



**HAL**  
open science

# Bayesian Parameterisation of a Regional Photovoltaic Model - Application to Forecasting

Yves-Marie Saint-Drenan, Stephan Vogt, Sven Killinger, Jamie M Bright, Rafael Fritz, Roland Potthast

► **To cite this version:**

Yves-Marie Saint-Drenan, Stephan Vogt, Sven Killinger, Jamie M Bright, Rafael Fritz, et al.. Bayesian Parameterisation of a Regional Photovoltaic Model - Application to Forecasting. Solar Energy, 2019, 188, pp.760-774. 10.1016/j.solener.2019.06.053 . hal-02174688

**HAL Id: hal-02174688**

**<https://hal.science/hal-02174688>**

Submitted on 5 Jul 2019

**HAL** is a multi-disciplinary open access archive for the deposit and dissemination of scientific research documents, whether they are published or not. The documents may come from teaching and research institutions in France or abroad, or from public or private research centers.

L'archive ouverte pluridisciplinaire **HAL**, est destinée au dépôt et à la diffusion de documents scientifiques de niveau recherche, publiés ou non, émanant des établissements d'enseignement et de recherche français ou étrangers, des laboratoires publics ou privés.

# Bayesian Parameterisation of a Regional Photovoltaic Model - Application to Forecasting

Yves-Marie Saint-Drenan<sup>a,\*</sup>, Stephan Vogt<sup>b</sup>, Sven Killinger<sup>c</sup>, Jamie M. Bright<sup>d</sup>, Rafael Fritz<sup>b</sup>, Roland Potthast<sup>e</sup>

<sup>a</sup>*MINES ParisTech, PSL Research University, O.I.E. Centre Observation, Impacts, Energy, 06904 Sophia Antipolis, France*

<sup>b</sup>*Fraunhofer Institute for Energy Economics and Energy System Technology (IEE), 34119 Kassel, Germany*

<sup>c</sup>*Fraunhofer Institute for Solar Energy Systems (ISE), 79110 Freiburg, Germany*

<sup>d</sup>*Fenner School of Environment and Society, The Australian National University, 2601 Canberra, Australia*

<sup>e</sup>*Deutscher Wetterdienst (DWD), Offenbach, Germany*

---

## Abstract

Estimating and forecasting photovoltaic (PV) power generation in regions—e.g. the area controlled by the transmission system operator (TSO)—is a requirement for the operation of the electricity supply system. An accurate calculation of this quantity requires detailed information of the installed PV systems within the considered region; however, this information is not publicly available making forecasting difficult. Therefore, approximating the undefined PV systems information for use in a PV power model (parameterization) is of critical interest. In this paper, we propose a methodological approach for parameterization using time series of aggregated PV power generation. A Bayesian approach is used to overcome the significant number of unknown parameters in the problem. It regularizes the linear system by imposing constraints on deviations from an initial-guess and covariance matrices; the initial guess uses available statistical distributions of PV system metadata. The performance of the proposed forecasting approach is evaluated using estimates of the regional PV power generation from three TSOs and meteorological data from the IFS forecast model (ECMWF). The proposed forecasting approach without the Bayesian parameterization has RMSE of 3.90%, 4.25% and 4.64%, respectively; including the Bayesian approach gives RMSE of 3.82%, 4.23% and 4.51%. For comparison, we also deployed a multiple linear regression which gave RMSE of 3.89%, 4.12% and 4.54%; however, there are considerable downsides to such an approach. Our approach is competitive with TSO forecasting systems despite using far fewer input data and simpler implementation of NWP prediction. This is particularly promising as there are many avenues for future development.

*Keywords:* PV system characteristics, regional PV power, forecast, Grid integration, inverse problem

---

---

\*Corresponding author

*Email address:* [yves-marie.saint-drenan@mines-paristech.fr](mailto:yves-marie.saint-drenan@mines-paristech.fr) (Yves-Marie Saint-Drenan)

# 1. Introduction

## 1.1. Background and motivation

With more than 400 GWp of installed photovoltaic (PV) capacity globally (IEA, 2018), the integration of the large amounts of solar energy in the electricity supply system is fundamental for modernization and maintaining grid reliability. The accurate estimation of power generated by a fleet of decentralized PV systems (hereafter referred to as regional PV power generation) is crucial at several stages of energy supply and network operations.

The objective of regional PV power estimates is to replicate the actual behaviour of the aggregated power production from all unknown PV systems installed in a given area; this can take advantage of all available information (power production measurements and/or PV system meta-data). Such systems have been described by Lorenz et al. (2011) or Schierenbeck and Graeber (2010). Estimation is made difficult because only a minority of systems continuously report their generation and few PV systems make their measurements publicly available—a serious issue that is the core subject of numerous studies (Bright et al., 2017; Lorenz et al., 2011; Shaker et al., 2015; Schierenbeck and Graeber, 2010).

A prominent application requiring regional PV power generation estimates is in the online and ex-post PV power analysis for grid monitoring and balancing-group settlement (Amprion, 2019). Grid operators are responsible for the estimation of aggregated PV power produced in their control area, as well as for publicly releasing the estimates as is often mandatory by law. An example of time series of the PV power generation estimated by most European transmission system operators (TSOs) are found on the ENTSO-E website<sup>1</sup>.

Another important application group is providing the day-ahead or short-term forecasts of regional PV power generation. Forecasts are essential for energy trading (or scheduling thermal power plants), planning needs for reserve power or mitigating possible network congestion, etc. Improving the accuracy of regional PV power forecast is key because it has a positive impact on the integration costs of RES as well as on the security of supply (Killinger, 2017). Since the actual value of the regional PV power generation remains unknown, forecasting error is typically evaluated against the aforementioned regional estimates as a reference. Hence, the true goal of regional PV forecasts is to accurately predict the estimates made by the grid operators. It would be logical if forecasting methodologies used identical information as is used for the estimates; unfortunately, the data and processes involved in the estimation of the regional PV power generation are typically confidential so forecast providers must evaluate the regional forecast without it.

Within these two critical applications, two sources of uncertainty must be addressed in order to improve regional PV power generation estimates that are applied to forecasts: (i) the uncertainty resulting from the weather prediction error, and (ii) the uncertainty due to a lack of information on the installed PV systems.

---

<sup>1</sup><https://transparency.entsoe.eu>

33 Each source of uncertainty represents a considerable field of research. The goal of this paper is to address  
34 the second source of uncertainty by proposing a method to infer the parameters of a regional PV model  
35 from times series of the aggregated PV power generation. Achieving this goal would enable the forecasts  
36 access to otherwise absent data, which, as will be demonstrated, can significantly improve the estimation.

### 37 *1.2. Related work*

38 Regional PV power forecasting research is maturing; it has particularly gained increasing interest in  
39 recent years. The approaches in literature are distinguishable principally by the strategies used to overcome  
40 uncertainty arising from a lack of information about the installed PV systems. In this literature review,  
41 research on regional estimates and forecast have intentionally been considered in tandem as the same algo-  
42 rithms are conventionally used in both; for absolute clarity, our methodology produces an estimate of PV  
43 power generation, we then assess improvements when the new estimate is used in a solar forecast.

44 A first approach is to assume that the PV power measurements of all systems installed in a region are  
45 known a-priori. Thus, the regional PV power forecast can simply be obtained by summing the forecasts from  
46 each PV system. This method is detailed by [Da Silva Fonseca et al. \(2014\)](#) where it is evaluated together  
47 with other methods in a benchmark analysis. Though this approach can be very insightful, it is difficult to  
48 make operational for two key reasons: poor access to PV power measurements, and linear computational  
49 scaling with increasing number of installations.

50 [Lorenz et al. \(2008, 2011\)](#) and [Schierenbeck and Graeber \(2010\)](#) proposed a pragmatic solution to the  
51 two aforementioned issues. The aggregated regional PV power generation is estimated from only a subset of  
52 the PV installations, limited to the most representative systems. The regional estimate is then reconstructed  
53 from the subset by means of an upscaling method. In [Da Silva Fonseca Junior et al. \(2014\)](#) and [Shaker  
54 et al. \(2015\)](#), the optimal subset of reference PV systems are determined mathematically using data-reduction  
55 techniques. A prerequisite of this method is access to an archive of all PV power measurements—a condition  
56 rarely satisfied. In [Lorenz et al. \(2008\)](#) and [Lorenz et al. \(2011\)](#), the choice of the most representative  
57 reference PV systems is based on a-priori knowledge on the fleet of PV system installed in the region as well  
58 as on spatial considerations. Whilst this latter technique is better suited for an operational implementation,  
59 it requires access to PV power measurements from a large number of installations, as well as a good knowledge  
60 of the metadata of installed PV systems in the considered region.

61 None of the approaches described previously can be implemented when too few PV power measurement  
62 data is available. Another kind of model can be used when lack of data is the barrier. As described by  
63 [Saint-Drenan et al. \(2017\)](#), the principle of this alternative method is to simulate the PV power generation of  
64 a limited number of commonly occurring PV system metadata configurations (in regard to capacity, tilt and  
65 azimuth) using meteorological data. The regional PV power estimate is then obtained by a weighted sum  
66 of the simulated power values, the weights corresponding to the frequency of occurrence of the considered

67 configurations. The unknown weights can whether be evaluated on the basis of authors' experience [Schubert](#)  
68 [\(2012\)](#); [Fonseca Junior et al. \(2015\)](#) or on the basis of a statistical analysis of PV system metadata [Saint-](#)  
69 [Drenan et al. \(2017\)](#); [Killinger et al. \(2018\)](#). A drawback of this approach is that possible differences between  
70 the linear coefficients chosen for the regional forecast and those corresponding to the regional estimates may  
71 penalize the forecast accuracy. This error can be minimized using model output statistics (MOS) techniques,  
72 which correct model outputs based on the information gathered from previous forecasts ([Wilks, 2011](#)). It is  
73 preferable, however, to directly use coefficients avoiding systematic errors; this is analyzed later.

74 Systematic differences between regional forecast and estimates can be avoided through use of supervised  
75 statistical methods, whereby the parameters of the model are trained using estimates of the regional PV  
76 power generation. A first example of this kind of approach can be found in the benchmark analysis by [Da](#)  
77 [Silva Fonseca et al. \(2014\)](#), where a support vector regression is realized using weather data as the input and a  
78 time series of the aggregated PV power generation is the output. In their work, the high-dimensionality of the  
79 input data penalizes the efficiency of the approach. [Da Silva Fonseca et al. \(2014\)](#) proved this by observing  
80 a noticeable improvement by using principal component analysis (PCA) of the entire weather information  
81 and accounting for 90% of the explained variance. A drawback of these types of method is that it requires  
82 important amounts of training data to learn the dependency between the input weather information and  
83 time series of aggregated power. Furthermore, whilst certain weather variables may account for significant  
84 variance in the aggregated power, that same variable may not have the same impact in different climates;  
85 hence, training data would always be required. The efficiency of the learning phase can be improved by  
86 using hybrid models, whereby known physical dependencies are considered and unknown model parameters  
87 trained from historical data. This possibility is shortly described in the "aggregated power curve" method  
88 proposed by [Nuno Martinez et al. \(2018\)](#) for the generation of stochastic solar area power forecast scenarios.  
89 Unfortunately, insufficient implementation and model performance information is provided by the authors.  
90 To the best of our knowledge, no other research investigating the potential of hybrid regional PV model for  
91 forecasting applications exists. This presents a significant opening for further research.

### 92 *1.3. Contribution*

93 Considering the lessons and outcomes identified in the literature review, we present a clear need for  
94 further investigation on the use of hybrid models for regional PV power forecasting applications. The  
95 objective of this paper is, therefore, to propose a physical regional PV power model whose parameters can  
96 be robustly inferred from estimates of the aggregate PV power production and to provide first results on its  
97 performance for forecasting the regional PV power generation.

98 The major benefit of achieving this goal is that regional PV power can be theoretically forecast without  
99 the uncertainties from a lack of reported metadata on the installed PV systems. Thus, regional PV power  
100 forecasts would be more accurate.

101 [section 2](#) is dedicated to the description of the regional PV power model that has been used for this work.  
102 We have chosen the same formulation as previously used in [Schubert \(2012\)](#), [Saint-Drenan et al. \(2017\)](#) and  
103 [Nuno Martinez et al. \(2018\)](#). This modelling approach is linear so that unknown parameters can be obtained  
104 using regression techniques as in [Nuno Martinez et al. \(2018\)](#). We have paid a particular attention to the  
105 implementation details of this method (choice of the reference configuration, spatial constraints) in order to  
106 limit the number of unknowns without impacting the modelling accuracy.

107 [section 3](#) focuses on the estimation of the model parameters. While the unknown parameters can easily  
108 be found by a simple regression, preliminary experiments have shown that a regression yields very high  
109 or negative parameters, which are physically meaningless and very sensitive to small variations in the  
110 training data set. These first observations, that contradict physical expectations, result from the “ill-  
111 conditioned nature” of the problem, which will penalize the power estimation accuracy and ultimately the  
112 forecast accuracy in application. To address this, we instead infer the parameters with a Bayesian method,  
113 which is a standard approach in inverse modelling. An additional benefit from a Bayesian approach is the  
114 integrating of an initial parametrization state such that previously defined iterations can be exploited to  
115 improve robustness.

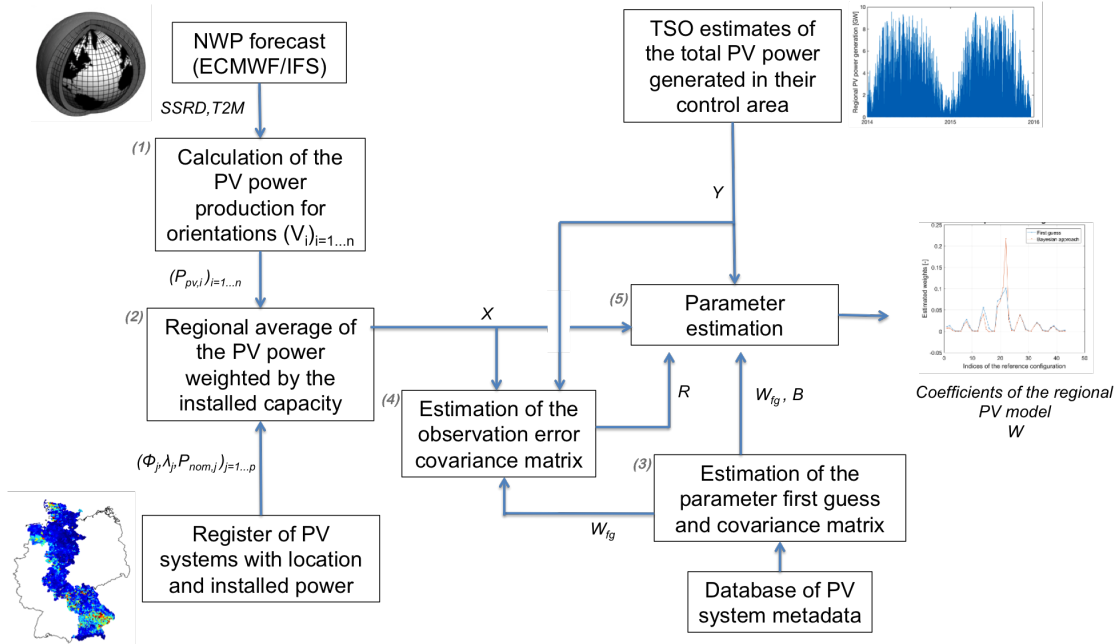
116 [section 4](#) contains the results obtained from the proposed methodology in application to forecasting. One  
117 year of regional PV power generation of three German TSOs is improved using Bayesian parameterisation  
118 before being applied to day-ahead forecasting using corresponding weather forecasts taken from the IFS  
119 numerical weather forecast model. The benefit resulting from a Bayesian approach are quantified and the  
120 performances of the obtained forecast are compared to alternative forecasting approaches.

121 Finally, a discussion on the potential of the proposed method and concluding remarks are given in  
122 [section 5](#).

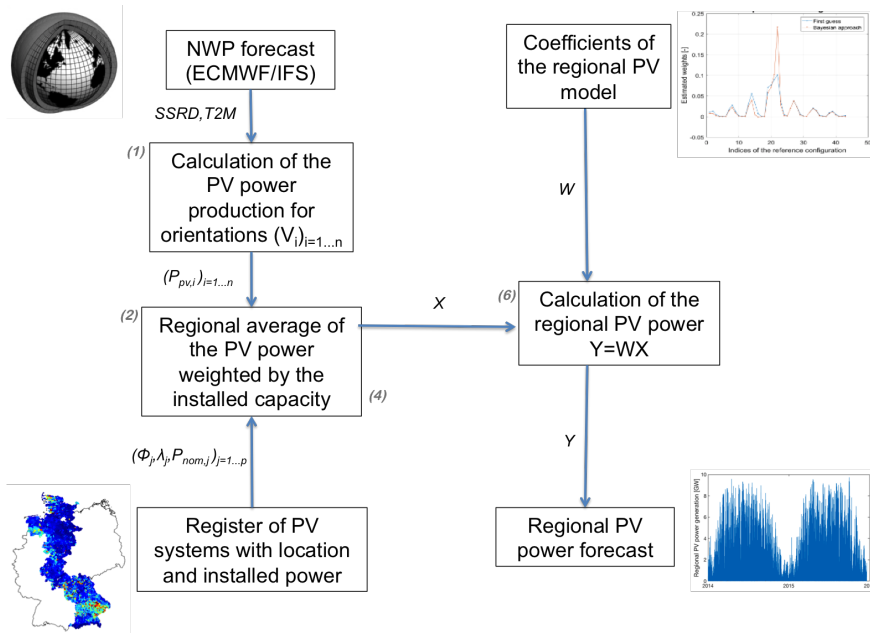
123 A flow chart summarizing our approach is given in [Figure 1](#). This figure is referred to throughout this  
124 work to guide the reader on the different calculation steps of our method. [Figure 1](#) is composed of two parts.  
125 The first part illustrates how the parameters of the regional model are inferred from the time series of the  
126 aggregated PV power production. The second part illustrates how the estimated parameters can be used to  
127 calculate regional PV power production using NWP data.

## 128 **2. Regional PV power model**

129 In this section, the regional PV power model used in this work is described in [section 2.1](#). It is a  
130 linear model that calculates regional PV power generation from meteorological data. A set of reference PV  
131 system configurations is needed; the selection process is described in [section 2.2](#). The loss of model accuracy  
132 resulting from our chosen set of configurations is then analyzed in [section 2.3](#). In [section 2.4](#), we describe  
133 the spatial constraints used to limit the number of unknowns due to the size of the regions. Finally, the



(a) Parameter estimation



(b) Calculation of the regional PV power forecast using the estimated parameters

Figure 1: Flowcharts summarizing (a) the calculation steps for the estimation of the parameters of the regional model, and, (b) calculation of the regional Pv power from NWP data using the estimated parameters.

134 modelling approach is summarized and matrix notation is introduced in [section 2.5](#).

135 *2.1. Description of the regional model*

136 The regional PV model used in this work is the same as in [Saint-Drenan et al. \(2017\)](#); it is very similar  
 137 to the approaches of [Schubert \(2012\)](#), [Fonseca Junior et al. \(2015\)](#) and [Nuno Martinez et al. \(2018\)](#). The  
 138 model is based on the idea that, instead of simulating each PV system individually, power production from  
 139 only those PV systems representing common PV system metadata configurations within a given region are  
 140 simulated. They are then scaled up to the total capacity within the region. This approach represents a  
 141 significant improvement in computational efficiency.

142 The regional power generated by a fleet of PV systems installed at location  $x$  can be expressed as the  
 143 sum of the power values calculated with characteristics  $V_i$  multiplied by the share  $w_i$  of the total capacity of  
 144 systems having the characteristics  $V_i$ . The simulated power is normalized by the corresponding peak power  
 145 and the weighted sum scaled to the actual value of the installed capacity at location  $x$ . The regional PV  
 146 power generation can then be expressed as:

$$P_{PV,t,region} = \underbrace{\sum_{x \in region} \underbrace{P_{installed,x}}_{\text{installed PV power}} \times \underbrace{\sum_{i=1}^{n_{ref}} \left( \underbrace{w_{i,x}}_{\text{weight}} \times \underbrace{f_{PV}(x,t,G_{x,t},T_{2m,x,t},V_i)}_{\text{power gen. of a ref. PV system}} \right)}_{\text{aggregated PV power at location } x \text{ (} P_{PV,t,x} \text{)}}}_{\text{aggregated PV power for the region}} \quad (1)$$

147 where:

- 148 •  $t$  is the time
- 149 •  $x$  represents the different locations within the considered region
- 150 •  $P_{PV,t,region}$  is an estimate of the aggregated power produced by all PV systems in the considered  
 151 region at time  $t$  [W]
- 152 •  $P_{installed,x}$  is the installed PV power at location  $x$  [ $W_p$ ]
- 153 •  $P_{PV,t,x}$  is an estimate of the aggregated power produced by all PV systems at location  $x$  at time  $t$   
 154 [ $W/W_p$ ]
- 155 •  $w_{i,x}$  is the weight of the  $i^{th}$  reference configuration at location  $x$  [-]
- 156 •  $f_{pv}(\dots)$  is a function representing the single PV system model used to calculate the normalized PV  
 157 power [ $W/W_p$ ]
- 158 •  $V_i$  is a vector with the configuration parameters of the  $i_{th}$  reference PV system
- 159 •  $G_{x,t}$  is the global horizontal irradiation at location  $x$  and time  $t$  [ $W/m^2$ ]
- 160 •  $T_{2m,x,t}$  is the air temperature at  $x$  and  $t$  [ $^{\circ}C$ ]

161 The function  $f_{PV}$  in Eq. 1 represents a single PV system model that needs to be chosen beforehand.  
 162 This function corresponds to the step (1) in the flow chart displayed in [Figure 1](#). [Saint-Drenan et al. \(2017\)](#)  
 163 demonstrated that a simple model with a limited number of input parameters performs well in regional



164 applications; thus, we select the same model for the present work. The calculation steps of this model are  
 165 illustrated in Figure 2 and a detailed description can be found in (Saint-Drenan, 2015).

166 With the chosen model, the set of characteristics  $V_i$  is only composed of the azimuth angle  $\alpha_{M,i}$  and  
 167 module tilt angle  $\gamma_{M,i}$ . Another important model parameter that is not explicitly considered here is the  
 168 total efficiency of the PV system. Though large variations of the efficiency may be observed among PV  
 169 systems (Killinger et al., 2018), we decided to use a constant. That said, variations of the system total  
 170 efficiency are implicitly considered in our regional model through the weights  $w_i$ . Therefore, these weights  
 171 not only reflect the distribution of capacity across all orientations, but also account for the efficiency of the  
 172 PV systems.

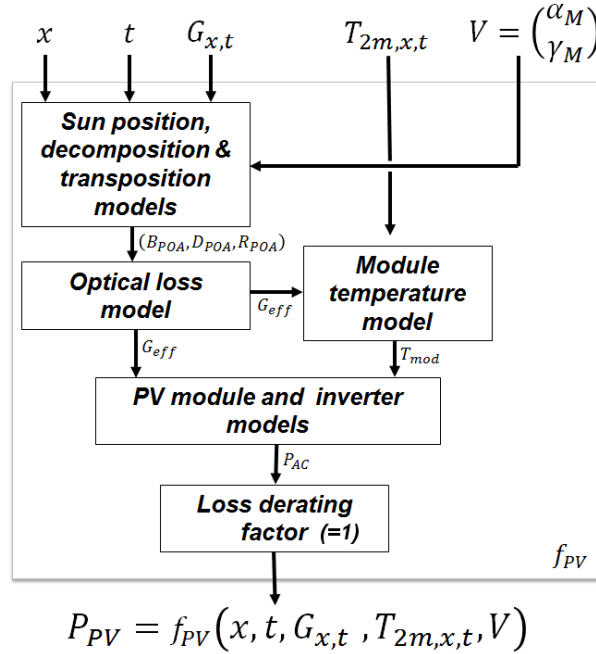


Figure 2: Flowchart of the single PV system power model.

## 173 2.2. Choice of the reference orientation combinations

174 In the single PV system model (see section 2.1), orientation of a PV system is defined by two param-  
 175 eters: azimuth and tilt angles. It is obvious that the power production simulated for two close module  
 176 orientations are highly correlated, and so we do not necessarily need every combination of tilt and azimuth,  
 177 only a representative set. Therefore, a smaller subset of orientation combinations that allows an accurate  
 178 calculation of the aggregated PV power generation is needed. We call this subset the ‘reference orientation  
 179 combinations’. Selecting too many combinations could lead to numerical issues during the search of the  
 180 model parameters due to the high number of unknowns and their co-linearity. Therefore, we use a regular

181 grid of tilt and azimuth angles with a coarse resolution of  $15^\circ$  in both dimensions to define the reference  
182 orientation combinations.

183 To further limit the number of unknown parameter of our model, we exclude uncommon combinations.  
184 Based on [Saint-Drenan \(2015\)](#) and [Killinger et al. \(2018\)](#), we limit the azimuth angles values between  $-45^\circ$   
185 and  $45^\circ$ , and tilt angles between  $0$  and  $45^\circ$ . Ultimately, we represent all systems by 22 reference orientation  
186 combinations.

### 187 *2.3. Uncertainty arising from the reference orientation configurations*

188 In the previous section, we opted for a coarse grid of reference orientation combinations. Though moti-  
189 vated by numerical considerations, the question arises of ‘how much loss of accuracy results from from this  
190 choice?’ To answer this question, we calculate the error when Eq. 1 is used to calculate the power produc-  
191 tion of a single PV system with arbitrary orientation. As errors within the regional model are expected to  
192 balance out with a increasing number of PV systems, the consideration of a single PV system should provide  
193 a worst-case indication of the error. To demonstrate, we attempted to rebuild the power production of an  
194 arbitrary orientation using our set of 22 reference orientation combinations. For this purpose, the power  
195 production was calculated using single PV power model for the 22 reference orientation combinations as  
196 well as for the single arbitrary test orientation. We used one year of 15-min global horizontal irradiation  
197 measurements and air temperature measurements from Fraunhofer IWES in Kassel. The weights of the  
198 linear model were then calculated by a multiple linear regression between the time series of the calculated  
199 power for the test orientation and the time series from the reference orientations.

200 The underlying assumption of Eq. (1) that the output of a PV system with an arbitrary orientation  
201 can be assessed by a linear combination of the outputs of different PV systems corresponding to reference  
202 orientations is illustrated in the left panel of [Figure 3](#): the squares represent each reference orientation;  
203 the circle represents the arbitrary test orientation; the numbers in boxes and the lines show the regression  
204 coefficients. In the right panel of [Figure 3](#), the power values calculated with the physical model are displayed  
205 against the linear combination of reference power values. The difference between the two power values are  
206 negligible.

207 This procedure is repeated for all integer value combinations of the azimuth angle (between  $-45^\circ$  and  $45^\circ$ )  
208 and the tilt angle (between  $0$  and  $45^\circ$ ). The resulting root mean square difference (RMSD) of the residual  
209 of the regression are represented by colored squares in [Figure 4](#).

210 All RMSD values are less than  $2.10^{-4} W/W_p$ ; this analysis shows that the power values calculated  
211 with the 22 reference orientation combinations is coarse enough to reconstruct the power for an arbitrary  
212 orientation (as long as this orientation lies in the domain covered by the reference orientations). We observe  
213 peaks in RMSD (though still small) at around  $40^\circ$  tilt and between  $45^\circ$  to  $30^\circ$  east and west. Since such  
214 orientations are infrequent in Germany, we consider these greater RMSD values to be of insignificant impact.

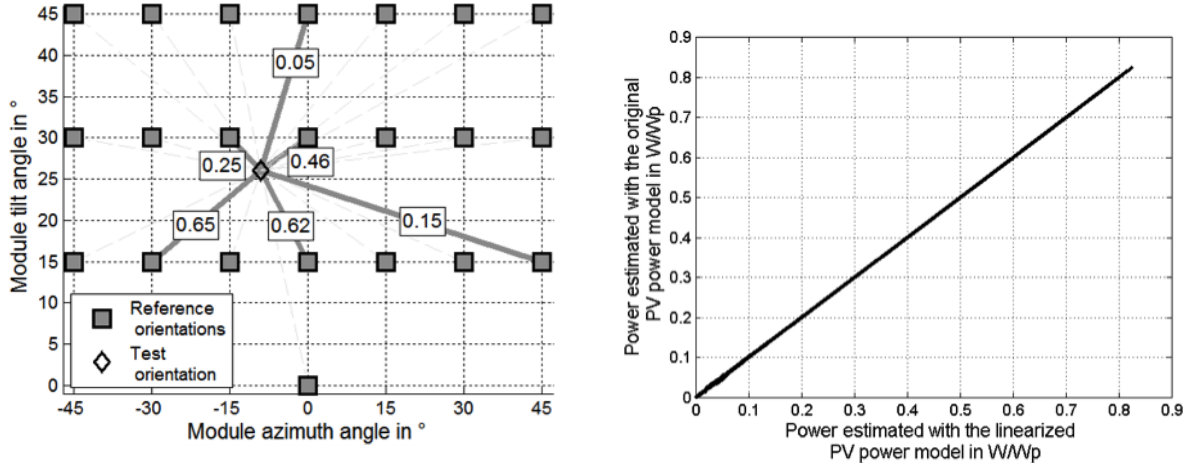


Figure 3: Left panel) Illustration of the approach used to estimate the PV power at a given orientation as derived from the PV power value estimated with the reference orientations using a linear approach. The reference and the test orientations are represented by grey squares and a circle, respectively. Weighted coefficients are indicated in each box. Small coefficients ( $\leq 0.01$ ) are not displayed. Right panel) Scatter plot of power values evaluated with the physical model against power values evaluated with the linear approach.

Control area or region	Number of PV systems	Size of the region
TenneT	646 968	140 500 $km^2$
50Hertz	125 696	109 000 $km^2$
Amprion	437 622	73 100 $km^2$
TransNetBW	281 420	34 600 $km^2$
Germany	1 491 706	357 168 $km^2$

Table 1: Number of installed PV systems in the four German control areas in August 2014.

#### 215 2.4. Spatial constraints

216 At this stage, there is a total of  $n_{ref} \times n_{loc}$  unknowns in our regional model, where  $n_{ref}$  is the number of  
 217 reference orientation combinations (22 orientations selected in the previous section) and  $n_{loc}$  is the number  
 218 of locations considered in the regions. Considering the size of a control area and the great number of PV  
 219 systems installed in a control area (Table 1), it is necessary to add spatial constraints to the problem in  
 220 order to limit the number of unknowns.

221 Given that the unknowns are inferred from regionally averaged PV power generation, it is unrealistic  
 222 to expect that local information can be retrieved from the method. In addition, the possible impact from  
 223 errors at the local scale are balanced when upscaled to a regional level; thus, they are deemed insignificant.

224 Therefore, we are mainly interested in large scale spatial trends that may have a significant impact on the  
 225 aggregated regional PV power generation estimate. It is possible to add spatial constraints by considering

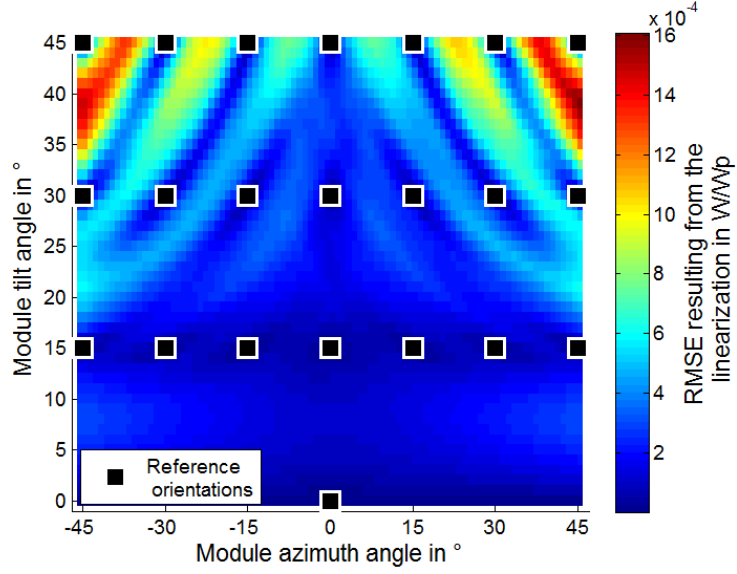


Figure 4: Root mean square difference (RMSD) between the PV model output and the corresponding value obtained with the 22 reference orientation combinations using a multilinear regression.

226 large scale regional trends; approaches exist such as spatial regularization techniques or the use of spatial  
 227 parametric model. We selected the simplest and most pragmatic method—dividing a region into a reasonable  
 228 number of sub-regions and assuming that unknown parameters are constant in each sub-region.

229 Regions are divided into sub-regions using a k-mean clustering algorithm on the coordinates of PV  
 230 systems installed in the greater region. Two examples are illustrated in Figure 5, where the control area of  
 231 TenneT is divided into 3 and 5 sub-regions, respectively. In this figure, light grey dots represent the known  
 232 installed PV systems and dark squares are centroids of the clusters (sub-regions). The borders between  
 233 different sub-regions are displayed by blue lines.

234 If  $n_{SbRg}$  is the number of sub-regions, the number of unknown of the linear system is now equal to  
 235  $n_{ref} \times n_{SbRg}$ , for example  $n_{ref} = 22$  and  $n_{SbRg} = 2...5$  tractable.

### 236 2.5. Summary

237 The weights  $w_{i,x}$  are assumed constant within each sub-region  $R_j$ . Thus, the model formulation in Eq. 1  
 238 can be simplified by introducing the smaller set of  $w_{i,R_j}$  for each configuration  $i$  and region  $R_j$ . The weights  
 239  $w_{i,x}$  are related to the weights  $w_{i,R_j}$  by the relationship  $w_{i,x} = w_{i,R_j} \forall x \in R_j$ . With this new variable, Eq. 1  
 240 is written:

$$\underbrace{P_{PV,t,region}}_{Y[t,1]} = \underbrace{\sum_{R_j} \sum_i}_{\sum_k} \underbrace{(w_{i,R_j})}_{W[k,1]} \times \underbrace{\left( \sum_{x \in R_j} P_{installed,x} \times f_{PV}(x,t, G_{x,t}, T_{2m,x,t}, V_i) \right)}_{H[t,k]} \quad (2)$$

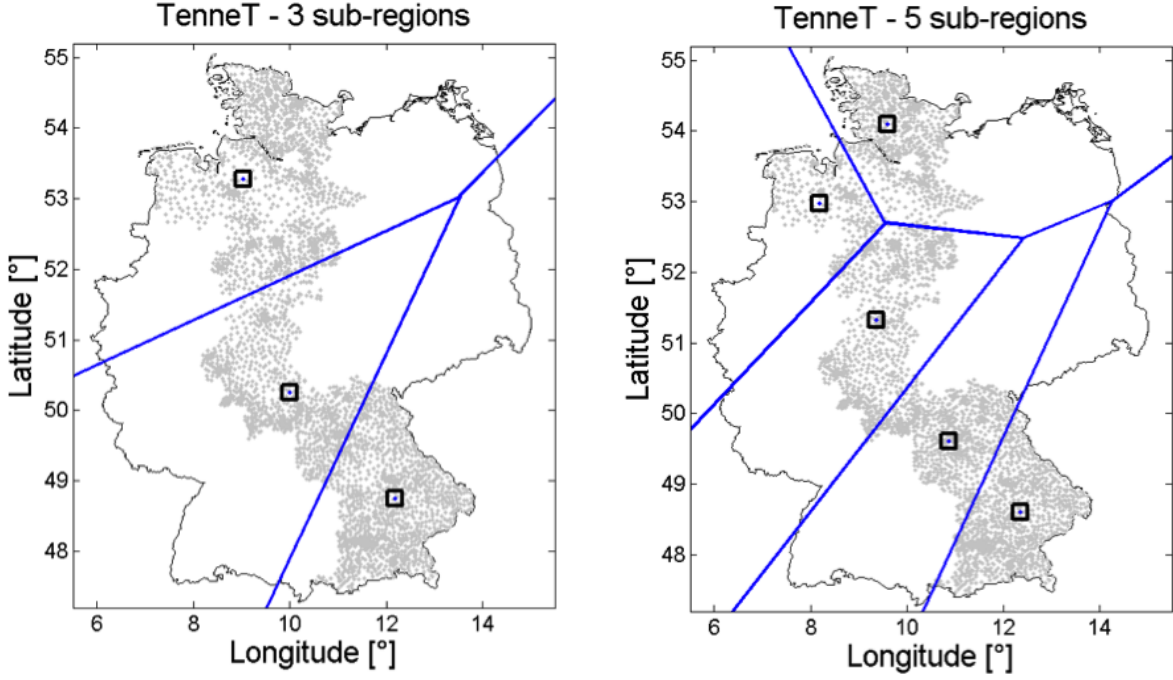


Figure 5: Illustration of the split of the TenneT control region into 3 sub-regions (left map) and 5 sub-regions (right map) using k-mean clustering. The centroids of each sub-regions are marked by a black square and the limits are indicated with blue lines.

241 For the implementation of Eq. 2, PV system installed capacity  $P_{installed,x}$  and meteorological information  
 242 ( $G_{x,t}, T_{2m,x,t}$ ) are needed for each location  $x$  of the considered region. With this information, the right-side  
 243 of Eq. 2 ( $H[t, k]$ ) can be calculated for each sub-region  $R_j$ . This quantity corresponds to the weighted sum of  
 244 the simulated power values for a configuration  $V_i$ , the weights being the installed capacities. The unknown  
 245 of the problem are the weights  $w_{i,R_j}$  for each reference configuration  $V_i$  and sub-region  $R_j$ .

246 It is convenient to express Eq. 2 in matrix notation  $\mathbf{Y} = \mathbf{H}\mathbf{W}$ . The column vector  $\mathbf{Y}$  contains aggregated  
 247 PV power generation in a given region at different time steps  $t$ . The vector  $\mathbf{W}$  contains the unknowns of  
 248 the problem  $w_{i,R_j}$ , and  $\mathbf{H}$  is a matrix containing the sum of the simulated PV power values for each  
 249 combination of sub-region and reference orientations (column) and time steps (row). For the sake of clarity,  
 250 the relationship between the summation and matrix notations are indicated by underbraces in Eq. 2. In  
 251 this illustration,  $k$  is an index that screens all combinations of reference orientations  $i$  and sub-regions  $j$ .

252 The above described summation of the simulated PV power weighted by the installed capacity to evaluate  
 253 the matrix  $\mathbf{X}$  corresponds to the step (2) of the flow chart given in Figure 1.

### 254 3. Estimation of the model parameters

#### 255 3.1. Approach

256 The problem described in [section 2.5](#) is inverse: we start with results (aggregate PV power generated  
257 within the region) then calculate the cause (weights for all PV systems with the reference orientation  
258 combinations). This is the opposite to previous methodologies ([Nuno Martinez et al., 2018](#); [Saint-Drenan  
259 et al., 2017](#)) where the authors start with the cause (the distribution of the PV systems according to the  
260 region and reference orientation combinations) and then calculate the result. Inverse problems feature  
261 heavily in literature ([Nakamura and Potthast, 2015](#); [Colton and Kress, 1997](#); [Engl et al., 2000](#)). They are  
262 common in many research fields, especially in meteorology: the retrieval of cloud optical characteristic, the  
263 assimilation of measurements in numerical weather prediction models, to name a few.

264 Solving inverse problems is non-trivial as they are typically ill-posed. Among the three Hadamard  
265 conditions for a well-posed problem ([Hadamard, 1902](#))—existence, uniqueness, and stability of the solution—  
266 the conditions of uniqueness and stability are often violated. Particularly with a great number of unknown  
267 parameters, an observation can be explained by several causes (violation of the uniqueness condition). In  
268 our application, multicollinearity of the input data poses an issue (see [section 2.3](#)). The linear dependency  
269 among regressors is reflected in the solution space, where many possible combinations of the solution vector  
270 components lie within a narrow valley (in high dimension) all positioned very close to the minimal residual  
271 sum of squares. In such cases, changes in the model output can instead be obtained by changes in the model  
272 input parameters.

273 Small perturbations in the input data can bring about noticeable changes in the solution (violation of  
274 the stability condition). Such problems can be addressed by using regularization techniques; for example,  
275 the generalized Tikhonov regularization, whereby the deviation of the solution from the initial guess is  
276 penalized ([Tikhonov, 1963](#)). This approach requires the choice of regularization parameters that balance  
277 the respective effects of the error and regularization terms ([Nakamura and Potthast \(2015\)](#), Chapter 3.1.6).  
278 There are many techniques for motivating the choice of the regularization parameter; however, selection  
279 remains a non-trivial and delicate issue. Alternatively, regularization parameter selection can be entirely  
280 avoided with a Bayesian approach.

281 In the Bayesian framework (see for example [Nakamura and Potthast \(2015\)](#) chapters 4.2 and 5.6, [Crisan  
282 and Bain \(2009\)](#) or [Crisan et al. \(2014\)](#)) our goal is to find the set of parameters  $\mathbf{W}$  giving the highest  
283 posterior probability  $P(\mathbf{W} | \mathbf{Y})$  given  $\mathbf{W}$ . To this end, we use the well known Bayes law:

$$\underbrace{P(\mathbf{W} | \mathbf{Y})}_{\text{posterior}} \propto \underbrace{P(\mathbf{W})}_{\text{prior}} \underbrace{P(\mathbf{Y} | \mathbf{W})}_{\text{likelihood}} \quad (3)$$

284 We assume that the prior  $P(\mathbf{W})$  can be approximated by a Gaussian distribution with mean value  $\mathbf{W}_{fg}$   
285 and covariance matrix  $\mathbf{B}$ . Similarly for the likelihood  $P(\mathbf{Y} | \mathbf{W})$ , we take a Gaussian distribution with

286 zero mean and covariance matrix  $\mathbf{R}$ . The covariance matrix  $\mathbf{B}$  can be interpreted as a quantification of  
 287 the possible variations of the resultant vector around the first guess; the second covariance matrix  $\mathbf{R}$ , a  
 288 quantification of the uncertainty of the observations  $\mathbf{Y}$ .

289 With these notations, the posterior probability can be written:

$$P(W | Y) \propto \exp\left(\frac{1}{2}(W - W_{fg})^T B^{-1}(W - W_{fg})\right) \times \exp\left(\frac{1}{2}(HW - Y)^T R^{-1}(HW - Y)\right). \quad (4)$$

290 Given that maximizing the likelihood is equivalent to minimizing the logarithm of the above expression  
 291 (Freitag and Potthast, 2013), the desired solution  $W_{opt}$  is to minimize the following functional:

$$J(W) = \frac{1}{2}(W - W_{fg})^T B^{-1}(W - W_{fg}) + \frac{1}{2}(HW - Y)^T R^{-1}(HW - Y). \quad (5)$$

292 The solution that minimizes the above cost function  $J$  has zero gradient ( $\nabla_W J = 0$ ). This condition  
 293 allows the explicit determination of the solution to Eq. 5:

$$W_{opt} = W_{fg} + [B^{-1} + H^T R^{-1} H]^{-1} \times (H^T R^{-1})(Y - HW_{fg}) \quad (6)$$

294 This relationship, which corresponds to the expression of the generalized Tikhonov regularization, can  
 295 now be used for estimating the weights that correspond to the different reference orientation combinations.  
 296 To this end, the estimation of the initial guess and covariance matrices ( $\mathbf{R}$  and  $\mathbf{B}$ ) is achievable—detailed  
 297 in the following subsections. It is illustrated by the step (5) in Figure 1.

### 298 3.2. Determination of the initial guess

299 As mentioned previously mentioned, the unknown parameters  $\mathbf{W}$  can be interpreted as the distribution  
 300 of PV systems according to the different possible orientations and to each sub-region  $R_j$ . This interpretation  
 301 can be exploited to evaluate the initial guess  $W_{fg}$ .

302 The same approach to evaluate our initial guess was presented by Saint-Drenan et al. (2017). A database  
 303 including module orientation angles for more than 20,000 systems is used to evaluate the share of the total  
 304 capacity corresponding to each of the reference orientation combinations. When evaluating the initial guess,  
 305 we neglected the potential geographical dependence of the parameters of our regional model. The distribu-  
 306 tions were thus evaluated using all PV systems in the database regardless of their location in Germany. If  
 307 sub-regions are later considered, we assume that the same distribution can be used as an initial guess for  
 308 each sub-region. Regions can and do present distinct differences in orientation; when considering a larger  
 309 aggregate statistic, these nuances can be ignored—Killinger et al. (2018) (fig. 3) visually demonstrates a  
 310 distinct north-south division in PV system orientation for France. In Figure 6, the components of the first  
 311 guess vector are represented by squares that are coloured as a function of the module azimuth and tilt  
 312 angles. The statistical analysis aiming at the estimation of the first guess is represented by the step (3) in  
 313 Figure 1.

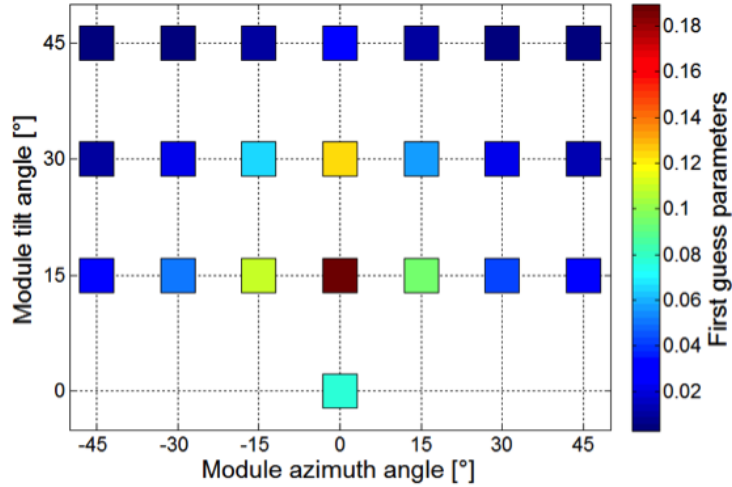


Figure 6: Values used for the initial guess (colour of the squares) as a function of the module azimuth and tilt angles.

### 3.3. Estimation of the background covariance matrix

The background covariance matrix  $\mathbf{B}$  quantifies the expected dispersion of the parameter vector around the initial guess  $W_{fg}$ . In order to evaluate this distribution, we generate a set of perturbed first guess vectors using the approach described in section 3.2:

$$\left\{ \left[ \tilde{w}_1^{(p)} \dots \tilde{w}_j^{(p)} \dots \tilde{w}_{n_{ref}}^{(p)} \right]^T, p = 1 \dots 10000 \right\} \quad (7)$$

The perturbation uses only 1000 randomly chosen PV systems instead of the total dataset. Indeed, a random sub-sampling of the set of PV systems brings about variations of the original distribution of PV systems according to the reference orientation combinations. However, we assume that with a sufficiently small number of samples and a sufficiently large number of iterations, the range of possible values taken by the resulting distribution is accounted for. This approach is re-iterated 10,000 times to populate a set of possible solutions that are subsequently used for the estimation of the background matrix  $\mathbf{B}$ .

The set of randomly evaluated parameter sets is then used to calculate each component  $b_{i,j}$  of the matrix  $\mathbf{B}$  using the following relationship:

$$b_{i,j} = Cov \left( \left[ \tilde{w}_1^{(i)} \dots \tilde{w}_k^{(i)} \dots \tilde{w}_{n_{ref}}^{(i)} \right]^T, \left[ \tilde{w}_1^{(j)} \dots \tilde{w}_l^{(j)} \dots \tilde{w}_{n_{ref}}^{(j)} \right]^T \right). \quad (8)$$

All PV systems are considered in the covariance estimation so that the resulting matrix reflects the true covariance of the model parameters for Germany (note that, when applying this methodology to a different country/region, a representative covariance matrix is required). We assume that the same covariance holds true for each sub-region and that the covariance of two parameters in different regions is null.

The estimation of the background covariance matrix using a dataset of PV system metadata corresponds to the step (3) in Figure 1.



332 *3.4. Estimation of the observation error covariance matrix*

333 The observation error covariance matrix quantifies the uncertainty of the observation vector  $\mathbf{R}$ . We  
334 assumes that this error is time-independent so that  $\mathbf{R}$  can be expressed as the product of an identity matrix  
335 with a scalar. We further assume that the variance of the error obtained with the initial guess  $\mathbf{W}_{fg}$  evaluated  
336 over a training time period can be used as an approximation of that of the error variance. As a result,  $\mathbf{R}$   
337 can be evaluated as follows:

$$\mathbf{R} = Var(\mathbf{H}\mathbf{W}_{fg} - \mathbf{Y}) \cdot \mathbb{I} \quad (9)$$

338 The estimation of the  $\mathbf{R}$  is subjected to several assumptions that can have a non negligible impact on the  
339 results. By choosing this simple and rough approach to gain insight on the benefit of the Bayesian approach  
340 applied to our problem, future work is needed to improve the estimation of this matrix. It is represented by  
341 the step (4) in [Figure 1](#).

342 **4. Results**

343 This section is separated into four parts. The data are firstly introduced in [section 4.1](#). The training  
344 period and validation procedure is described in [section 4.2](#). Model performance looking particularly impacts  
345 from regularization and changing number of sub-regions is presented in [section 4.3](#). Comparisons of the  
346 newly proposed forecasting approaches are compared to persistence, initial guess, and the TSO forecast in  
347 [section 4.4](#).

348 *4.1. Data*

349 The first type of information needed for implementing our method is an estimation of the aggregated PV  
350 power produced in the considered area. In this work, we use TSO estimates of the aggregated PV power  
351 generation for three German control areas: TenneT, Amprion and 50Hertz. For each control area, time series  
352 of the estimated aggregated PV power generation as well as the installed PV system metadata (location,  
353 peak power and time of installation) are available. The time series of the regional PV power estimates have  
354 a temporal resolution of 15 min and are available for the years 2014 and 2015. A short description of these  
355 data provided by the German TSOs can be found in [Saint-Drenan et al. \(2017\)](#).

356 The second input required by our model is meteorological data and more precisely the global horizontal  
357 irradiation and air temperature. The most natural choice is to select the most accurate gridded data, as for  
358 example reanalysis and/or satellite data. As reported by [Ineichen \(2014\)](#) or [Boilley and Wald \(2015\)](#), these  
359 data have their own sources of error that can add to forecast error of the NWP data used for the predic-  
360 tion. We therefore chose the pragmatic option of using directly the NWP forecast data for the parameter  
361 estimation in a way to avoid double penalty error described above. Forecasts from the Integrated Forecast

362 System (IFS) model run by the European Center for Medium-Range Weather Forecast (ECMWF) are there-  
363 fore used for the calculation of the PV power generation for each region and for each reference orientation  
364 (matrix  $\mathbf{H}$ ). Hourly gridded 2m ambient temperature (2T) and solar surface radiation downwards (SSRD)  
365 are available at  $0.125 \times 0.125^\circ$  spatial resolution; they are collected, prepared and linearly interpolated to a  
366 15 min temporal resolution corresponding to the PV power. Only forecast data with a lead time between  
367 24 and 48h for the run starting at 00:00 (GMT) are used to evaluate day-ahead forecast.

368 Finally, we used the operational forecast published by the three TSOs on their website to benchmark  
369 the output of our model. These forecast are calculated by the TSOs using several forecasts of the regional  
370 PV power generation from private forecast providers. These individual forecasts are optimally mixed and  
371 calibrated using the reference PV production value.

372 In this work, we decided to show the potential of our approach using real-world operational data: the  
373 parameters of the regional model are estimated using TSO estimates of the PV power production and the  
374 model output is compared to TSO day-ahead forecasts. If this approach allows demonstrating the practical  
375 relevancy of our work, it has also some drawbacks: uncertainties in the TSO data <sup>2</sup> used for training  
376 and validation can bring about errors, that are impossible to differ from the actual error of our method.  
377 An alternative to avoid these undesired effects would have consisted in generating synthetic data. Yet, we  
378 preferred evaluating the performances of our method in real conditions because the motivation of the present  
379 work is strongly linked to the targeted application. As a consequence, it is important to bear the real-world  
380 settings of the validation in mind to properly interpret the results presented in this chapter.

#### 381 *4.2. Training and validation setup*

382 An adaptive approach has been chosen to train the model (see Figure 7). The model parameters are  
383 evaluated and tested for the year 2015 on a monthly basis using 12 months of training data preceding  
384 each test time-period. This restrains various issues that could lead to differences in characteristics between  
385 training and testing data, such as regular improvement of NWP models, change of the TSO estimates (e.g.  
386 extension of the set of reference systems with time), or, new installed PV systems over time.

387 The presence of snow on PV systems results in negligible PV power generation that would not be  
388 appropriately represented in our proposed system without additional complexity. Since this effect is not  
389 considered, days marked by the presence of snow in Germany must be excluded from the training and testing  
390 data (see Figure 7). To do this, days with snow have been identified using snow height data at 200 German  
391 SYNOP stations collected from the OGIMET website ([www.ogimet.com](http://www.ogimet.com)). The network of SYNOP station  
392 is coarse; we decided to exclude each day when snow is reported at more than one meteorological station as

---

<sup>2</sup>The major sources of uncertainty are the uncertainty of the regional estimates and the information on the installed PV capacity.



Figure 7: Graphical representation of the training and testing setup (upper plot) and number of days available for training and testing (lower plot).

393 a conservative measure. Should large NWP forecast errors be present at the training phase of the model, the  
 394 whole approach will be unjustly penalized. To avoid this, all days were excluded from the training dataset  
 395 only if at least one value of the first guess forecast error exceed  $0.2W/W_p$ . This latter criterion is not applied  
 396 to the testing data set.

397 Finally, training and testing routines were conducted at three control areas with a varying number of  
 398 sub-regions (1 to 5).

399 For each forecast considered in this work, all standard error metrics were evaluated using only daytime  
 400 values; these results are presented in Table 2. In the two next sections, the performances of the forecast  
 401 methods will be discussed with a focus on the RMSE values. Indeed, this metric is widely used by TSOs and  
 402 forecast providers to assess the forecast accuracy as it well reflects the expectations on the performances of  
 403 the forecasting methods for grid integration mechanism.

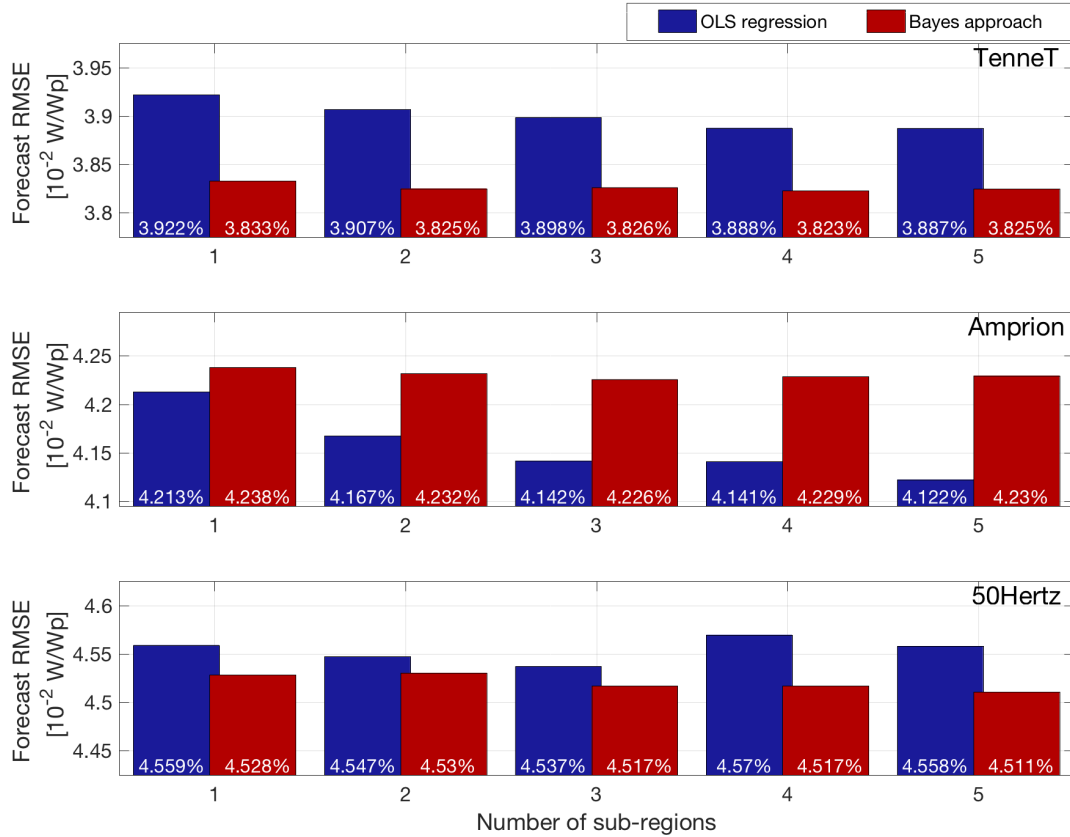


Figure 8: Effect of the number of sub-regions on the model performances for the three considered control areas (the limits of the y-axis have been scaled to illustrate best the differences between models).

#### 4.3. Model performance impact from regularization and varying number of sub-regions

In Figure 8, RMSE values obtained with an ordinary least square regression (OLS) and the Bayesian method are displayed for 1 to 5 sub-regions. Considering Figure 8, it is remarkable that the two approaches yield very different results based on the data of the three transmission system operators. Disregarding the influence of regions and focusing only on the average differences between the OLS regression and the Bayesian method, we confirm that the added value of the regularization is not systematic. A beneficial and a moderate improvement of regularization can be observed for TenneT and 50 Hertz, respectively, as indicated by smaller RMSE values from the Bayes method; regularization brings about an increase of the RMSE values for Amprion. There are likely many causes for the observed differences; however, the most probable explanation is that the initial guess—and to a lower extent the covariance matrix  $B$ —computed for the whole Germany is not representative for the three control areas. The initial guess and the covariance matrix approach was motivated by the limited number of information about the PV systems installed in Germany. A calculation with a initial guess and covariance matrix  $B$  specific to each region would be

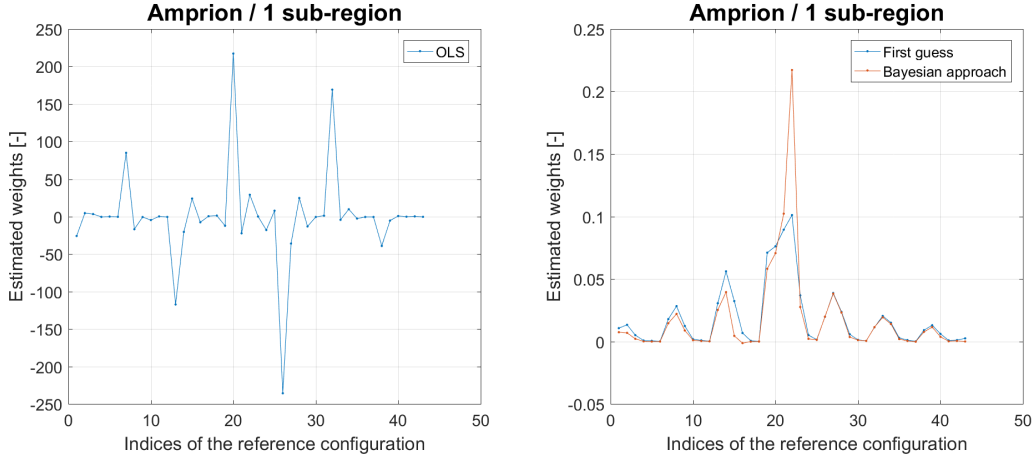


Figure 9: Comparison of the weights evaluated with the ordinary least-square regression (OLS) and those obtained with the Bayesian method with one sub-region for the Amprion control area.

417 of interest at a later data as soon as more information is available on the installed PV systems in the  
 418 respective regions. Even though RMSE from the Bayes method is worse than that of the OLS regression in  
 419 Amprion, the resulting parameters are still more robust due to the nature of a Bayes methodology. This is  
 420 observed in Figure 9 where the parameters obtained from the OLS regression are compared to the regularized  
 421 counterpart. While the Bayes method parameters closely resemble the actual values, those obtained through  
 422 OLS regression can be very high and even negative—physically impossible. Under these circumstances, if  
 423 testing data were to exist outside the parameterized domain from the training data, it must be accepted  
 424 that a forecast trained with OLS regression may yield meaningless results due to the collinearity of the input  
 425 features —this will not happen with the Bayesian method. Physically meaningful parameters are preferable  
 426 —even at the cost of a lower performance— in operational context, where a reliable forecast system is needed.

427 The evolution of the RMSE with increasing number of sub-regions is also very different across the three  
 428 main regions. While increasing the number of sub-regions results in a small but steady RMSE decrease with  
 429 the Bayes method for the three control areas, the effect of the number of regions on the performances of the  
 430 OLS method is very different from one control area to another. In the TenneT and Amprion control areas,  
 431 it is observed that the RMSE decrease for OLS regression is relatively more consistent and pronounced. By  
 432 contrast in the 50 Hertz control area, RMSE from OLS regression decrease from 1 to 3 sub-regions, but  
 433 increase again thereafter; this is presumably the result of a lack of generalization from too many degrees  
 434 of freedom with more sub-regions (overtraining). With too many parameters, the model has a tendency to  
 435 learn information that corresponds too closely to the training data set (such as rare events or even noise);  
 436 therefore, it fails to generalize the trained dependencies to testing data. If inferring model parameters for  
 437 1 to 3 sub-regions can appear less considering the numerous PV systems installed in the different regions,  
 438 it should be pinpointed that the spatial differentiation is made using only one time series of the aggregated

439 PV power production. It can therefore be expected that modifying the approach to integrate more data  
440 (DSO aggregated power, time series of single PV systems, dataset of PV metadata...) would enable a better  
441 spatial characterization of the unknown fleet of PV systems.

#### 442 4.4. Benchmarking model performance against alternative forecasts

443 In this section, the forecasting performances from the Bayesian method is compared to alternative fore-  
444 casts methods in order to evaluate any added value. We evaluated the performance of 5 different forecasting  
445 approaches at each of the three control areas. RMSE results are displayed in [Figure 10](#).

446 The first comparison is against smart persistence, whereby the current day's TSO estimate is used  
447 as the forecast for the following day. The second comparison is from the initial guess (without training  
448 phase), which corresponds to the approach described in [Saint-Drenan et al. \(2017\)](#). The performances of  
449 the Bayesian method and OLS regression with 1 to 5 sub-regions are also considered in this evaluation  
450 for additional insight. The final comparison is against the publicly available day-ahead forecasts from the  
451 German TSOs.

452 It is generally considered good practice to include persistence forecast in a benchmark as it shows  
453 the added value of more advanced method ([Yang et al., 2017](#)). Persistence can yield remarkable results  
454 with stable weather conditions over several days, however leads to significant errors under variability. For  
455 day-ahead prediction, it is well established that NWP-based forecasts outperform persistence forecasting—  
456 [Figure 10](#) corroborates this.

457 The exceptional performance of the initial guess, Bayesian method and OLS regression forecasting at  
458 Amprion is not representative and is now a subject of discussion with the German TSO; clearly the published  
459 forecast not the best forecasting approach, even though it has been made public by Amprion. Similar  
460 observations were already made by [Saint-Drenan et al. \(2017\)](#) in a previous work. We believe it would be  
461 a poor representation of the described approach to consider Amprion in our analysis as there are evidently  
462 issues with the published Amprion data; as such, Amprion data are not considered in the present discussion.

463 Comparing the Bayesian method performance against the initial guess quantifies the benefits of train-  
464 ing the parameters on historical data. We clearly demonstrate that including the training stage reduces  
465 the RMSE by 0.08 and  $0.12 \cdot 10^{-2} W/W_p$  for TenneT and 50 Hertz, respectively, representing relative im-  
466 provements of 2.05 and 2.59%—a small but clear improvement. Ultimately, training the model parameters  
467 offers real added value to the forecasting skill. We also note that the initial guess is already a reasonable  
468 approximation of the true model parameters.

469 Comparing to the RMSE from the initial guess, the Bayesian method and OLS regression corroborates  
470 our initial assumption that the parameters of the regional model can be inferred from the aggregated PV  
471 power generation. For TenneT, OLS regression RMSE is greater the Bayesian method and initial guess  
472 RMSE results. This indicates that Bayesian regularization avoids issues from overtraining. For 50 Hertz,

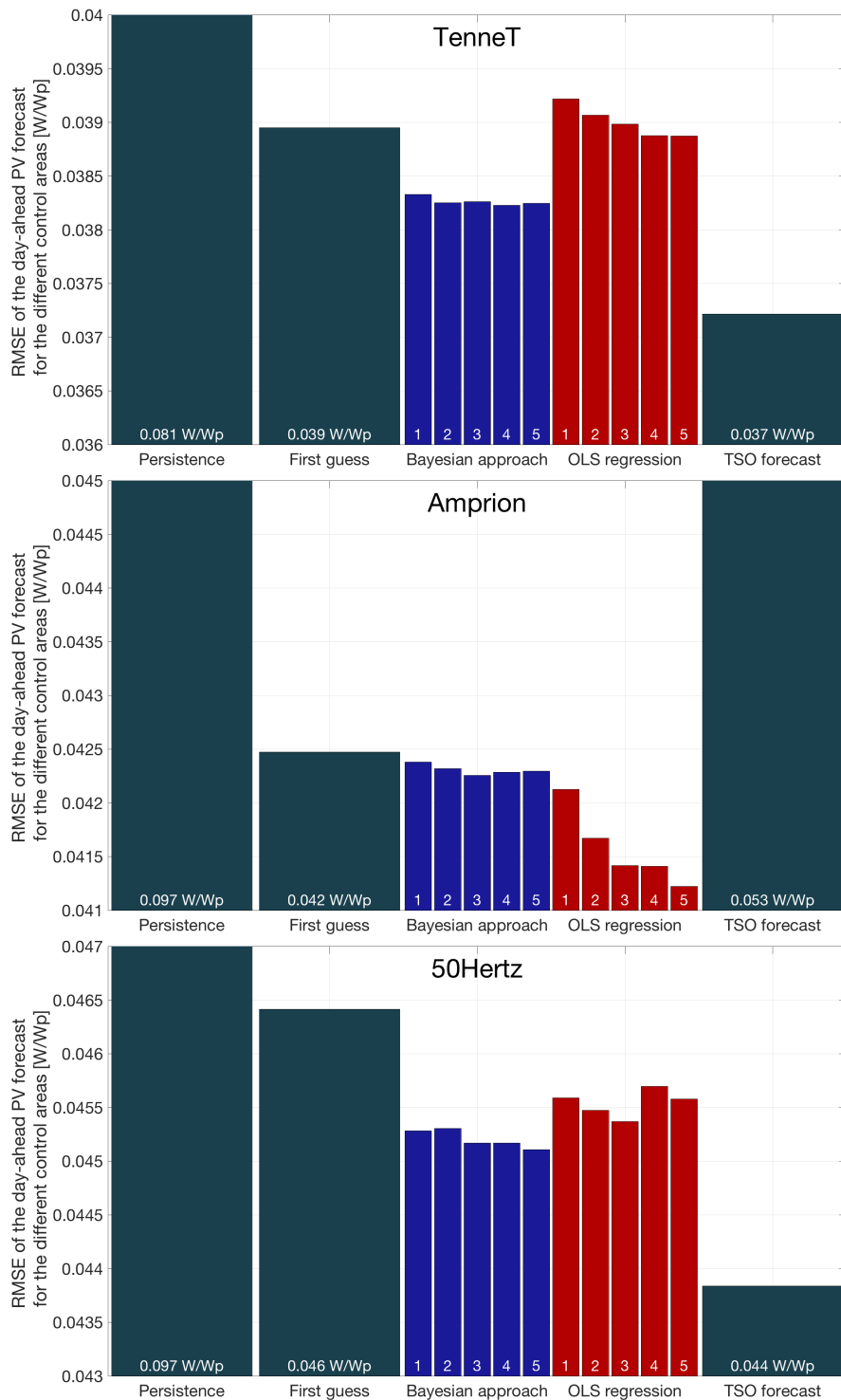


Figure 10: Comparison of the RMSE obtained with the different forecasts for TenneT (upper plot), Amprion (middle plot) and 50Hertz (lower plot)

473 OLS regression RMSE is better than initial guess RMSE. Thus, it is possible than the initial guess is sub-  
474 optimal but is balanced by the Bayesian inference, which ultimately yields a smaller RMSE than simple  
475 regression.

476 Our Bayesian method model output is finally compared to the TSO day-ahead forecast. TSO forecast  
477 RMSE values (excl. Amprion) are noticeably smaller than the alternative forecasts. TSO forecasts are the  
478 output of an optimized ensemble from several weather models; the ensemble is optimized daily by means  
479 of an adaptive method accounting for the actual value of the power production. In contrast, our simpler  
480 forecast only uses a single weather model that is optimized on a 12 month period. Given the substantial  
481 differences in required input data and adjustment techniques between the TSO and Bayesian forecasts, it  
482 is unsurprising that TSO forecasts have a better skill. Outperforming the TSO forecast is an unrealistic  
483 expectation with the proposed methodology. Instead, TSO forecasts are indicative expected performance  
484 from highly optimized forecasting approaches; this makes them a useful benchmark for the present evaluation.  
485 [Figure 10](#) and [Table 2](#) show that the difference between RMSE values obtained with the Bayesian method  
486 and TSO data decrease from 0.18 % (3.90 % - 3.72 %) with the initial guess to 0.10 % (3.82 % - 3.72 %) with  
487 the Bayesian method for Tennesse; and from 0.26% (4.64 % - 4.38 %) to 0.14% (4.52 % - 4.38 %) at 50 Hertz.  
488 This result clearly demonstrates the potential of the proposed methodology with respect to the previous  
489 versions ([Saint-Drenan et al., 2017](#)). As the Bayesian method has performance close to TSO forecasting,  
490 we expect significant rewards from further methodological development. Potential improvement avenues  
491 are—but not limited to—consideration of an ensemble of NWP models, and incorporation of advanced MOS  
492 technique, such as exponential smoothing, Kalman Filter to track change of the fleet of PV systems, and  
493 mitigate NWP forecast error.

## 494 5. Discussion and conclusions

495 In this paper, a Bayesian method to parameterize a regional PV power forecasting model is proposed  
496 using only time series of regional PV power generation. This work addresses a need of forecast providers who  
497 must compete with transmission system operators (TSOs) estimates of the regional PV power generation  
498 without having access to critical information about the PV installations within the region.

499 The proposed approach is based on a linear regional PV model previously defined in literature (e.g.  
500 [Saint-Drenan et al. \(2017\)](#)). We made considerable effort to minimize the number of unknowns, especially  
501 when calibrating the reference orientation combinations and the selecting spatial constraints. Restricting  
502 of the number of predictors—or regularization by discretization—is interpreted as a projection method.  
503 We proposed a regularization that is based on a Bayesian problem formulation and describe a method to  
504 approximate the initial parameters as well as the covariance matrices required.

505 There is the possibility to consider different number of sub-regions in our model, whereby a larger region



506 is considered as the sum of all sub-regions. We found out that the sensitivity of the forecast performance on  
507 the number of sub-regions is very different for the 3 TSOs considered; while better performance was obtained  
508 from using more sub-regions for Amprion and TenneT, it is moderate to negative for 50 Hertz. This finding  
509 is probably a result of the actual PV system characteristics in the wider region as well as from the quality  
510 of the first guess and covariance matrices. Finally, it is clear that the model is prone to overtraining with an  
511 increasing number of degrees of freedom. Overtraining issues can be efficiently addressed by the proposed  
512 regularization technique. We illustrated this by exposing the coefficients as obtained by the OLS regression;  
513 they were frequently unrealistically high or even negative values. In parallel, coefficients from the Bayesian  
514 method were reasonable, positive and physically likely, resulting in a more robust parameterization.

515 Data from three German TSOs was used to evaluate forecasting performance. We find that, as long as the  
516 initial guess and covariance matrices are well-defined, the proposed approach offers significant added value  
517 with respect to persistence, an OLS regression, or to using only the initial guess. Interestingly, we identified  
518 that the Amprion TSO forecast was operating at significantly sub-optimal performance; for representation  
519 of our methodology, we excluded Amprion from our overall conclusions. A comparison of forecasts within  
520 TenneT and 50 Hertz regions showed that all approaches were better than persistence. The initial guess from  
521 our method was better than OLS regression with only  $< 4$  sub-regions; however, OLS regression was better  
522 than the initial guess forecast at 50 Hertz. The Bayesian method was better than OLS regression and the  
523 initial guess. This is indicative of added value from Bayesian parametrization and from a parametrization  
524 approach.

525 With our method, the initial parameters and covariance matrices are inferred from a dataset containing  
526 metadata of numerous PV systems, however, information on installed PV systems in terms of installed  
527 capacity and spatial distribution are still required. Although our method has a reduced dependency on data  
528 availability than alternative methods, this initial parametrization is compensated by Bayesian parameter  
529 estimation; regardless, this overarching PV data is still needed. An interesting solution arises from the  
530 impressive work of collection, analysis and synthesis of PV system metadata statistics presented by [Killinger  
531 et al. \(2018\)](#), which constitutes exemplary research that can potentially improve the performances of regional  
532 forecasts in the future by removing the data dependency.

533 We trained the model parameters on a NWP weather forecast, selected because the targeted application  
534 is day-ahead regional PV power forecasting. It would be interesting to train the model parameters on alter-  
535 native datasets, such as reanalysis or satellite data, to minimize the input error influence on the parameter  
536 estimation. Furthermore, integrate robust regression methods in the proposed approach to limit sensitivity  
537 to large input data errors is of interest.

538 In this work, we focused on the prediction of the regional PV power generation for the TSO but our  
539 approach can be used in any application where a regionally distributed fleet of PV systems needs to be  
540 modeled and information on the aggregated PV power production is available.

541 Finally, our approach could be extended to evaluate model parameters using different sources of informa-  
542 tion of different natures. The proposed framework indeed offers the possibility - with a minimal adaptation  
543 - to assimilate, in addition to aggregated power estimates, power measurements of individual PV systems  
544 or even aggregated power estimation of sub-regions (DSO estimates) to better assess the model paramete-  
545 rs. This would allow getting more reliable information on the fleet of installed PV system but also better  
546 characterizing the PV systems by using more sub-regions.

## 547 Acknowledgements

548 The authors would like to warmly thank Amprion, TenneT and 50 Hertz for the pleasant cooperation  
549 and for the highly interesting discussions on the RES integration mechanisms.

550 This research was partly funded by the ‘Bundesministerium für Wirtschaft und Energie’ and conducted  
551 within the project EWeLiNE (Erstellung innovativer Wetter-und Leistungsprognosemodelle für die Netzin-  
552 tegration wetterabhängiger Energieträger 0325500B).

553 J.M. Bright is funded by the Australian Renewable Energy Agency (ARENA, Research and Development  
554 Programme Funding G00854).

## 555 References

556 Amprion, 2019. Balancing group management. Accessed: 2019-02-27.

557 URL <https://www.amprion.net/Energy-Market/Balancing-Groups/Balancing-Group-Management/>

558 Boilley, A., Wald, L., 2015. Comparison between meteorological re-analyses from era-interim and merra and measurements of  
559 daily solar irradiation at surface. *Renewable Energy* 75, 135 – 143.

560 URL <http://www.sciencedirect.com/science/article/pii/S0960148114006077>

561 Bright, J. M., Killinger, S., Lingfors, D., Engerer, N. A., 2017. Improved satellite-derived PV power nowcasting using real-time  
562 power data from reference PV systems. *Solar Energy*.

563 URL <http://www.sciencedirect.com/science/article/pii/S0038092X17309714>

564 Colton, D., Kress, R., 1997. Inverse Acoustic and Electromagnetic Scattering Theory. Applied Mathematical Sciences. Springer  
565 Berlin Heidelberg.

566 URL <https://books.google.fr/books?id=U5wemv1DnNwC>

567 Crisan, A., Bain, D., 2009. Fundamentals of Stochastic Filtering, springer n Edition.

568 Crisan, D., Hambly, B., Zariphopoulou, T., 2014. Stochastic analysis and applications 2014: In honour of terry lyons.

569 URL <http://dx.doi.org/10.1007/978-3-319-11292-3>

570 Da Silva Fonseca, J. G., Oozeki, T., Ohtake, H., Shimose, K. I., Takashima, T., Ogimoto, K., 2014. Forecasting regional  
571 photovoltaic power generation - A comparison of strategies to obtain one-day-ahead data. *Energy Procedia* 57, 1337–1345.

572 URL <http://dx.doi.org/10.1016/j.egypro.2014.10.124>

573 Da Silva Fonseca Junior, J. G., Oozeki, T., Ohtake, H., ichi Shimose, K., Takashima, T., Ogimoto, K., 2014. Regional forecasts  
574 and smoothing effect of photovoltaic power generation in Japan: An approach with principal component analysis. *Renewable*  
575 *Energy* 68, 403–413.

576 URL <http://dx.doi.org/10.1016/j.renene.2014.02.018>

577 Engl, H., Hanke, M., Neubauer, A., 2000. Regularization of Inverse Problems. Mathematics and Its Applications. Springer  
578 Netherlands.  
579 URL <https://books.google.fr/books?id=VuEV-Gj1GZcC>

580 Fonseca Junior, J. G. d. S., Oozeki, T., Ohtake, H., Takashima, T., Ogimoto, K., oct 2015. Regional forecasts of photovoltaic  
581 power generation according to different data availability scenarios: a study of four methods. Progress in Photovoltaics:  
582 Research and Applications 23 (10), 1203–1218.  
583 URL <http://doi.wiley.com/10.1002/pip.2528>

584 Freitag, M., Potthast, R., 2013. Synergy of inverse problems and data assimilation techniques. Radon Series on Computational  
585 and Applied Mathematics. De Gruyter, Germany.

586 Hadamard, J., 1902. Sur les problèmes aux dérivées partielles et leur signification physique. Princeton Uni. Bull. (1902) 13,  
587 49–52.

588 IEA, 2018. 2018: Snapshot of global photovoltaic markets: Report IEA PVPS T1-33:2018.  
589 URL [http://www.iea-pvps.org/fileadmin/dam/public/report/statistics/IEA\\_PVPS-A\\_Snapshot\\_of\\_Global\\_PV-1992-2017.pdf](http://www.iea-pvps.org/fileadmin/dam/public/report/statistics/IEA_PVPS-A_Snapshot_of_Global_PV-1992-2017.pdf)

591 Ineichen, P., 2014. Long term satellite global, beam and diffuse irradiance validation. Energy Procedia 48, 1586 – 1596,  
592 proceedings of the 2nd International Conference on Solar Heating and Cooling for Buildings and Industry (SHC 2013).  
593 URL <http://www.sciencedirect.com/science/article/pii/S187661021400441X>

594 Killinger, S., 2017. A probabilistic approach to the estimation of regional photovoltaic power generation using meteorological  
595 data: Application of the Approach to the German Case. Ph.D. thesis, Karlsruher Institut für Technologie (KIT), Karlsruhe,  
596 Germany.  
597 URL [http://publica.fraunhofer.de/eprints/urn\\_nbn\\_de\\_0011-n-4942211.pdf](http://publica.fraunhofer.de/eprints/urn_nbn_de_0011-n-4942211.pdf)

598 Killinger, S., Lingfors, D., Saint-Drenan, Y.-M., Moraitis, P., van Sark, W., Taylor, J., Engerer, N., Bright, J., Oct. 2018.  
599 On the search for representative characteristics of PV systems: Data collection and analysis of PV system azimuth, tilt,  
600 capacity, yield and shading. Solar Energy 173, 1087 – 1106.  
601 URL <https://hal.archives-ouvertes.fr/hal-01882680>

602 Lorenz, E., Hurka, J., Karampela, G., Heinemann, D., Beyer, H., Schneider, M., 2008. Qualified Forecast of ensemble power  
603 production by spatially dispersed grid-connected PV systems. 23rd European Photovoltaic Solar Energy Conferenc (January),  
604 3285–3291.  
605 URL <http://tasks.iea-shc.org/publications/downloads/23rd{ }EU{ }PVSEC{ }5A0.8.6{ }lorenz.pdf>

606 Lorenz, E., Scheidsteger, T., Hurka, J., Heinemann, D., Kurz, C., nov 2011. Regional PV power prediction for improved grid  
607 integration. Progress in Photovoltaics: Research and Applications 19 (7), 757–771.  
608 URL <http://doi.wiley.com/10.1002/pip.1033>

609 Nakamura, G., Potthast, R., 2015. Inverse Modeling: An Introduction to the Theory and Methods of Inverse Problems and  
610 Data Assimilation. IOP expanding physics. Institute of Physics Publishing.  
611 URL [https://books.google.fr/books?id=toR\\_jgEACAAJ](https://books.google.fr/books?id=toR_jgEACAAJ)

612 Nuno Martinez, E., Koivisto, M. J., Cutululis, N. A., Sorensen, P., 2018. On the simulation of aggregated solar PV forecast  
613 errors. IEEE Transactions on Sustainable Energy, 1–10.

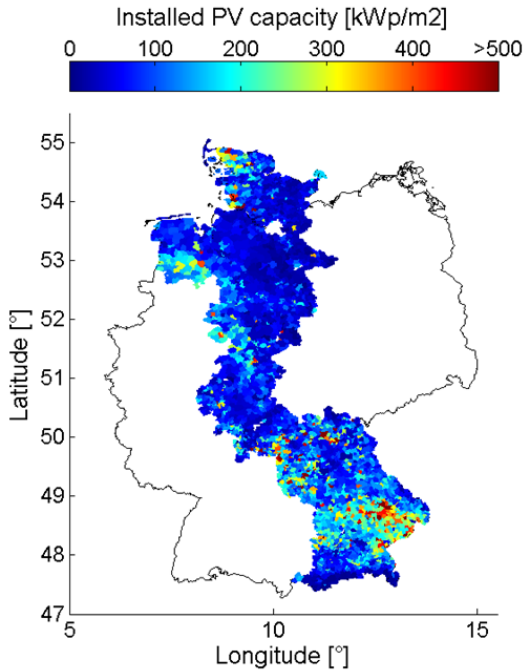
614 Saint-Drenan, Y., Good, G., Braun, M., 2017. A probabilistic approach to the estimation of regional photovoltaic power  
615 production. Solar Energy 147, 257–276.

616 Saint-Drenan, Y.-M., 2015. A probabilistic approach to the estimation of regional photovoltaic power generation using  
617 meteorological data: Application of the Approach to the German Case. Ph.D. thesis, University of Kassel, Kassel, Germany.  
618 URL <https://kobra.bibliothek.uni-kassel.de/bitstream/urn:nbn:de:hebis:34-2016090550868/3/DissertationYMSaintDrenan.pdf>

- 620 Schierenbeck, S., Graeber, D., 2010. Ein distanzbasiertes Hochrechnungsverfahren für die Einspeisung aus Photovoltaik.  
621 Energiewirtschaftliche tagesfragen 60 (12), 60–64.  
622 URL [http://www.transnetbw.de/downloads/kennzahlen/erneuerbare-energien/2010-fachartikel-graeber-semmig-schierenbeck-weber-](http://www.transnetbw.de/downloads/kennzahlen/erneuerbare-energien/2010-fachartikel-graeber-semmig-schierenbeck-weber-pdf)  
623 [pdf](#)
- 624 Schubert, G., 2012. Modeling hourly electricity generation from PV and wind plants in Europe. 9th International Conference  
625 on the European Energy Market, EEM 12, 1–7.
- 626 Shaker, H., Zareipour, H., Wood, D., 2015. A Data-Driven Approach for Estimating the Power Generation of Invisible Solar  
627 Sites. IEEE Transactions on Smart Grid PP (99), 1–11.
- 628 Tikhonov, A. N., 1963. Solution of Incorrectly Formulated Problems and the Regularisation Method.pdf. Soviet Mathematics  
629 Doklady 4 (4), 1035–1038.
- 630 Wilks, D. S., 2011. Statistical methods in the atmospheric sciences. Elsevier Academic Press, Amsterdam; Boston.  
631 URL [https://www.amazon.com/Statistical-Atmospheric-Sciences-International-Geophysics/dp/0123850223/ref=pb\\_](https://www.amazon.com/Statistical-Atmospheric-Sciences-International-Geophysics/dp/0123850223/ref=pb_xgxy_14_img_3?_encoding=UTF8&psc=1&refRID=ESPQQOR2PB1TP1VJSGCZ)  
632 [bxgy\\_14\\_img\\_3?\\_encoding=UTF8&psc=1&refRID=ESPQQOR2PB1TP1VJSGCZ](#)
- 633 Yang, D., Quan, H., Disfani, V. R., Rodríguez-Gallegos, C. D., 2017. Reconciling solar forecasts: Temporal hierarchy. Solar  
634 Energy 158, 332–346.

635 Appendix A. Description of the control area of the German TSOs

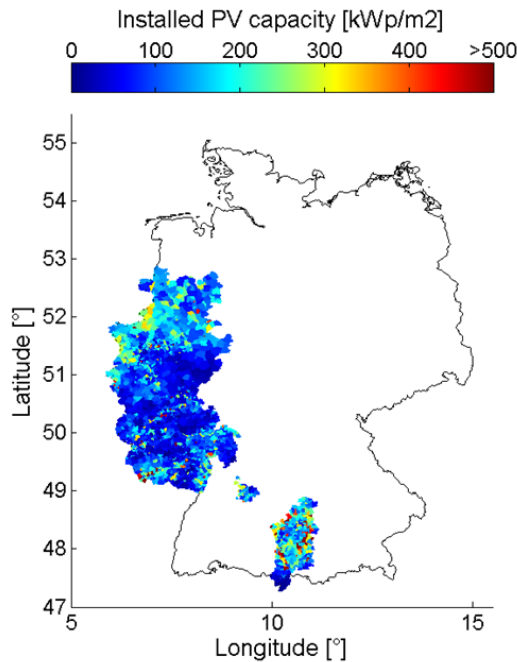
Description of the control area of TenneT



<b>Main characteristics of the TenneT control area</b>		
<i>(Update time: 05/2016)</i>		
Area		140.000 km <sup>2</sup>
Time of the last input		12-May-2016
Number of plants	0-30kW <sub>p</sub>	595 415 (88.7 %)
	30-100 kW <sub>p</sub>	63 347 (9.4 %)
	100-1000 kW <sub>p</sub>	11 346 (1.7 %)
	>1000 kW <sub>p</sub>	1 143 (0.2 %)
	<b>Total</b>	<b>671 251</b> <b>(100 %)</b>
Installed capacity	0-30kW <sub>p</sub>	6 569.2 MW <sub>p</sub> (42.8 %)
	30-100 kW <sub>p</sub>	3 075.3 MW <sub>p</sub> (20.0 %)
	100-1000 kW <sub>p</sub>	2 591.7 MW <sub>p</sub> (16.9 %)
	>1000 kW <sub>p</sub>	3 115.9 MW <sub>p</sub> (20.3 %)
	<b>Total</b>	<b>15 352.0 MW<sub>p</sub></b> <b>(100 %)</b>

636

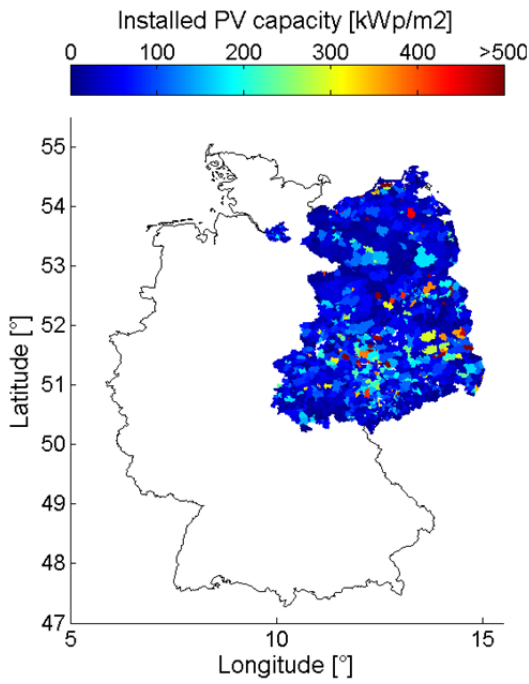
Description of the control area of Amprion



<b>Main characteristics of the Amprion control area</b>		
<i>(Update time: 05/2016)</i>		
Area		73.100 km <sup>2</sup>
Time of the last input		03-May-2016
Number of plants	0-30kW <sub>p</sub>	411 271 (89.3 %)
	30-100 kW <sub>p</sub>	40 387 (8.8 %)
	100-1000 kW <sub>p</sub>	8 623 (1.9 %)
	>1000 kW <sub>p</sub>	501 (0.1 %)
	<b>Total</b>	<b>460 782</b> <b>(100 %)</b>
Installed capacity	0-30kW <sub>p</sub>	4 004.5 MW <sub>p</sub> (43.1 %)
	30-100 kW <sub>p</sub>	2 039.2 MW <sub>p</sub> (21.9 %)
	100-1000 kW <sub>p</sub>	1 927.7 MW <sub>p</sub> (20.7 %)
	>1000 kW <sub>p</sub>	1 329.2 MW <sub>p</sub> (14.3 %)
	<b>Total</b>	<b>9 300.6 MW<sub>p</sub></b> <b>(100 %)</b>

637

Description of the control area of 50 Hertz



<b>Main characteristics of the 50 Hertz control area</b>		
<i>(Update time: 05/2016)</i>		
<i>Area</i>		109.360 km <sup>2</sup>
<i>Time of the last input</i>		20-Apr-2016
<i>Number of plants</i>	0-30kW <sub>p</sub>	118 296 (85.3%)
	30-100 kW <sub>p</sub>	12 910 (9.3 %)
	100-1000 kW <sub>p</sub>	5923 (4.3%)
	>1000 kW <sub>p</sub>	1565 (1.1%)
	<b>Total</b>	<b>138 694</b> <b>(100%)</b>
<i>Installed capacity</i>	0-30kW <sub>p</sub>	983.2 MW <sub>p</sub> (11.0%)
	30-100 kW <sub>p</sub>	626.5 MW <sub>p</sub> (7.0 %)
	100-1000 kW <sub>p</sub>	1 777.5 MW <sub>p</sub> (19.9 %)
	>1000 kW <sub>p</sub>	5 533.2 MW <sub>p</sub> (62.0 %)
	<b>Total</b>	<b>8920.4 MW<sub>p</sub></b> <b>(100 %)</b>

**a) TenneT control area**

Number of values: 12 473

Average power: 0,2091  $W/W_p$

	Persistence	First guess	Bayesian method min / max	OLS regression min / max	TSO forecast
correlation [-]	0,8815	0,9724	0.9737 / 0.9738	0.9723 / 0.9729	0,9806
bias [ $W/W_p$ ]	-0,0008	-0,0038	-0.0059 / -0.0057	-0.006 / -0.0053	0,0131
std [ $W/W_p$ ]	0,0808	0,0388	0.0378 / 0.0379	0.0384 / 0.0388	0,0348
MAE [ $W/W_p$ ]	0,0538	0,0263	0.0258 / 0.0259	0.0264 / 0.0266	0,0249
RMSE [ $W/W_p$ ]	0,0808	0,0390	0.0382 / 0.0383	0.0389 / 0.0392	0,0372

**b) Amprion control area**

Number of values: 12 451

Average power: 0,2307  $W/W_p$

	Persistence	First guess	Bayesian method min / max	OLS regression min / max	TSO forecast
correlation [-]	0,8601	0,9731	0,9730 / 0,9732	0,9735 / 0,9746	0,9638
bias [ $W/W_p$ ]	0,0001	-0,0034	-0,0015 / -0,0013	-0,0032 / -0,0028	0,0013
std [ $W/W_p$ ]	0,0972	0,0423	0,0422 / 0,0424	0,0411 / 0,0420	0,0527
MAE [ $W/W_p$ ]	0,0637	0,0286	0,0284 / 0,0285	0,0280 / 0,0287	0,0360
RMSE [ $W/W_p$ ]	0,0972	0,0425	0,0423 / 0,0424	0,0412 / 0,0421	0,0527

**c) 50 Hertz control area**

Number of values: 12 111

Average power: 0,2380  $W/W_p$

	Persistence	First guess	Bayesian method min / max	OLS regression min / max	TSO forecast
correlation [-]	0,8710	0,9712	0,9720 / 0,9722	0,9716 / 0,9720	0,9736
bias [ $W/W_p$ ]	0,0022	0,0018	-0,0046 / -0,0045	-0,0048 / -0,0043	-0,0011
std [ $W/W_p$ ]	0,0973	0,0464	0,0449 / 0,0451	0,0451 / 0,0454	0,0438
MAE [ $W/W_p$ ]	0,0646	0,0302	0,0299 / 0,0300	0,0299 / 0,0302	0,0285
RMSE [ $W/W_p$ ]	0,0973	0,0464	0,0451 / 0,0453	0,0454 / 0,0457	0,0438

Table 2: Different error metrics obtained with the various forecast approaches for the control areas of a) TenneT, b) Amprion and c) 50 Hertz