–27. doi:10.1016/j.eneco.2013.04.004.
- [73] J. Nowotarski, R. Weron, On the importance of the long-term seasonal component in day-ahead electricity price forecasting, *Energy Economics* 57 (2016) 228–235. doi:10.1016/j.eneco.2016.05.009.
- [74] F. Lisi, M. Pelagatti, Component estimation for electricity market data: Deterministic or stochastic?, *Energy Economics* 74 (2018) 13–37. doi:10.1016/j.eneco.2018.05.027.
- [75] G. Marcjasz, B. Uniejewski, R. Weron, On the importance of the long-term seasonal component in day-ahead electricity price forecasting with NARX neural networks, *International Journal of Forecasting* 35 (4) (2019) 1520–1532. doi:10.1016/j.ijforecast.2017.11.009.
- [76] K. Hubicka, G. Marcjasz, R. Weron, A note on averaging day-ahead electricity price forecasts across calibration windows, *IEEE Transactions on Sustainable Energy* 10 (1) (2019) 321–323. doi:10.1109/tste.2018.2869557.
- [77] G. Marcjasz, T. Serafin, R. Weron, Selection of calibration windows for day-ahead electricity price forecasting, *Energies*

- 11 (9) (2018) 2364. doi:10.3390/en11092364.
- [78] S. Mujeeb, N. Javaid, M. Akbar, R. Khalid, O. Nazeer, M. Khan, Big data analytics for price and load forecasting in smart grids, in: *Lecture Notes on Data Engineering and Communications Technologies*, Springer, 2018, pp. 77–87. doi:10.1007/978-3-030-02613-4_7.
- [79] X. Xie, W. Xu, H. Tan, The day-ahead electricity price forecasting based on stacked CNN and LSTM, in: *Lecture Notes in Computer Science*, Springer International Publishing, 2018, pp. 216–230. doi:10.1007/978-3-030-02698-1_19.
- [80] U. Ugurlu, O. Tas, A. Kaya, I. Oksuz, The financial effect of the electricity price forecasts’ inaccuracy on a hydro-based generation company, *Energies* 11 (8) (2018) 2093. doi:10.3390/en11082093.
- [81] J. K. Kolberg, K. Waage, Artificial intelligence and Nord Pool’s intraday electricity market Elbas: a demonstration and pragmatic evaluation of employing deep learning for price prediction: using extensive market data and spatio-temporal weather forecasts, Master’s thesis, Norwegian School of Economics (2018).
- [82] J. Xu, R. Baldick, Day-ahead price forecasting in ERCOT market using neural network approaches, in: *Proceedings of the Tenth ACM International Conference on Future Energy Systems*, 2019, pp. 486–491. doi:10.1145/3307772.3331024.
- [83] J.-H. Meier, S. Schneider, I. Schmidt, P. Schüller, T. Schönfeldt, B. Wanke, ANN-based electricity price forecasting under special consideration of time series properties, in: *Information and Communication Technologies in Education, Research, and Industrial Applications*, Springer International Publishing, 2019, pp. 262–275. doi:10.1007/978-3-030-13929-2_13.
- [84] Z. Chang, Y. Zhang, W. Chen, Effective adam-optimized LSTM neural network for electricity price forecasting, in: *Proceedings of the 2018 IEEE International Conference on Software Engineering and Service Science*, 2018, pp. 245–248. doi:10.1109/icsess.2018.8663710.
- [85] H. Jahangir, H. Tayarani, S. Baghali, A. Ahmadian, A. Elkamel, M. Aliakbar Golkar, M. Castilla, A novel electricity price forecasting approach based on dimension reduction strategy and rough artificial neural networks, *IEEE Transactions on Industrial Informatics* (2019) 1–1doi:10.1109/TII.2019.2933009.
- [86] W. Ahmad, N. Javaid, A. Chand, S. Y. R. Shah, U. Yasin, M. Khan, A. Syeda, Electricity price forecasting in smart grid: A novel E-CNN model, in: *Web, Artificial Intelligence and Network Applications*, Springer International Publishing, 2019, pp. 1132–1144. doi:10.1007/978-3-030-15035-8_109.
- [87] D. Aineto, J. Iranzo-Sánchez, L. G. Lemus-Zúñiga, E. Onaindia, J. F. Urchueguía, On the influence of renewable energy sources in electricity price forecasting in the Iberian market, *Energies* 12 (11) (2019) 2082. doi:10.3390/en12112082.
- [88] S. Schnürch, A. Wagner, Machine learning on EPEX order books: Insights and forecasts, arXiv preprint (2019). arXiv:1906.06248.
- [89] J. Zhang, Z. Tan, C. Li, A novel hybrid forecasting method using GRNN combined with wavelet transform and a GARCH model, *Energy Sources, Part B: Economics, Planning, and Policy* 10 (4) (2015) 418–426. doi:10.1080/15567249.2011.557685.
- [90] J.-L. Zhang, Y.-J. Zhang, D.-Z. Li, Z.-F. Tan, J.-F. Ji, Forecasting day-ahead electricity prices using a new integrated model, *International Journal of Electrical Power & Energy Systems* 105 (2019) 541–548. doi:10.1016/j.ijepes.2018.08.025.
- [91] H. Varshney, A. Sharma, R. Kumar, A hybrid approach to price forecasting incorporating exogenous variables for a day ahead electricity market, in: *Proceedings of the 2016 IEEE International Conference on Power Electronics, Intelligent Control and Energy Systems*, 2016, pp. 1–6. doi:10.1109/icpeices.2016.7853355.
- [92] L. Xiao, W. Shao, M. Yu, J. Ma, C. Jin, Research and application of a hybrid wavelet neural network model with the improved cuckoo search algorithm for electrical power system forecasting, *Applied Energy* 198 (2017) 203–222. doi:10.1016/j.apenergy.2017.04.039.
- [93] R. Bisoi, P. K. Dash, P. P. Das, Short-term electricity price forecasting and classification in smart grids using optimized multikernel extreme learning machine, *Neural Computing and Applications*doi:10.1007/s00521-018-3652-5.
- [94] M. K. Kim, Short-term price forecasting of nordic power market by combination Levenberg–Marquardt and cuckoo search algorithms, *IET Generation, Transmission & Distribution* 9 (13) (2015) 1553–1563. doi:10.1049/iet-gtd.2014.0957.
- [95] A. Pourdaryaei, H. Mokhlis, H. A. Illias, S. H. A. Kaboli, S. Ahmad, Short-term electricity price forecasting via hybrid backtracking search algorithm and ANFIS approach, *IEEE Access* 7 (2019) 77674–77691. doi:10.1109/access.2019.2922420.
- [96] H. Ebrahimian, S. Barmayoon, M. Mohammadi, N. Ghadimi, The price prediction for the energy market based on a new method, *Economic Research-Ekonomska Istraživanja* 31 (1) (2018) 313–337. doi:10.1080/1331677x.2018.1429291.
- [97] O. Abedinia, N. Amjady, M. Shafie-khah, J. Catalão, Electricity price forecast using combinatorial neural network trained by a new stochastic search method, *Energy Conversion and Management* 105 (2015) 642–654. doi:10.1016/j.enconman.2015.08.025.
- [98] S. Itaba, H. Mori, An electricity price forecasting model with fuzzy clustering preconditioned ANN, *Electrical Engineering in Japan* 204 (3) (2018) 10–20. doi:10.1002/eej.23094.
- [99] M. Ghofrani, R. Azimi, F. M. Najafabadi, N. Myers, A new day-ahead hourly electricity price forecasting framework, in: *Proceedings of the 2017 North American Power Symposium*, 2017, pp. 1–6. doi:10.1109/naps.2017.8107269.
- [100] S. Itaba, H. Mori, A fuzzy-preconditioned GRBFN model for electricity price forecasting, *Procedia Computer Science* 114 (2017) 441–448. doi:10.1016/j.procs.2017.09.010.
- [101] ENTSO-E transparency platform, <https://transparency.entsoe.eu/>, last accessed: 2019-09-15.
- [102] PJM website, www.pjm.com.
- [103] Elia, Grid data, <http://www.elia.be/en/grid-data/dashboard>, last accessed: 2019-09-19.
- [104] RTE, Grid data, <https://data.rte-france.com/>, last accessed: 2017-05-15.
- [105] Amprion website, <https://www.netztransparenz.de/>.
- [106] Information platform of the German transmission system operators, <https://www.amprion.net/>.

- [107] J. Bergstra, R. Bardenet, Y. Bengio, B. Kégl, Algorithms for hyper-parameter optimization, in: *Advances in Neural Information Processing Systems*, 2011, pp. 2546–2554.
- [108] T. Hastie, R. Tibshirani, J. Friedman, *The Elements of Statistical Learning*, Springer Series in Statistics, Springer New York Inc., New York, NY, USA, 2001. doi:10.1007/978-0-387-21606-5.
- [109] F. Ziel, R. Steinert, S. Husmann, Efficient modeling and forecasting of electricity spot prices, *Energy Economics* 47 (2015) 98–111. doi:10.1016/j.eneco.2014.10.012.
- [110] B. Efron, T. Hastie, I. Johnstone, R. Tibshirani, Least angle regression, *Annals of Statistics* 32 (2) (2004) 407–499. doi:10.1214/009053604000000067.
- [111] Y. Yao, L. Rosasco, A. Caponnetto, On early stopping in gradient descent learning, *Constructive Approximation* 26 (2) (2007) 289–315. doi:10.1007/s00365-006-0663-2.
- [112] F. Pedregosa, G. Varoquaux, A. Gramfort, V. Michel, B. Thirion, O. Grisel, M. Blondel, P. Prettenhofer, R. Weiss, V. Dubourg, J. Vanderplas, A. Passos, D. Cournapeau, M. Brucher, M. Perrot, E. Duchesnay, Scikit-learn: Machine learning in Python, *Journal of Machine Learning Research* 12 (2011) 2825–2830.
- [113] F. Chollet, Keras (2015).
URL <https://github.com/fchollet/keras>
- [114] R. J. Hyndman, Errors on percentage errors (2014).
URL <https://robjhyndman.com/hyndsight/smape/>
- [115] R. J. Hyndman, A. B. Koehler, Another look at measures of forecast accuracy, *International Journal of Forecasting* 22 (4) (2006) 679–688. doi:10.1016/j.ijforecast.2006.03.001.
- [116] R. J. Hyndman, G. Athanasopoulos, *Forecasting: principles and practice*, OTexts, 2018.
- [117] F. X. Diebold, R. S. Mariano, Comparing predictive accuracy, *Journal of Business & Economic Statistics* 13 (3) (1995) 253–263. doi:10.1080/07350015.1995.10524599.
- [118] F. X. Diebold, Comparing predictive accuracy, twenty years later: A personal perspective on the use and abuse of Diebold-Mariano tests, *Journal of Business and Economic Statistics* 33 (1) (2015) 1–9. doi:10.1080/07350015.2014.983236.
- [119] S. Bordignon, D. W. Bunn, F. Lisi, F. Nan, Combining day-ahead forecasts for British electricity prices, *Energy Economics* 35 (2013) 88–103. doi:10.1016/j.eneco.2011.12.001.
- [120] J. Nowotarski, E. Raviv, S. Trück, R. Weron, An empirical comparison of alternative schemes for combining electricity spot price forecasts, *Energy Economics* 46 (2014) 395–412. doi:10.1016/j.eneco.2014.07.014.
- [121] T. Serafin, B. Uniejewski, R. Weron, Averaging predictive distributions across calibration windows for day-ahead electricity price forecasting, *Energies* 12 (13) (2019) 2561. doi:10.3390/en12132561.
- [122] G. Marcjasz, J. Lago, R. Weron, B. D. Schutter, Neural networks in day-ahead electricity price forecasting: single vs. multiple outputs, *Energy Conversion and Management* (Submitted).
- [123] R. Giacomini, H. White, Tests of conditional predictive ability, *Econometrica* 74 (6) (2006) 1545–1578. doi:10.1111/j.1468-0262.2006.00718.x.
- [124] R. Giacomini, B. Rossi, Forecasting in macroeconomics, in: *Handbook of Research Methods and Applications in Empirical Macroeconomics*, Edward Elgar Publishing, 2013, pp. 381–408. doi:10.4337/9780857931023.00024.
- [125] N. N. A. N. Ibrahim, I. A. W. A. Razak, S. S. M. Sidin, Z. H. Bohari, Electricity price forecasting using neural network with parameter selection, in: *Intelligent and Interactive Computing*, Springer Singapore, 2019, pp. 141–148. doi:10.1007/978-981-13-6031-2_33.
- [126] I. P. Panapakidis, A. S. Dagoumas, Day-ahead electricity price forecasting via the application of artificial neural network based models, *Applied Energy* 172 (2016) 132–151. doi:10.1016/j.apenergy.2016.03.089.
- [127] N. K. Singh, A. K. Singh, M. Tripathy, Short-term load/price forecasting in deregulated electric environment using ELMAN neural network, in: *Proceedings of the 2015 International Conference on Energy Economics and Environment*, 2015, pp. 1–6. doi:10.1109/energyeconomics.2015.7235086.
- [128] S. S. Reddy, C.-M. Jung, K. J. Seog, Day-ahead electricity price forecasting using back propagation neural networks and weighted least square technique, *Frontiers in Energy* 10 (1) (2016) 105–113. doi:10.1007/s11708-016-0393-y.
- [129] J. Nascimento, T. Pinto, Z. Vale, Day-ahead electricity market price forecasting using artificial neural network with spearman data correlation, in: *Proceedings of the 2019 IEEE PowerTech Conference*, 2019, pp. 1–6. doi:10.1109/ptc.2019.8810618.
- [130] D. Kotur, M. Zarkovic, Neural network models for electricity prices and loads short and long-term prediction, in: *Proceedings of the 2016 International Symposium on Environmental Friendly Energies and Applications*, 2016, pp. 1–5. doi:10.1109/efea.2016.7748787.
- [131] C. Monteiro, I. Ramirez-Rosado, L. Fernandez-Jimenez, P. Conde, Short-term price forecasting models based on artificial neural networks for intraday sessions in the Iberian electricity market, *Energies* 9 (9) (2016) 721. doi:10.3390/en9090721.
- [132] C. Monteiro, L. Fernandez-Jimenez, I. Ramirez-Rosado, Explanatory information analysis for day-ahead price forecasting in the Iberian electricity market, *Energies* 8 (9) (2015) 10464–10486. doi:10.3390/en80910464.
- [133] Anamika, N. Kumar, Market-clearing price forecasting for Indian electricity markets, in: *Proceeding of International Conference on Intelligent Communication, Control and Devices*, Springer Singapore, 2016, pp. 633–642. doi:10.1007/978-981-10-1708-7_72.
- [134] A. F. Atiya, Why does forecast combination work so well?, *International Journal of Forecasting* 36 (1) (2020) 197–200. doi:10.1016/j.ijforecast.2019.03.010.
- [135] Y. Li, J. Zhu, L1-norm quantile regression, *Journal of Computational and Graphical Statistics* 17 (2008) 163–185. doi:10.1198/106186008x289155.