

# Short-Term Forecasting of CO<sub>2</sub> Emission Intensity in Power Grids by Machine Learning

Kenneth Leerbeck<sup>a,\*</sup>, Peder Bacher<sup>a</sup>, Rune Grønberg Junker<sup>a</sup>, Goran Goranović<sup>a</sup>, Olivier Corradi<sup>b</sup>, Razgar Ebrahimi<sup>a</sup>, Anna Tveit<sup>a</sup>, Henrik Madsen<sup>a,c</sup>

<sup>a</sup>Technical University of Denmark

<sup>b</sup>Tmrow IVS

<sup>c</sup>Norwegian University of Science and Technology

---

## Abstract

A machine learning algorithm is developed to forecast the CO<sub>2</sub> emission intensities in electrical power grids in the Danish bidding zone DK2, distinguishing between average and marginal emissions. The analysis was done on data set comprised of a large number (473) of explanatory variables such as power production, demand, import, weather conditions etc. collected from selected neighboring zones. The number was reduced to less than 30 using both LASSO (a penalized linear regression analysis) and a forward feature selection algorithm. Three linear regression models that capture different aspects of the data (non-linearities and coupling of variables etc.) were created and combined into a final model using Softmax weighted average. Cross-validation is performed for debiasing and autoregressive moving average model (ARIMA) implemented to correct the residuals, making the final model the variant with exogenous inputs (ARIMAX). The forecasts with the corresponding uncertainties are given for two time horizons, below and above six hours. Marginal emissions came up highly independent of any conditions in the DK2 zone, suggesting that the marginal generators are located in the neighbouring zones.

The developed methodology can be applied to any bidding zone in the European electricity network without requiring detailed knowledge about the zone.

**Keywords:** CO<sub>2</sub> emission forecasting; electrical power grids; machine learning; feature selection; LASSO; ARIMA.

---

## 1. Introduction

Consumption of electricity contributes heavily to the CO<sub>2</sub> emissions, [1]. The intensity with which it does so depends on the proportion of renewable sources (eg. solar, wind) vs. nonrenewable sources (eg. coal, gas, nuclear). The proportion, and hence the CO<sub>2</sub> emission, fluctuates with time based on power market mechanisms and weather conditions. Ideally, in the future, electricity users (the demand) would respond to the renewable power generation in attempt to lower the emissions. Proposed solutions include scheduling of storage (e.g. batteries, fuel cells, hydro reservoirs, thermal) and flexible demand (e.g. heat pumps, electric cars), [2]. This paper presents additional methodology aimed to fulfill the goal of lowering the emissions.

Forecasting of the power grid is essential for these solutions and exist in various forms already. For renewables, new forecasting methods are developed on a regular basis to increase revenue. Most recently in [3], a multi-step ahead deterministic forecasting on wind power is done (many prior models only did one-step ahead), using complex multi-stage machine learning (kernel-based) algorithms for error correction. The deterministic model, however, cannot reliably account for the volatile wind speeds. Hence a probabilistic model (distribution forecast) of wind speeds based on robust machine learning algorithms was developed in [4]. The probabilistic models quantify uncertainty of a forecast, crucial for risk management. On the other hand, multi-step forecasts are important in any scheduling application, such as market bids and flexibility. Examples of the combination of the two types of models are the multi-step ahead probabilistic solar power forecasts developed in [5] and [6] with horizons of 36 and 72 hours ahead respectively.

---

\*Corresponding Author

Email address: kenle@dtu.dk

Postal Address: Anker Engelunds vej 1, Building 101A, 2800 Kongens Lyngby, Denmark

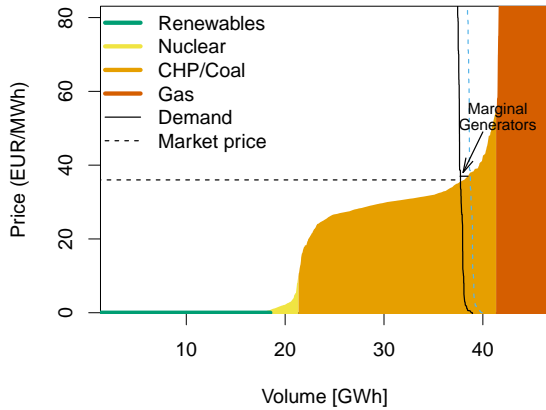


Figure 1: The merit order illustrated with a supply and demand curve example. The x-axis is the accumulated generators in the power system and the y-axis is their corresponding costs. The highest generator in the merit order is the one crossing with the demand curve - a coal generator in this example. The average emissions are a weighted average from all activated generators. The **marginal generator** is the generator that will be activated by moving the demand line slightly to the right (dashed blue line). Data source: Nord Pool AS.

The generator responding to small changes in demand is called the marginal generator and is not weather dependent (must generate on-demand) - hence, does not include wind turbines and PV-panels<sup>1</sup>. A good estimate of the marginal generator is achieved by using price signals, Figure 1, where the marginal generator is, in this case, a coal fired generator.

Because price-based control can be both economically and environmentally beneficial, forecasting of the day-ahead spot prices has been proposed in recent papers [7, 8], using different methodologies (e.g. ARIMA, Neural Networks, machine learning). There is a problem though with the spot prices, known as the merit order emission dilemma, illustrated for the German-Austrian power market [9]: The price for coal is low but the emissions are high. A price based-control, therefore, only leads to low emissions if there is surplus of renewable energy (more renewable energy than needed) - otherwise coal is favored.

Precise estimates and forecasts of the marginal CO<sub>2</sub> emissions are needed to correctly implement storage

<sup>1</sup>In modern electricity networks with high proportions of renewable sources, a weather-dependent generator *can* become marginal. For example, during overproduction from wind turbines, demand cannot keep up and pushes down the electricity prices; the wind turbines can thus be shut-off to down-regulate the production. According to data supplied by Energinet Denmark, these situations occurred for 12.5% of all the hours in 2017 and 2018 in Western Denmark (bidding zone DK1), where the marginal generator could be argued to be a wind turbine. This calls for a (future) refinement of the marginal CO<sub>2</sub> emission estimation.

and flexible demand into the power grid. Discussing CO<sub>2</sub> emissions, two distinct concepts are used; average and marginal emission intensities ( $\frac{\text{kgCO}_2\text{-eq}}{\text{MWh}}$ ). Average emissions correspond to the overall, e.g. region-wide, electricity production including net imports. The marginal reflects the emissions of the marginal generator. The concepts are compared in [10] and the importance of distinguishing between the two is highlighted due to their very opposing patterns.

Both concepts have been estimated in prior studies [10–19]. Early studies of the marginal emissions, [11–14], all estimate the highest generator in the merit order of the power generation system. However, this is rarely the only generator responding to a change in demand, as addressed in [15], where these early approaches are discussed and a new empirical approach is presented; By estimating the contribution of all power generators to a specific change in system demand, using linear regression on historical data including outputs from major power producers in the UK to estimate the average response from each generation technology class to changes in demand. The power plants were disaggregated to investigate the impact of plant turnovers (switching old power plants for new ones). Traditionally treated as dummy variables, the imports have been treated explicitly in a recent study that focused on the average emissions in the Nordic European countries, [18]. The study showed the interplay between the imports and the average emissions and how both vary from one bidding zone to another. Incorporating the imports in the marginal emissions, the company Tmrow IVS has developed a new empirical approach using machine learning on historical data that follows the chain of imports (the so-called flow tracing, originally introduced in [20, 21]) to assess the impact of a specific generator or load on the power system [19]. This is a large scale solution using data from the majority of bidding zones around the world.

The just mentioned studies provide methodologies for marginal CO<sub>2</sub> emission *estimates*. Also, the *long-term* (e.g. annual) forecasting of CO<sub>2</sub> emissions are widely conducted for promoting green energy, e.g. [22]. However, the more accurate short-term emission forecasting methodologies, currently unavailable in the literature, are needed to implement the flexible demand. In this study, short-term (24h ahead) forecasts with uncertainty margins (95% prediction intervals) of both the average and marginal CO<sub>2</sub> emission intensities from the power generation in bidding zone DK2 (Sealand region, Denmark) are developed. These enable flexible consumers (electric cars, heat pumps, etc.) to sched-

ule for optimal electricity usage i.e minimal CO<sub>2</sub> emissions, be it for regulatory or branding purposes. The methodology can be applied to other bidding zones in the European electricity network without requiring detailed knowledge of response and explanatory variables.

The forecasting in this study models the given response variables (average and marginal CO<sub>2</sub> emission intensities) in terms of so-called explanatory variables - e.g. power production, demand, import, weather conditions, etc. - that are represented by generalized basis functions (splines). Examples of the explanatory variables with respect to the response variables are shown in Sec. 2. The data is divided into two sets: one set available for  $\leq 6$  hours and another available for horizons  $> 6$  hours. The machine learning techniques that include trend extraction (seasonality, nonlinearity, interaction terms), feature selection (LASSO, forward feature selection), residual correction (ARIMA) and cross-validation are elaborated in Sec. 3 and 4, where also three different models are built and combined with a Softmax weighted average into the overall forecasting model. The most significant variables and forecasting results are highlighted in Sec. 5, and concluding remarks are found in Sec. 6. A list of all the variables can be reviewed in Appendix A where it is indicated which variables are being used for which set of data.

## 2. Data analysis: examples

The CO<sub>2</sub> emissions in Denmark are interesting when scheduling flexible consumers due to large amounts of wind power production - 48% of the total electricity production in 2017. The country also has a variety of good trading options with its neighboring countries - e.g. Germany (DE), Sweden (SW) and Norway (NO). Indirectly, influence can also come from countries further away, but the scope of this study is limited to selected bidding zones DE, DK1, DK2, NO2, SW3 and SW4. The focus area, DK2, has direct transmission cables to DK1, SW4 and DE.

Linear relationships between the response and explanatory variables are detected first. The power production in DK2 and net import in SW4 from SW3 show the highest correlations to the response variables as illustrated in Figure 2. The average emissions are highly correlated to the power production in DK2 because it is a response to the demand and this will mainly activate coal and gas-fueled generators. Net import in SW4 from SW3 show the highest correlation to the marginal emissions and is an indirect influence - this shows that the marginal generator is often located in SW3 where all Sweden's nuclear power is produced [23]. This is

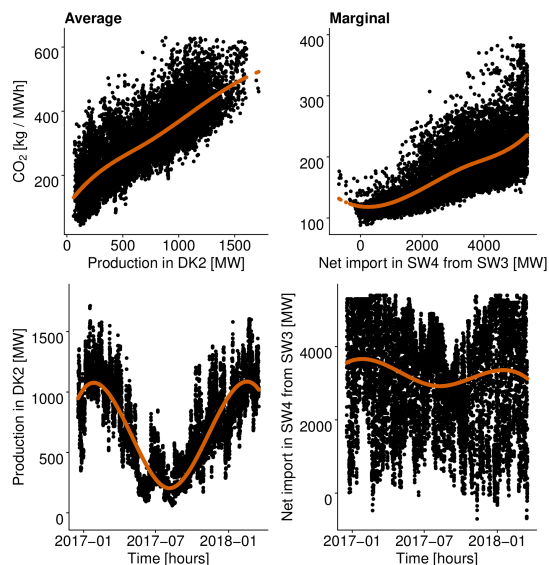


Figure 2: Top left: The average CO<sub>2</sub> emission intensity vs the power production in DK2. Top right: The marginal CO<sub>2</sub> emission intensity vs the net import in SW3 from SW4. Bottom: The power production in DK2 and the net import in SW3 from SW4 vs time (hourly resolution).

because nuclear power is cheaper than the local options in DK2 e.g. coal and gas. Figure 2 also shows a yearly seasonality in the local power production varying about 1,000 MW - lowest in the summer. The import show the same yearly seasonality but less significant because nuclear power serves the baseload.

Next, a linear regression model can be fitted onto the average and marginal emissions using the discussed production and import respectively to reveal other important variables. The residuals from these models will represent the average and marginal emissions that are independent of these variables. This reveals the non-linear relationships shown in Figure 3.

Interestingly, the average emissions are the highest when the net import from SW4 is zero, indicating trades usually happen when cheap non-polluting electricity is generated - e.g. nuclear and renewables - and the corresponding CO<sub>2</sub> emissions are proportionally lowered. This extra power, produced e.g. in Sweden, can then be imported for use and serves as an indicator of the lower CO<sub>2</sub> emissions. The marginal emissions are highly depending on the net import in DK1 from DE but only when DK1 exports to DE. The higher the export the more domestic generators must serve as marginals. The net import in DK2 from SW4 shows a yearly seasonality being the highest in the summer - recall the domestic production being the lowest here - since demand is gen-

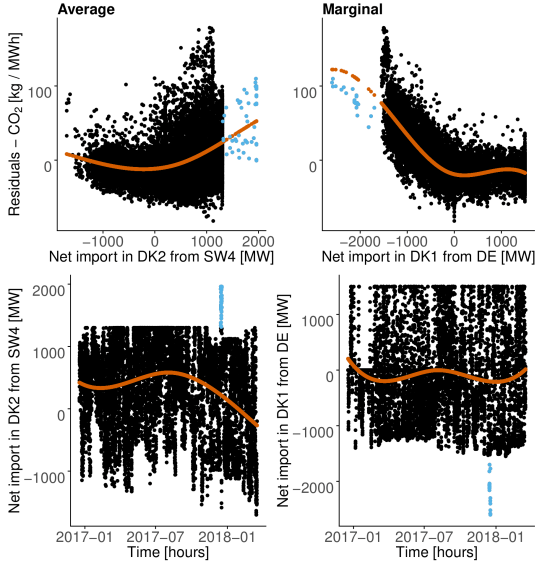


Figure 3: Top left: The residuals from the linear fit of the domestic production (DK2) onto the average CO<sub>2</sub> emission intensity vs the net import in DK2 from SW4. Top right: The residuals from the linear fit of the import in SW4 from SW3 onto the marginal CO<sub>2</sub> emission intensity vs the net import in DK1 from DE. Bottom: The net import in DK2 from SW4 and the net import in DK1 from DE vs time (hourly resolution). Note the outliers (light blue); over three days, the import seemingly overloaded both the transmission cables to Germany and Sweden. For this study, these data points are modified to the maximum capacity (1,300 and 1,550 MW in these cases)

erally low the proportion of nuclear power is large and Denmark can, therefore, import it rather than produce electricity locally from gas or coal.

### 3. Regression models and basis functions

#### 3.1. Linear Regression Models (*lm*)

The CO<sub>2</sub> emissions are modeled using a multivariate linear regression model

$$\mathbf{Y} = \mathbf{X}\beta + \epsilon, \quad (1)$$

for  $\epsilon \sim N(0, \sigma^2 I)$ ,

where  $\mathbf{X}$  is the input matrix (explanatory variables),  $\mathbf{Y}$  is the output vector (response variables: CO<sub>2</sub> emissions) and  $\beta$  is a vector of regression coefficients to be found.  $\epsilon$  represents the normally distributed errors in the model.

The least square regression is performed to minimize

$$S(\beta) = \|\mathbf{Y} - \mathbf{X}\beta\|^2$$

and obtain the ordinary least-squares solution

$$\beta = (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \mathbf{Y}. \quad (2)$$

Once  $\beta$  is obtained on training data, the response variables  $\hat{y}$  can be forecasted on new input data (the test data)

$$\hat{\mathbf{y}}_{t+h} = \left[ \mathbf{x}_{t+h}^F \quad MA(\mathbf{x}_{t+h}^{RT})_{24} \quad MA(\mathbf{x}_{t+h}^{RT})_{48} \quad lag(\mathbf{y})_{t+h} \right] \beta \equiv \mathbf{z}_{t+h}^* \beta, \quad (3)$$

where F and RT refer to Forecast and Real-Time, and MA to Moving Average. The terms in the brackets (the matrix elements) are provided by Tmrow:  $\mathbf{x}_{t+h}^F$  are the input variables forecasted  $h$ -hours ahead ( $h \leq 6$  hours for all explanatory variables except weather data which is available for  $h \geq 6$  hours). The moving average is constructed to translate the real-time input variables into the forecasting format

$$MA(\mathbf{x}_{t+h}^{RT})_n = \frac{1}{n} \sum_{i=0}^{n-1} \mathbf{x}_{t+h-i}^{RT}, \quad (4)$$

where  $n$  is the length of the averaging period - either 24 or 48 hours. Finally,  $lag(\mathbf{y})_{t+h}$  in Eq. 3 is just  $\mathbf{y}_t$ , the last available observation at time  $t$ . The idea is to include all available information to obtain the forecast. This is an Auto-Regression (regression on lagged values of  $\mathbf{y}$ ) model with eXogeneous inputs (regression on functions of  $\mathbf{x}$ ) or ARX for short.

Both  $\mathbf{X}$  and  $\mathbf{Y}$  are given as time series.  $\mathbf{X}$  consists of 418 variables - for  $h \leq 6$  hours - (related both to weather and power system from six bidding zones; 66 listed in Appendix A) and  $\mathbf{Y}$  of 10,897 observations. The forecasted input variables at  $t+h$  constitute a new input matrix  $\mathbf{Z}$  that will be expanded with more columns (next three subsections).

#### 3.2. Periodic variations: Fourier Series

Periodic variations (seasonality) of the average and marginal CO<sub>2</sub> emission intensities are investigated using Fourier Series defined as

$$FS(n, period)_t = \hat{\mathbf{y}}_{FS(n, period), t} = \quad (5)$$

$$A_0 + \sum_{i=1}^{n_{series}} A_i \cdot \sin\left(i \frac{2\pi t}{period}\right) + B_i \cdot \cos\left(i \frac{2\pi t}{period}\right),$$

where  $\mathbf{y}$  is the response variable,  $t$  is the time,  $A$  and  $B$  are linear regression coefficient matrices,  $n$  is the order of the Fourier Series and 'period' is the length of the seasonality period.  $n$  is adjusted to find the best fit. In Figure 4 the daily, weekly, and yearly patterns are estimated using data ranging from January 2015 to January 2018 with  $n = [2, 1, 2]$  respectively.

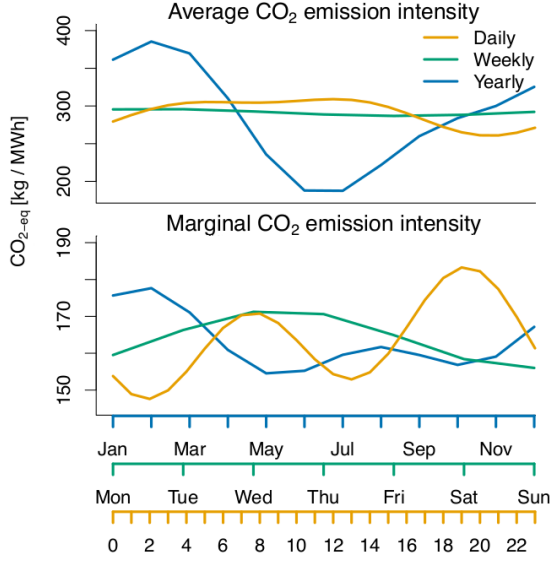


Figure 4: Fourier Series showing the daily, weekly and yearly patterns for the average and marginal CO<sub>2</sub> emission intensities.

It is observed that average and marginal emissions follow completely different patterns. The yearly pattern shows the largest variations for the average emissions, as already discussed in Section 2. The daily average emission varies very little in comparison to the yearly emissions. On the other hand, the daily marginal emissions are the highest when the average emissions are the lowest, illustrating the importance of using the correct emission measure. The weekly pattern is only of importance to the marginal emissions and is lowest on the weekends. The seasonality components are added to the input matrix  $\mathbf{Z}$  of Eq. 3 such that

$$\mathbf{z}_t = \left[ \mathbf{z}_t^* \quad \exp^{\mathbf{z}_t^*} \quad FS(2, 24)_t \quad FS(1, 7)_t \quad FS(2, 12)_t \right], \quad (6)$$

Note that an exponential term of  $\mathbf{Z}$  is added, too, which makes  $Z$  a (10,897x951) matrix (both  $\mathbf{Z}^*$  and  $\exp(\mathbf{Z}^*)$  have 474 columns after the clean-up of non-available points and constant variables).

### 3.3. Non-linearities: Splines

Non-linearities are captured by using splines, the local polynomials between specified points called knots [24]. Splines are implemented in R with the built-in functions  $bs()$  (base splines) and  $ns()$  (natural splines). The former are basis functions and increasing their number in the expansion of a function improves the fitting procedure at the risk of overfitting. In this

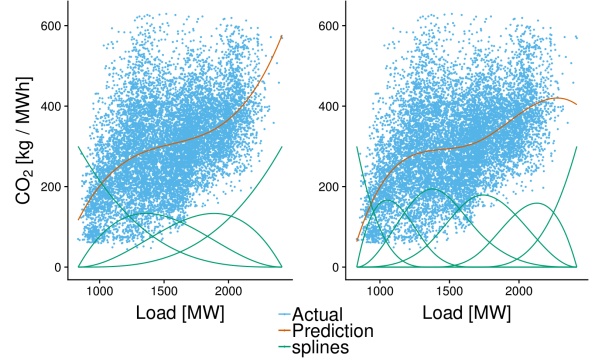


Figure 5: Estimation of the average CO<sub>2</sub> emission intensity (orange) based on four and six splines, respectively. The splines are scaled for illustration.

study, four base splines with knots located at the quantiles (default settings in R), are used to avoid overfitting

$$bs(\mathbf{z}_t) = \left[ bs_0(\mathbf{z}_t) \quad bs_1(\mathbf{z}_t) \quad bs_2(\mathbf{z}_t) \quad bs_3(\mathbf{z}_t) \right]^T. \quad (7)$$

The number is justified by the assumption that the relationships between the explanatory and the response variables are stationary. The least-square coefficients associated with the base splines are labeled by the vector  $\beta^{bs}$ . To refine the fitting, natural splines  $ns(\mathbf{z}_t)$  are also used and represented analogously to Eq. 7.

Figure 5 shows the estimated mean (orange) of the average CO<sub>2</sub> emissions in DK2 calculated on the basis of the demand, based on four and six splines (green).

### 3.4. Interaction terms

The interactions [25] between the explanatory variables is modeled as

$$IA(\mathbf{z}_{i,t}, \mathbf{z}_{j,t}) = \left[ \mathbf{z}_{i,t} \quad \mathbf{z}_{j,t} \quad \mathbf{z}_{i,t} \cdot \mathbf{z}_{j,t} \right]^T, \quad (8)$$

where the product  $\mathbf{z}_i \cdot \mathbf{z}_j$  denotes the mutual interaction of the  $(i, j)$  pair. The coefficients are denoted as vector  $\beta^{IA}$ .

Interactions were also represented by splines to refine the non-linearity, i.e.

$$bs(IA(\mathbf{z}_{i,t}, \mathbf{z}_{j,t})) = \begin{bmatrix} bs_0(\mathbf{z}_{i,t}) & bs_0(\mathbf{z}_{j,t}) & bs_0(\mathbf{z}_{i,t} \cdot \mathbf{z}_{j,t}) \\ bs_1(\mathbf{z}_{i,t}) & bs_1(\mathbf{z}_{j,t}) & bs_1(\mathbf{z}_{i,t} \cdot \mathbf{z}_{j,t}) \\ bs_2(\mathbf{z}_{i,t}) & bs_2(\mathbf{z}_{j,t}) & bs_2(\mathbf{z}_{i,t} \cdot \mathbf{z}_{j,t}) \\ bs_3(\mathbf{z}_{i,t}) & bs_3(\mathbf{z}_{j,t}) & bs_3(\mathbf{z}_{i,t} \cdot \mathbf{z}_{j,t}) \end{bmatrix}^T, \quad (9)$$

with the corresponding coefficients the matrix  $\beta^{bs(IA)}$ .

Note that explanatory variables generally change in time; for example, the production in DK2 has a clear daily pattern and its corresponding linear regression coefficient will vary accordingly. This can be expressed

as interactions with the time variables (hour, week, month). A separate matrix  $\tau$  is thus defined to group the periodic as well as nonlinear character of the time variables (Appendix B).

#### 4. Statistical selection and refinement of models

##### 4.1. Cross validation strategy

Throughout the study rolling forward cross-validation is used, where data is divided into eight sets each consisting of training, validation and testing data [26]. There are eight rounds of cross-validation, where the training data is increased by one set in each round (the validation set of a previous round becomes part of the training data in the next one), and the validation and testing sets are always new independent data sets.

The cross-validation is done by averaging the Root Mean Squared Error (RMSE) over the eight validation sets. The final comparison of models is made on the test sets. The data range is 15 months, from 19th December 2016 to 18th March 2018. This is a smaller period from the one mentioned in seasonality, Sec. 3.2. The reason is that many of the explanatory variables have limited historical data and hence cannot be cross-validated further back.

##### 4.2. Feature selection techniques

Feature selection is necessary to remove co-linearity and overfitting in linear regression. The co-linearity happens when two or more explanatory variables are linearly dependent or highly correlated. When this happens the condition number of the matrix  $\mathbf{X}$  is lowered, making the determinant of it close to zero, and thus inverting it results in large numerical errors. The overfitting is related to a large number of model parameters used to fit training data, which then causes poor model performance on test data, as seen from the yellow points in Figure 6. In the extreme, high order polynomials are found to perfectly fit scattered data that in reality follow a simple, say, linear trend.

Two methods are used for feature selection. In the first, Least Absolute Shrinkage and Selection Operator (LASSO) algorithm, a penalty term is added to the objective function  $S$  that is to be minimized

$$S(\beta) = \|\mathbf{y} - \mathbf{X}\beta\|^2 + \lambda\|\beta\|_1, \quad (10)$$

where  $\lambda$  is the penalty parameter and the subscript 1 indicates the L1 norm; the larger the L1 norm of the coefficient  $\beta$ , the larger the penalty [27]. This reduces parameter estimates to zero, hence its name "Shrinkage".

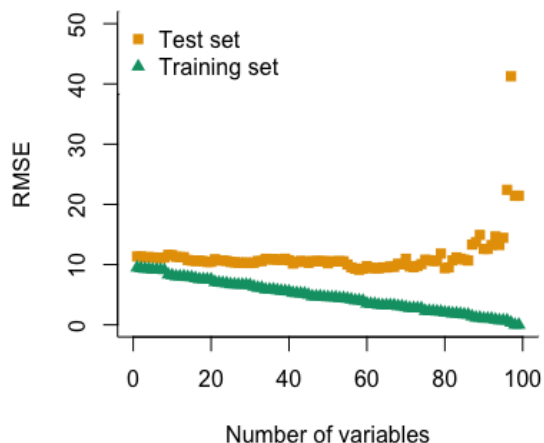


Figure 6: RMSE vs number of variables. Generated using 200 observations from 100 independently simulated explanatory variables, all with a certain degree of correlation to a simulated response variable.

The penalty term  $\lambda$  is tuned using the average RMSE from the validation sets defined in the above eight-fold cross-validation strategy. The higher the value of  $\lambda$  the more the  $\beta$ 's shrink towards zero and therefore fewer variables will be selected.  $\lambda$  is slowly decreased from an initial high value until the optimal value is reached when model performance stops improving. If the performance starts to decrease, the model is over-fitted.

Before applying LASSO, highly correlated variables are removed manually to reduce the computation time of the LASSO regression.

The second method is the forward feature selection algorithm in which new variables are added to the best models and tested for improvement, Appendix C. These methods are used for this study in the following order

1. Highly correlated variables ( $\rho > 0.99$ ) are removed.
2. LASSO regression is applied.
3. The forward selection algorithm is applied.
4. Step two and three repeated with updated variables until convergence (number of variables does not decrease anymore).

##### 4.3. Residual correction

When the model does not predict well the test data, residuals are non-random and become auto-correlated

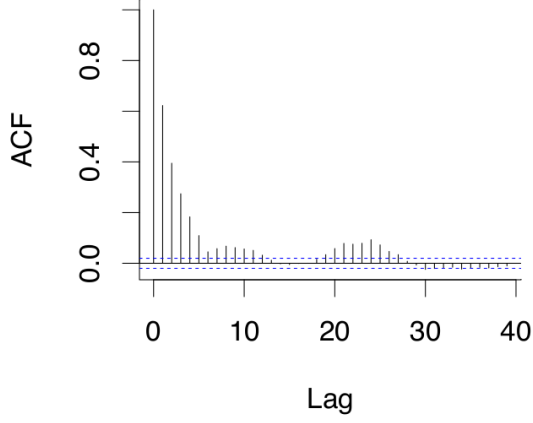


Figure 7: Auto Correlation Function of the residuals from the combined model -  $y$  is the average CO<sub>2</sub> emission intensity and  $h = 6$  hours - on the training set from CV set 8.

(correlated to lagged versions of itself). Thus corrections to the model need to be made to account for the correlation.

Residual auto-correlations are checked with the ACF (Auto Correlation Function) that reveals any linear dependencies in the residuals. In Figure 7, residuals of the compound model (see later, Sec. 4.5) of the average CO<sub>2</sub> emission intensity on six-hour horizon is shown. There are high correlations up until the lag of six hours, and also a smaller correlation around 24 hours due to seasonality.

The residual was modeled independently with Auto Regressive Moving Average (ARIMA) model [28]

$$\mathbf{Y}_t = \epsilon_t + \psi_1 \epsilon_{t-1} + \dots + \psi_q \epsilon_{t-q} - \phi_1 \mathbf{Y}_{t-1} - \dots - \phi_p \mathbf{Y}_{t-p}, \quad (11)$$

containing  $p$  lagged values (AR part) and  $q$  errors of previous observation of the moving average (MA). Considering an AR model the prediction errors are obtained as

$$\mathbf{Y}_{t+h} - \hat{\mathbf{Y}}_{t+h|t} = \epsilon_{t+h} + \psi_1 \epsilon_{t-1+h} + \dots + \psi_{h-1} \epsilon_{t+1}. \quad (12)$$

The models of this type are denoted  $ARIMA_{p,d,q}$  process, where parameters  $p$ ,  $d$ ,  $q$  refer to the AR, I and MA part, respectively - I is an integrating term used to make the data stationary (mean and variance are constant over time). In this study an extension is used,  $ARIMA_{(p,d,q)(P,D,Q)_M}$ , where  $P, D$  and  $Q$  are the seasonal

parameters and  $M$  is the length of the season. The parameters are fitted using information from the ACF and using the built-in function `auto.arima()` in R, which automatically selects the model with the best fit. In this case, it is found to be seasonal  $ARIMA_{(3,0,0)(0,1,2)_{24}}$  and  $ARIMA_{(5,1,0)(2,0,0)_{24}}$  models that removes all significant correlations for the average and marginal emission models with  $h = 6$  hours.

The prediction of the residual model is added to Eq. 3 to obtain the final prediction

$$\mathbf{y}_{t+h|t} = \hat{\mathbf{y}}_{t+h} + ARIMA(\mathbf{y}_{1:t} - \hat{\mathbf{y}}_{1:t}, h) + \epsilon_{t+h|t} \quad (13)$$

where  $\mathbf{y}_{t+h|t}$  and  $\epsilon_{t+h|t}$  are the response variable and error at time  $t + h$  given (residual) information at  $t$ , and  $\hat{\mathbf{y}}_{t+h}$  is the prediction from the linear model.  $ARIMA(\mathbf{y}_{1:t} - \hat{\mathbf{y}}_{1:t}, h)$  is the ARIMA predicted error in line with Equation 12, however with an extended version. Note that once obtained, the ARIMA model is used for predicting the residuals at the horizon  $h$ . Finally, the uncertainty of the model is evaluated by applying the 95% prediction interval. This is in accordance to the equations defined in [28], and is applied easily in R through built-in options in `lm()` and `arima()`.

The described residual correction can improve both the forecast for the specific horizon and the forecast for the lower horizons at the same time. Besides this, it also gives more consistent results since the cross-validation sets do not stand out (the variance of the errors becomes smaller).

#### 4.4. Base models

The formalism of Sections 3, 4.1 and 4.2 is used to assemble different models of increasing complexity, listed below. Starting point is the model

- **M0**

$$\mathbf{y}_{t+h} = \beta_0 + \sum_{i=1}^{477} \beta_i(\mathbf{z}_{i,t+h}) + \epsilon_{t+h},$$

where the 477 refer to the columns in  $\mathbf{Z}$ .

- **M1**

$$\begin{aligned} \mathbf{y}_{t+h} = & \beta_0 + \sum_{i=1}^{951} \beta_i(\mathbf{z}_{i,t+h}) + \sum_{i=1}^{15} \beta_i^{\tau}(\tau_{i,t+h}) \\ & + \sum_{i=1}^{477} \beta_i^{bs} \cdot bs(\mathbf{z}_{i,t+h}) + \sum_{i=1}^{477} \beta_i^{ns} \cdot ns(\mathbf{z}_{i,t+h}) + \epsilon_{t+h} \end{aligned}$$

includes time variables  $\tau$  and non-linearities (splines).

- **M2**

$$\mathbf{y}_{t+h} = \beta_0 + \sum_{i=1}^{50} \sum_{j=1 \wedge j \neq i}^{50} \beta_{i,j,k}^{IA} \cdot IA(\mathbf{z}_{i,t+h}, \mathbf{z}_{j,t+h})_k + \sum_{i=1}^{50} \sum_{k=1}^{15} \beta_{i,k,l}^{IA} \cdot IA(\mathbf{z}_{i,t+h}, \tau_{k,t+h})_l + \epsilon_{t+h}$$

is based on the reduced number of features (maximum 50, obtained from the LASSO regression of the Model 0 and ranked based on the size of their linear regression coefficients), and their interactions defined by the vector  $IA$ , Equation 8.

- **M3**

$$\mathbf{y}_{t+h} = \beta_0 + \sum_{i=1}^{50} \sum_{j=1 \wedge j \neq i}^{50} \beta_{i,j,l,m}^{bs(IA)} \cdot bs(IA(\mathbf{z}_{i,t+h}, \mathbf{z}_{j,t+h}))_{l,m} + \sum_{i=1}^{50} \sum_{k=1}^{21} \beta_{i,k,l,m}^{bs(IA)} \cdot bs(IA(\mathbf{z}_{i,t+h}, \tau_{k,t+h}))_{l,m} + \epsilon_{t+h},$$

is also based on the reduced features, but with the interactions defined via matrix  $bs(IA)$  of Equation 9.

The feature selection procedure defined in Sec. 4.2 is applied to **M1**, **M2** and **M3** and reduces the number of variables to 10-30 depending on the horizon.

#### 4.5. Weighted average model

The final model was the weighted average of the above Models 1-3. The weights were based on the performance of models on eight validation sets. The Softmax function was used

$$w_i = \left( \frac{\exp x_i}{\sum_{j=1}^n \exp x_j} \right),$$

where  $w$  is the weight vector,  $x$  is a vector representing the average  $RMSE$  scores of the included models on the validation sets and  $n$  is the number of models included. Compared to the flat weight  $w_i = \left( \frac{x_i}{\sum_{j=1}^n x_j} \right)$ , the Softmax function gives more weight to the good models and almost neglects the bad ones due to the exponential term.

In Table 1 the performance of the three models are shown for the average and marginal  $CO_2$  emission intensity (the response variable  $\mathbf{y}$ ) when using the forecast horizon of  $h = 6$  hours. Listed are  $RMSE$ 's..

Model **M2** with the linear interaction terms is the best model for both the average and marginal emissions implying the importance of variable interactions. The marginal emissions have lower errors compared to the average emissions, suggesting that the marginal value is

Average	M1	M2	M3	M <sub>WA</sub>
<b>Validation</b>	39.63	38.97	39.19	<b>37.87</b>
<b>Test</b>	41.13	39.54	40.37	<b>38.45</b>
<b>Weights</b>	<b>0.22</b>	<b>0.43</b>	<b>0.35</b>	
<hr/>				
Marginal				
<b>Validation</b>	11.06	8.77	10.03	<b>8.57</b>
<b>Test</b>	11.94	9.94	10.83	<b>9.63</b>
<b>Weights</b>	<b>0.07</b>	<b>0.73</b>	<b>0.2</b>	

Table 1: Root-Mean-Squared Error [RMSE] of the average and marginal emissions for models **M1-M3**, with a prediction horizon of 6 hours. The weights are calculated using the RMSE values from the validation sets. **M<sub>WA</sub>** is the resulting weighted average model. Units:  $\left[ \frac{\text{kgCO}_2\text{-eq}}{\text{MWh}} \right]$ .

easier to predict - it is less influenced by highly uncertain (weather dependent) variables as already remarked.

The weighted average model **M<sub>WA</sub>** is constructed by combining the models with the Softmax weights: the  $RMSE$  in the *test* set becomes 38.56 and 9.63  $\left[ \frac{\text{kgCO}_2\text{-eq}}{\text{MWh}} \right]$  for the average and marginal emissions, respectively. This is only a slight improvement to the  $RMSE$  compared to **M2**, because of the large weight assigned to it.

#### 4.6. The compound model

The pure **M<sub>WA</sub>** model of the previous section - i.e., without the residual correction of Sec. 4.3 - was used for forecasts on each individual horizon  $h = 1, 2, \dots, 24$  hours (denoted **M<sub>WA1...</sub>**). Trials showed that ARIMA model performed on  $h = 6$  hours corrects the residuals on all earlier horizons; however, individual **M<sub>WA1,2</sub>** as well as **M<sub>WA7-24</sub>** outperform the **M<sub>ARIMA6</sub>** on their corresponding horizons for the average emissions. This is so both because individual models designed for specific horizons may perform better, and also ARIMA prediction converge towards the average as the prediction horizon increases, thus being less suitable for longer horizons. The models for different horizons with the corresponding  $RMSE$ s for the average and marginal emissions are summarized in Table 2.

Note from the table that  $RMSE$  during 7-24 hours becomes almost stationary. The reason is that on longer horizons, current information, say, on production data (available through short-term forecasts, Appendix A) affects the predictions and associated uncertainties much less than the available long-term information of e.g. weather. On shorter horizons, 0-6 hours,  $RMSE$  depends on short-term data and gradually increases in time until reaching the stationary value. The features can be seen in Figure 8.



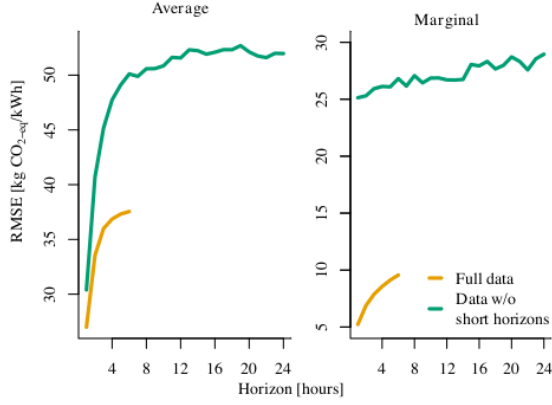


Figure 8: RMSE of CO<sub>2</sub> emission predictions vs prediction horizon. When short term forecasts are excluded from the data set (green line), the error increases (31% on the 6th hour horizon for the average emissions). The error stays at around the same level for  $h > 6$  hours (49-52  $\left[\frac{\text{kgCO}_2\text{-eq}}{\text{MWh}}\right]$ ). In the marginal emissions, the short term forecasts are more important for the predictions and without them the error is high even on the short horizons.

	Average			Marginal	
	$\mathbf{M}_{\text{WA}1,2}$	$\mathbf{M}_{\text{ARIMA}6}$	$\mathbf{M}_{\text{WA}7-24}$	$\mathbf{M}_{\text{ARIMA}6}$	$\mathbf{M}_{\text{WA}7-24}$
Horizons [h]	1-2	3-6	7-24	1-6	7-24
RMSE $\left[\frac{\text{kgCO}_2\text{-eq}}{\text{MWh}}\right]$	27.0-	36.0-	49.9-	5.2-9.6	26.2-
	33.5	37.6	52.0		29.0

Table 2: Performance of the final best models for particular horizons on test set. The average emission model is thus on  $h = 6$  improved from 38.45 to 37.6  $\left[\frac{\text{kgCO}_2\text{-eq}}{\text{MWh}}\right]$  and the marginal is improved from 9.7 to 9.6  $\left[\frac{\text{kgCO}_2\text{-eq}}{\text{MWh}}\right]$  (from Table 1).

The final compound model  $\mathbf{M}_{24}$  used for the 24-hour forecast of the emissions is the sum of the corresponding triplets of Table 2. When the short-term forecast data is included in  $\mathbf{M}_{24}$ , the RMSE are naturally smaller than when it is excluded (yellow vs. green line, Figure 8). The improvement is by 31% (within the 0-6 hour horizon, for which short-term forecasts are available).

## 5. Selected results

### 5.1. The interaction coefficients $\beta^{IA}$

Of the three base models,  $\mathbf{M}_2$  that includes the interaction terms is the most accurate. The parameters and results of this model are discussed here for both the average and marginal emissions. Since the compound model puts the most weight on  $\mathbf{M}_2$ , the analysis applies to this model as well.

$\mathbf{z}_i$	$\mathbf{z}_j$	$\beta_i^{IA}$	$\beta_j^{IA}$	$\beta_{i,j}^{IA}$
<b>Average</b>				
Production in DK2.	Spline (midday).	57.3	1.7	-4.7
Production in DK2.	Daily pattern.	57.3	3.3	-2.3
Net export from DK1 to DE .	Spline (midday).	-11.4	1.7	3.8
Wind speed in DK2.	Net export from DK2 to DE.	-11.8	-8.2	4.2
Offshore wind in DK2.	Net export from SW3 to SW4.	-16.7	-14.1	2.6
<b>Marginal</b>				
Net export to DE from DK1 - exp. solar radiation in DE.	Wind speed in DE.	10.9	-6.1	-1.9
solar radiation in DE.	Net export from DK1 to NO2 - exp.	-1.7	-1.3	1.8
solar radiation in DE.	Net export from DK1 to SW3 - exp.	-1.7	-3.7	1.4
Net export from DK1 to SW3.	Net export from SW3 to SW4 - exp.	-8.0	0.54	-1.5
Net export from SW3 to SW4.	Demand in SW4	11.4	2.7	-0.9

Table 3: Strongest interactions in Model 2 for the average and marginal emission.  $z, i$  and  $j$  correspond to the terms in Equation 8. Note the original variables in Appendix A use "net import" rather than "net export". Switched here, to make it easier to relate to.

The five largest  $\beta$  coefficients of  $\mathbf{M}_2$  are featured in Table 3, both for the average and the marginal emissions.  $\beta_i^{IA}$ ,  $\beta_j^{IA}$  and  $\beta_{i,j}^{IA}$  from the table refer to the first, second and third coefficient of the interaction vector  $\beta^{IA}$  in Equation 8. The minus sign indicates opposing trends. To understand the table properly please realize the following for interaction terms:

$$\beta_i^{IA} \mathbf{z}_i + \beta_j^{IA} \mathbf{z}_j + \beta_{i,j}^{IA} \mathbf{z}_{i,j} \quad (14)$$

$$= (\beta_i^{IA} + \beta_{i,j}^{IA} \mathbf{z}_j) \cdot \mathbf{z}_i + \beta_j^{IA} \mathbf{z}_j \quad (15)$$

$$= (\beta_j^{IA} + \beta_{i,j}^{IA} \mathbf{z}_i) \cdot \mathbf{z}_j + \beta_i^{IA} \mathbf{z}_i \quad (16)$$

This means the final coefficient for e.g.  $\mathbf{z}_i$  becomes  $(\beta_i^{IA} + \beta_{i,j}^{IA} \mathbf{z}_j)$ . The column for  $\beta_i^{IA}$  is thus the coefficient for  $\mathbf{z}_i$  if  $\mathbf{z}_j$  is zero and vice versa. The column for  $\beta_{i,j}^{IA}$

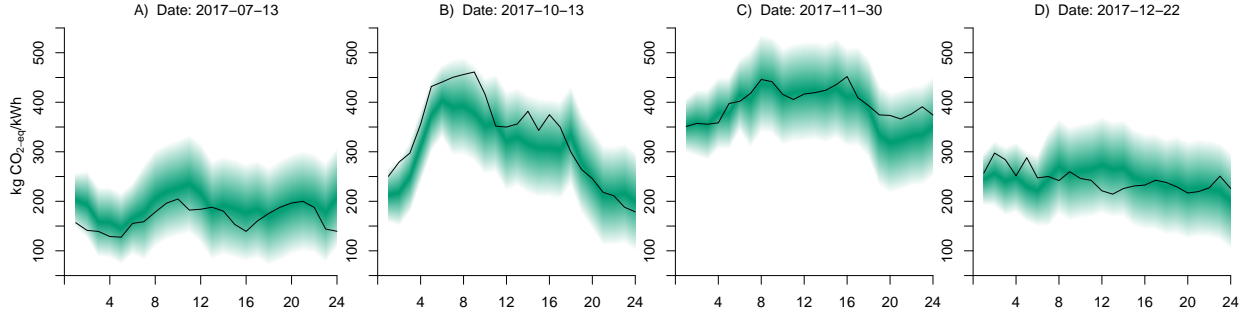


Figure 9: Average CO<sub>2</sub> emissions: The final forecast submitted at midnight with a 24 hour horizon for 8 different days. The real CO<sub>2</sub> emission intensity is the thick line, and the colored areas are bounding the 95% confidence interval and the point prediction.

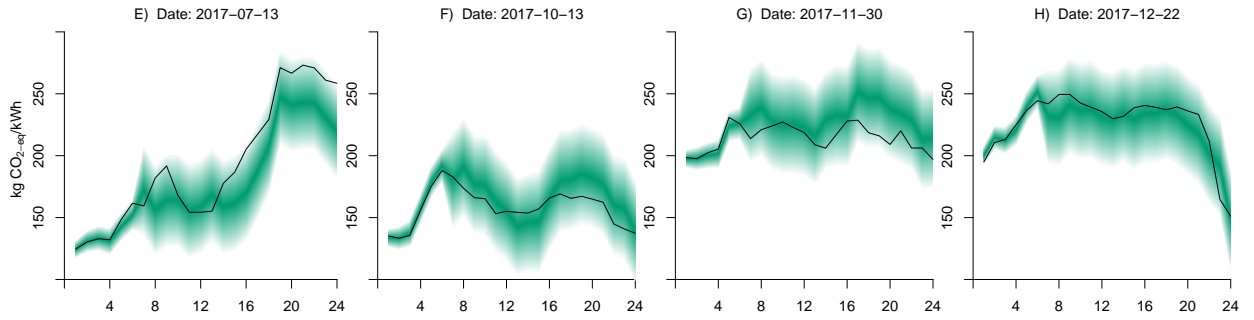


Figure 10: Marginal CO<sub>2</sub> emissions: The final forecast submitted at midnight with a 24 hour horizon for 8 different days. The real CO<sub>2</sub> emission intensity is the thick line, and the colored areas are bounding the 95% confidence interval and the point prediction.

is the linear coefficient that explains the change in the final coefficient for both  $\mathbf{z}_i$  and  $\mathbf{z}_j$ .

In the **average** emission model, the production in DK2 has the largest coefficient, 57.26: more production  $\rightarrow$  higher emissions, as also concluded in Figure 2. It is interacting with the daily periodicities: the mid-day spline (row one in Table 3, see Appendix B), and the daily pattern from Figure 4 (row five). The negative interaction coefficients  $\beta_{i,j}^{IA}$ , -4.7 and -2.3, via Equation 15, mean the impact from the power production decays and reaches a minimum at around noon. This is because other factors start to influence the emissions more.

The midday spline is again interacting with the net export from DK1 to DE (row three) that has a negative coefficient -11.4; here,  $\beta_{i,j}^{IA}$  is positive, 3.8, and suggests the least impact during the midday, concluded from Equation 15.

The wind speed in DK2 (row four) has a negative coefficient: the average emissions will decrease as more wind power enters the grid. The wind speed is positively interacting with net export from DK2 to DE (4.2) hence the final wind speed coefficient becomes larger, approaching zero, because the wind power is being exported rather than being used to decrease the emissions in DK2.

The offshore wind power production in DK2 (row five) also has a negative coefficient, -16.7, and is interacting with the net export from SW3 to SW4 - this is primarily a one way inter-connector explicitly exporting to SW4. The emissions in DK2 decrease as exports increase - this is expected because SW3 is in possession of all Sweden's nuclear power [23], and that is exported further into DK2.  $\beta_{i,j}^{IA}$  is positive and the final wind power coefficient will approach zero as the export increase, Equation 15: the nuclear power takes part in the overall emissions in DK2 making the wind power contribution account for relatively less.

In the **marginal** emission model, note that none of the listed variables describe the grid in DK2. All are external influences from neighboring bidding zones and include net exports in all interactions. They are created because the net exports often only impact the emissions in one trading direction. Recall from Figure 3, only the exports from DK1 to DE has an impact on the emissions. In this model, the coefficient of that trading pair is positive (row six) fitting the exports. To force the coefficient closer to zero during import, it is negatively interacting with the wind speed in DE,  $\beta_{i,j}^{IA} = -1.9$ . DK1 will import if DE has wind power to offer and the final coefficient for the net export decrease.

The solar radiation in DE is both interacting with net exports from DK1 to NO2 and net exports DK1 to SW3 ( $\beta_{i,j}^{IA} = 1.8$  and  $1.4$  respectively). The two interactions explain the same phenomenon: the emissions decrease as DK1 exports to SW3 and NO2 because  $\beta_j = -1.3$  is negative. However, when solar radiation in DE is high DK1 will again import from DE and the final coefficients for the net exports increase approaching zero because of the positive interaction coefficients.

The net export from DK1 to SW3 (row 9) is interacting with the net export from SW3 to SW4 (exponential term) too with  $\beta_{i,j}^{IA} = -1.5$ : export from SW3 to SW4 decreases the impact from the net export from DK1 to SW3. That is because the marginal emissions in SW3 increase as their total export increases. The net export from SW3 to SW4 is the most significant variable and is negatively interacting with the demand in SW4  $\beta_{i,j}^{IA} = -0.9$ : the final coefficient for the net import in SW4 from SW3 decreases as the demand increases, most likely because in this case power is consumed in SW4 rather than exported further into DK2.

### 5.2. CO<sub>2</sub> emissions: the forecast of $M_{24}$

A few examples for the average and marginal emission forecasts by model  $M_{24}$  are shown in Figure 9 and Figure 10 to illustrate the performance, where a 24-hour horizon forecast is released at midnight. Note, the dates in the plots for the average and marginal emissions are identical.

Comparing the plots (**average** emissions) to the daily pattern in Figure 4, plot B and C fit best with the highest emissions during the day. Plot B peaks already in the morning slightly higher than the predictions and in plot C there is expected a lower decrease than observed in the evening. To a certain degree, daily patterns are often expected, so when the real observations differ too much, the accuracy decreases: in plot D, the emissions had a downward going trend all day, but it was expected to peak at around noon and then decrease. Plot A illustrates a day with irregularities too where the trend is captured to an adequate degree.

The **marginal** emissions have a much slimmer confidence interval than the average emissions due to the higher accuracy. Plot F differs the least from the average daily pattern and the prediction shows this too. The predictions in Plot G had a low accuracy because of irregular and small spikes. Plot E and H are good examples for a control mechanism: in plot A, the emissions are expected to increase in the evening, so it is encouraged to schedule flexible demand as early as possible before the emissions increase. In plot D, the opposite

is seen: here, the demand should be shifted to the late evening where the emissions decreased.

## 6. Conclusion

From data collected and supplied by Tmrow IVS, new forecasting models for average and marginal CO<sub>2</sub> emissions in the European electricity grid are developed using linear regression and residual correction. A machine learning methodology to systematically select the important variables that best fit the desired variable is presented.

It is found that interactions between the explanatory variables are important: large coefficients are found for net imports, and time-dependence is least pronounced during midday or midnight. Interestingly, none of the most important variable related to the marginal emissions in DK2 were local (DK2) variables - all contributions came from neighboring bidding zones (DK1, DE, SW3 (indirect) and SW4). This suggests that the marginal generator is effectively supplied from the import, in agreement with [19] which found mainly import from SW4.

The study aimed to provide a tool that can help electricity consumers schedule their load to minimize CO<sub>2</sub> emissions. This was accomplished by forecasts of emissions 24 hours ahead, which provides a basis for decision making for load scheduling. The average and marginal emissions follow different patterns that can be exploited for different applications. The marginal CO<sub>2</sub> emissions are valid for small changes in demand and are therefore the signal to use when scheduling home appliances. The average emissions are useful for evaluating total electricity system emissions but should not be used as a control signal.

To evaluate the usefulness of the marginal emission forecast, testing on various flexible applications, e.g. heat pumps, electric cars, etc. should be conducted in the future. Results from this can indicate if there is a need for further model improvements. The marginal emission estimates used in this study cover most situations but there are still limitations as mentioned in the footnote 1 (page 1). Further studies are needed to incorporate the weather dependent generators as marginals to fully understand the concept.

## 7. Acknowledgement

We are thankful for Tmrow IVS who has provided the data used in this study (including emission calculations for the bidding zone DK2). The work is supported

through the project Smart Cities Accelerator 2016-2020 funded by the EU program Interreg resund-Kattegat-Skagerrak, the European Regional Development Fond and the CITIES project (DSF1305-00027B).

## References

- [1] CO<sub>2</sub> emissions statistics, <https://www.iea.org/>, [Online; accessed 6-Nov-2019] (2017).
- [2] Status of power system transformation 2019: Power system flexibility, [https://www.iea.org](https://www.iea.org/), [Online; accessed 6-Nov-2019] (2019).
- [3] Z. Liang, J. Liang, C. Wang, X. Dong, X. Miao, Short-term wind power combined forecasting based on error forecast correction, *Applied Energy* 119 (2016) 215–226.
- [4] J. Wang, T. Niu, H. Lu, Z. Guo, W. Yang, P. Du, An analysis-forecast system for uncertainty modeling of wind speed: A case study of large-scale wind farms, *Applied Energy* 211 (2018) 492–512.
- [5] P. Bacher, H. Madsen, H. A. Nielsen, An analog ensemble for short-term probabilistic solar power forecast, *Solar Energy* 10 (2009) 1772–1783.
- [6] S. Alessandrini, L. D. Monache, S. Sperati, G. Cervone, An analog ensemble for short-term probabilistic solar power forecast, *Applied Energy* 157 (2015) 95–110.
- [7] D. Keles, J. Scelle, F. Paraschiv, W. Fichtner, Extended forecast methods for day-ahead electricity spot prices applying artificial neural networks, *Applied Energy* 162 (2016) 218–230.
- [8] Z. Yang, L. C. and Li 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.
- [9] U. Wagner, W. Mauch, S. von Roon, Das merit-order-dilemma der emissionen, Tech. rep., Forschungsstelle fr Energiewirtschaft e.V (2002).
- [10] A. Regett, F. Baing, J. Conrad, Emission assessment of electricity: Mix vs. marginal power plant method, 15th International Conference on the European Energy Market (EEM) (2018).
- [11] Voorspools, D’haeseleer, An evaluation method for calculating the emission responsibility of specific electric applications, *Energy Policy* 28 (2006) 967–980.
- [12] Voorspools, D’haeseleer, The influence of the instantaneous fuel mix for electricity generation of the corresponding emissions, *Energy* 25 (2000) 1119–1138.
- [13] Marnay, Fisher, Murtishaw, Phadke, Price, Sathaye, Estimating carbon dioxide emissions factors for the california electric power sector, Tech. rep., Lawrence National Laboratory, Berkley USA (2002).
- [14] R. Bettle, C. Pout, E. Hitchin, Interactions between electricity-saving measures and carbon emissions from power generation in england and wales, *Energy Policy* 34 (2006) 3434–3446.
- [15] A. Hawkes, Estimating marginal co<sub>2</sub> emissions rates for national electricity systems, *Energy Policy* 38 (2010) 5977–5987.
- [16] Rekkas, Uk marginal powerplant and emissions factors, Master’s thesis, Imperial College London (2005).
- [17] Hadland, Marginal emissions factors for the united kingdom electricity system, Master’s thesis, Imperial College London (2009).
- [18] J. Clau, S. Stinner, C. Solli, K. B. Lindberg, H. Madsen, L. Georges, Evaluation method for the hourly average co<sub>2</sub>-eq intensity of the electricity mix and its application to the demand response of residential heating, *energies* 12 (2019).
- [19] O. Corradi. <https://medium.com/electricitymap/using-machine-learning-to-estimate-the-hourly-marginal-carbon-intensity-of-electricity-49eade43b421> [online] (2018).
- [20] J. Bialek, Tracing the flow of electricity, Vol. 143, IEEE, 1996.
- [21] D. Kirschen, R. Allan, G. Strbac, Contributions of individual generators to loads and flows, *Transactions on Power System* 12 (1997) 52–60.
- [22] A. Heydari, D. A. Garcia, F. Keynia, F. Bisegna, L. D. Santoli, Renewable energies generation and carbon dioxide emission forecasting in microgrids and national grids using grnn-gwo methodology, *Applied Energy* 159 (2019) 154–159.
- [23] Elomrde 1-4 (sn1-4) - statistik per mnad 2017), <https://www.svk.se/>, [Online; accessed 1-Nov-2019] (2017).
- [24] R. T. Trevor Hastie, J. Friedman, *The Elements of Statistical Learning*, Springer Series in Statistics, 2017.
- [25] Porter, A. S, West, Multiple regression: Testing and interpreting interactions, *Journal of the Royal Statistical Society: Series D (The Statistician)* 43 (1994) 453.
- [26] M. Y. Hu, G. Zhang, C. X. Jiang, B. E. Patuwo, A cross-validation analysis of neural network out-of-sample performance in exchange rate forecasting, *Decision Sciences* 30 (1999) 197–215.
- [27] R. Tibshirani, Regression shrinkage and selection via the lasso: a retrospective, *Journal of the Royal Statistical Society. Series B (Statistical Methodology)* 73 (2011) 273–282.
- [28] H. Madsen, *Time Series Analysis*, Chapman & Hall/CRC - Taylor & Francis Group, 2007.

## Appendix A. Explanatory variables

Here, all explanatory variables are listed by data set for bidding zone DK2. For each variable in the data sets, it is indicated whether the variable is used to create models for  $h \leq 6$  hours or  $h > 6$  hours.

### Short Term Forecasts

	$h \leq 6$ hours	$h > 6$ hours
dewpoint	X	
precipitation	X	
solar	X	
temperature	X	
price	X	
production	X	
consumption	X	
wind_speed	X	
wind_direction_x	X	
wind_direction_y	X	
power_net_import_DK-DK1	X	
power_net_import_DE	X	
power_net_import_SE-SE4	X	
power_net_import_SE	X	

6 hours Forecasts (updated every sixth hour) provided by Tmrow IVS. Imports are originally coming from ENTSO-E who collect the data from individual Transmission System Operators (TSO's).

### Weather Forecasts

	$h \leq 6$ hours	$h > 6$ hours
dewpoint_mean_value		X
precipitation_mean_value		X
solar_mean_value		X
temperature_mean_value		X
wind_mean_value		X

52 hours Forecasts (updated every sixth hour) provided by Tomorrow. Originally created by GFS - Global Forecasting System

### Market Data (Nordpool)

	$h \leq 6$ hours	$h > 6$ hours
solar_power	X	X
wind_power_offshore	X	X
wind_power_onshore	X	X
production		X
consumption		X
spot_price		X

Published at 6pm CET covering the whole next day. Reported in at 12pm CET, and technically available at that time.

## Real Time Data

	$h \leq 6$ hours	$h > 6$ hours
carbon_intensity	X	X
carbon_intensity_production	X	X
carbon_intensity_import	X	X
carbon_rate	X	X
total_production	X	X
total_storage	X	X
total_discharge	X	X
total_import	X	X
total_export	X	X
total_consumption	X	X
power_origin_%_fossil	X	X
power_origin_%_renewable	X	X
power_production_biomass	X	X
power_production_coal	X	X
power_production_gas	X	X
power_production_hydro	X	X
power_production_nuclear	X	X
power_production_oil	X	X
power_production_solar	X	X
power_production_wind	X	X
power_production_geo	X	X
power_production_unknown	X	X
power_origin_%_biomass	X	X
power_origin_%_coal	X	X
power_origin_%_gas	X	X
power_origin_%_hydro	X	X
power_origin_%_nuclear	X	X
power_origin_%_oil	X	X
power_origin_%_solar	X	X
power_origin_%_wind	X	X
power_origin_%_geo	X	X
power_origin_%_unknown	X	X
power_origin_%_hydro	X	X
carbon_origin_%_biomass	X	X
carbon_origin_%_coal	X	X
carbon_origin_%_gas	X	X
carbon_origin_%_hydro	X	X
carbon_origin_%_nuclear	X	X
carbon_origin_%_oil	X	X
carbon_origin_%_solar	X	X
carbon_origin_%_wind	X	X
carbon_origin_%_geo	X	X
carbon_origin_%_unknown	X	X
carbon_origin_%_hydro	X	X
power_net_discharge_hydro	X	X
power_net_import_DK-DK1	X	X
power_net_import_DE	X	X
power_net_import_SE-SE4	X	X
power_net_import_SE	X	X

Provided by Tomorrow and originally created by GFS.

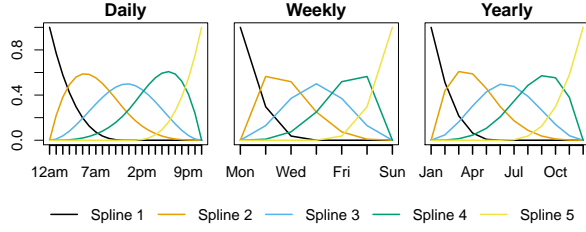


Figure B.11: Periodic splines - Daily, weekly and yearly. Spine 3 in Daily is referred to as the midday spline. The outer splines (1 and 5) are not included in the model because they can be misleading.

## Appendix B. Periodic time variables

Time variable matrix  $\tau$  is defined as:

$$\begin{aligned}
 \tau_{\text{hour}} &= t \bmod 24 \\
 \tau_w &= \text{weekday}(t) \\
 \tau_m &= \text{month}(t) \\
 \tau_{\sin, \text{hour}} &= \sin(\tau_{\text{hour}}) \\
 \tau_{\sin, w} &= \sin(\tau_w) \\
 \tau_{\sin, m} &= \sin(\tau_m) \\
 \tau_{bs, \text{hour}, t} &= \left[ bs_0(\tau_{\text{hour}, t}) \quad bs_1(\tau_{\text{hour}, t}) \quad \dots \quad bs_{n, df-1}(\tau_{\text{hour}, t}) \right] \\
 \tau_{bs, w, t} &= \left[ bs_0(\tau_w, t) \quad bs_1(\tau_w, t) \quad \dots \quad bs_{n, df-1}(\tau_w, t) \right] \\
 \tau_{bs, m, t} &= \left[ bs_0(\tau_m, t) \quad bs_1(\tau_m, t) \quad \dots \quad bs_{n, df-1}(\tau_m, t) \right],
 \end{aligned}$$

where hour, w and m denote the hour, weekday and month of the datetime  $t$ , respectively.  $\tau_{bs, \text{hour}}$ ,  $\tau_{bs, w}$  and  $\tau_{bs, m}$  each represent five columns corresponding to their underlying splines.  $n_{df}$  is the number of splines which is set to 5 in this case. The periodic splines are illustrated in Figure B.11.

## Appendix C. Forward feature selection

The forward selection algorithm selects the best variables for a model and requires a good cross validation strategy to avoid overfitting.

---

### Algorithm 1 Forward feature selection

---

1. 1) Find the variable that best describes the response variable. This can be done with any best fit criteria (BIC, AIC or RMSE). This study relies on the RMSE value calculated on the validation sets of Sec. 4.1. Call this  $\text{Model}_{\text{best}}$ .

2. 2) Add a new variable  $\mathbf{x}_i$

$$\mathbf{y}_{t+h} = \beta_0 + \beta_1 \mathbf{x}_{\text{best}} + \beta_2 \mathbf{x}_i. \quad (\text{C.1})$$

Call this  $\text{Model}_{\text{new}}$ .

3. 3) Evaluate the model. If  $\text{Model}_{\text{new}}$  is better than  $\text{Model}_{\text{best}}$ , keep the newly added variable and update:  $\text{Model}_{\text{best}} = \text{Model}_{\text{new}}$ .
  4. 4) Repeat step 2 and 3 until all variables have been tested.
-