

Probabilistic prediction of solar power supply to distribution networks, using forecasts of global horizontal irradiation

F. von Loeper^{a,*}, P. Schaumann^a, M. de Langlard^a, R. Hess^b, R. Bäsman^c,
V. Schmidt^a

^aUniversity Ulm, Institute of Stochastics, Helmholtzstraße 18, 89081 Ulm, Germany

^bDeutscher Wetterdienst, Frankfurter Straße 135, 63067 Offenbach, Germany

^cN-ERGIE Netz GmbH, Sandreuthstraße 21, 90441 Nürnberg, Germany

Abstract

This paper presents a mathematical model for the prediction of the probabilities of reverse power flow exceeding predefined critical thresholds at feed-in points of a distribution network. The parametric prediction model is based on hourly forecasts of global horizontal irradiation and uses copulas, a tool for modeling the joint probability distribution of two or more strongly correlated random variables with non-Gaussian (marginal) distributions. The model is used for determining the joint distribution of forecasts of global horizontal irradiation and measured solar power supply at given feed-in points, where respective sample datasets were provided by Deutscher Wetterdienst and the N-ERGIE Netz GmbH. It is shown that the fitted model replicates important characteristics of the data such as the corresponding marginal densities. The validation results highlight strong performance of the proposed model. The copula-based model enables to predict the distribution of solar power supply conditioned on the forecasts of global horizontal irradiation, thus anticipating great fluctuations in the distribution network.

Keywords: Probabilistic prediction model, Global horizontal irradiation, Solar power supply, Mixed beta distribution, Archimedean copula

*Corresponding author

Email address: `freimut.von-loeper@uni-ulm.de` (F. von Loeper)

1. Introduction

In the recent decade, the global annually installed capacity of solar power increased rapidly, reaching 98 Gigawatts in 2017 alone. Compared to other power generation sources connected to electricity networks, solar power has the greatest capacity installed in 2017, followed by wind power with 52 Gigawatts, gas power with 38 Gigawatts, coal power with 35 Gigawatts and various other sources adding up to 37 Gigawatts. Although the global capacity of solar power installed in 2017 already exceeded most expectations, solar analysts predict even further increase of the annually installed capacity for the future. Thus, it comes as no surprise that the worldwide installed capacity of solar power is estimated to exceed 1 Terrawatt by 2022 (SolarPower Europe, 2017).

The increase in solar penetration causes greater fluctuations in the power supply, which might result in increasing overloading problems and voltage violations (Karimi et al., 2016). To tackle the upcoming challenges for distribution network operators, suitable computer-based models are developed. Since solar power is a variable power source, better forecasts of solar power supply are useful to predict oversupply events and reduce the curtailment (Bird et al., 2014). Moreover, planning maintenance becomes easier and bids in electricity markets can be better estimated. For a deeper understanding of the economical value of forecasting solar power supply we refer to Antonanzas et al. (2016).

Most prediction models of solar power supply rely on weather forecasts as input (Antonanzas et al., 2016). Usually, these weather forecasts are the result of a complex modeling process starting with numerical weather prediction models based on data assimilation that use all kinds of meteorological observations in real time including e.g. temperature, wind, pressure, humidity, as well as geospatial data as input information (Coiffier, 2011). Based on this data, the task of numerical weather prediction models such as the high resolution version COSMO-DE of the Consortium for Small-scale Modeling (COSMO) run by Deutscher Wetterdienst (DWD) is to solve differential equations for modeling the spatio-temporal evolution of meteorological variables and simulate physi-

31 cal processes (Baldauf et al., 2011). However, numerical weather prediction
32 models are subject to systematic errors and not all weather phenomena are
33 simulated and need to be interpreted therefore. For these reasons, statistical
34 post-processing methods are applied. For example, DWD runs Model Output
35 Statistics (MOS) techniques (Hess et al., 2015), which are based on multiple
36 linear and logistic regression models (Jobson, 1991). The regression models fit
37 historical data of the direct output of numerical models to corresponding obser-
38 vations from synoptic stations. In operational use, the fitted regression models
39 calibrate and interpret the output of numerical weather predictions resulting
40 in deterministic and probabilistic forecasts for various meteorological variables
41 (Heinemann et al., 2006).

42 Since energy conversion by solar plants mainly depends on direct and dif-
43 fuse radiation, our prediction model is based on deterministic and statistically
44 post-processed forecasts of global horizontal irradiation (GHI). Other meteoro-
45 logical variables such as ambient temperature, wind velocity, humidity and dust,
46 may also influence the energy conversion and have been used in the literature
47 (Kaldellis et al., 2014; Mekhilef et al., 2012). Moreover, probabilistic models
48 which take into account the spatio-temporal correlation between these variables
49 were also proposed. Almeida et al. (2015) developed a non-parameteric quantile
50 regression forest model which takes as training input temperature, wind speed,
51 wind direction, humidity, sea level pressure and cloud cover at different levels.
52 Bessa et al. (2015) proposed a vector auto-regressive model with time series
53 information collected at different locations on a smart grid as input. In Zhang
54 et al. (2016) a Gaussian conditional random field model was applied, where
55 historical forecasts and solar power measurements were considered at many so-
56 lar sites. Huang and Perry (2016) estimated prediction intervals based on a
57 k -nearest neighbor regression and added to deterministic forecasts computed by
58 gradient boosting, where weather variables such as solar radiation, temperature,
59 cloud ice water content, wind speed were considered. Solar power measurements
60 at adjacent solar farms were included as explanatory regression variables. In
61 contrast we apply a copula model which only considers the GHI forecasts, the

62 most important explanatory variable for energy conversion. This simplifies the
63 fitting procedure of the model due to a smaller number of parameters, while
64 still capturing the correlation structure between weather conditions and solar
65 power supply.

66 Copulas are a mathematical tool to model the joint distribution of two or
67 more random variables. In the context of renewable energies, they were first
68 applied for the probabilistic prediction of wind power generation, see e.g. Pa-
69 paefthymiou and Kurowicka (2009); Wang et al. (2014); Lu et al. (2014). Re-
70 cently, copula models have also been used for statistical analysis of data on
71 solar power generation. For example, Golestaneh et al. (2016a,b) applied quan-
72 tile regression to non-parametrically compute conditional marginal densities of
73 solar power supply for neighboring solar plants, given numerical weather pre-
74 diction forecasts. Then, in a next step, multivariate Gaussian copulas were used
75 in Golestaneh et al. (2016a) to determine the joint conditional distribution of
76 solar power supply at neighboring plants with the previously computed non-
77 parametric conditional marginal densities. Golestaneh and Gooi (2017) com-
78 pared Gaussian copulas with multivariate R-vine copulas. However, in both
79 papers copulas were merely used to model the spatial relationship between so-
80 lar power supply at neighboring plants. More recently, Panamtash et al. (2020)
81 proposed a similar copula-based model, where bivariate copulas are applied to
82 improve prior probabilistic forecasting done by traditional forecasting methods
83 such as multiple linear regression, artificial neural networks, gradient boost-
84 ing, random forests and autoregressive integrated moving average. The results
85 showed that copulas capture the joint probability distribution of solar power
86 and temperature effectively. In Panamtash et al. (2020), ambient temperature
87 has been considered as input variable, while we focus on the correlation between
88 GHI forecasts and solar power supply.

89 Given a weather forecast, our proposed model computes conditional proba-
90 bilities of reverse power flow exceeding predefined critical thresholds at feed-in
91 points of a distribution network. To implement the prediction model, the first
92 step is to fit univariate (so-called marginal) probability distributions using his-

93 torical data of hourly GHI forecasts and hourly averages of measured solar
94 power supply. The second step is to model the joint probability distribution of
95 GHI forecast and solar power supply by applying copula theory (Durante and
96 Sempi, 2015; Nelsen, 2006; Joe, 2014). Finally, the third step is to compute the
97 conditional probability distribution of solar power supply based on the fitted
98 joint and marginal distributions. Taking a real-time GHI forecast as input, a
99 probabilistic prediction of solar power supply for the same time horizon as the
100 weather forecast can be computed by the prediction model.

101 The rest of this paper is organized as follows. In Section 2, the data used in
102 this paper is described, including its pre-processing and analysis. The model and
103 its fitting procedure is explained in Section 3. The fitted model characteristics
104 and the validation of the prediction model are discussed in Section 4, where also
105 the performance of the proposed model is compared with that of the quantile
106 regression technique, one of the most frequently used probabilistic prediction
107 method (see, for instance, Bacher et al. (2009); Zamo et al. (2014); Massidda and
108 Marrocu (2018); Lauret et al. (2017); Golestaneh et al. (2016b); Alessandrini
109 et al. (2015)). Finally, Section 5 concludes.

110 **2. Data**

111 The modeling approach for the prediction of solar energy supply, proposed in
112 the present paper, is a parametric probabilistic model which is based on copulas
113 (Durante and Sempi, 2015; Nelsen, 2006; Joe, 2014). Compared to conventional
114 photovoltaic performance models, a probabilistic model is more flexible, but
115 needs historical data as input (Antonanzas et al., 2016). In particular, the gen-
116 eral modeling idea is not concerned with physical attributes of the datasets, e.g.
117 the locations of the measurement points. Our probabilistic modeling approach
118 and its application are illustrated by using suitable sample datasets provided
119 by DWD and the N-ERGIE Netz GmbH (NNG). For calibration and validation
120 of our model, the time frame covering the months May, June and July of the
121 years 2015 till 2017 is considered, resulting in 273 days.

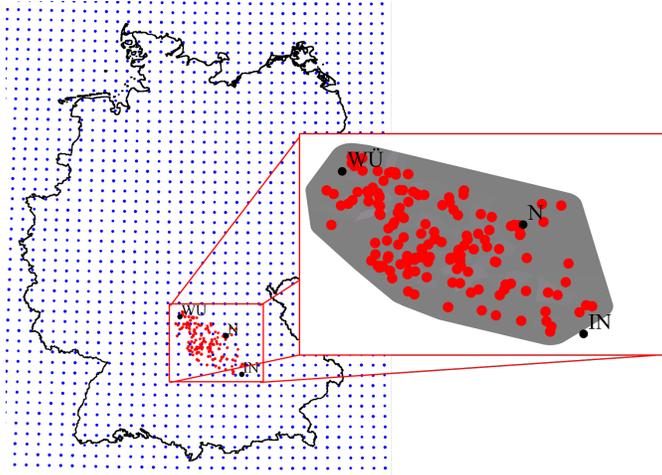


Figure 1: Forecast grid of DWD (blue) and feed-in points of NNG (red). The (appropriately dilated) convex hull (grey) of all feed-in points in the zoom-in of the red box is the part of Germany used for visualizing our results. To illustrate the geographical location of this area, the cities of Würzburg (WÜ), Nürnberg (N) and Ingolstadt (IN) are depicted (black).

122 *2.1. Description of data*

123 The fitting of our model is based on two datasets, namely GHI forecasts
 124 provided by DWD and measured solar power supply provided by NNG.

125 The first dataset consists of hourly GHI forecasts (in kJ/m^2), statistically
 126 interpreted based on synoptic observations and numerical forecasts of COSMO-
 127 DE-EPS, the ensemble system of COSMO-DE at DWD. The forecasts are issued
 128 every three hours with forecast lead times up to 19 hours. For the time frame
 129 mentioned above, the forecasts are available on a $20\text{ km} \times 20\text{ km}$ grid covering
 130 Germany and parts of neighboring countries, see Figure 1. However, there is
 131 no GHI forecast generated for grid points and forecasts times with a local solar
 132 elevation angle of less than 5 degrees at the beginning or end of the forecast
 133 hour.

134 The second dataset, provided by NNG, records the amount of electricity,
 135 which was generated by solar plants, supplied to the distribution network and
 136 measured at its feed-in points. These amounts of solar power supply are 15-
 137 minute average values. Each of the 168 feed-in points, considered in this paper,

138 is connected with at least one solar plant with nominal capacity of the connected
139 solar plants ranging from 0.25 up to 10 Megawatts.

140 In Figure 1 the considered feed-in points of NNG are visualized. The corre-
141 sponding supply area covers ca. 8000 km².

142 *2.2. Data preprocessing*

143 At first temporal and spatial compatibility between the datasets has to be
144 established. Temporal compatibility can be easily realized by calculating the
145 averages of solar power supply for each hour. For spatial compatibility we con-
146 sider two different hierarchy levels in the distribution network. On the one hand
147 we are interested in the GHI forecasts and amounts of solar power supply at the
148 feed-in points, on the other hand we want to apply our model to communities,
149 i.e. sets of neighboring feed-in points, as well. Therefore, we match GHI fore-
150 casts and amounts of solar power supply to feed-in points and communities, see
151 Section 2.2.1 and 2.2.2.

152 *2.2.1. Spatial compatibility for feed-in points*

153 In this section we consider the problem to match the amount of solar energy
154 measured at every feed-in point to a single GHI forecast. For simplicity, the
155 locations of solar plants are assumed to coincide with the locations of their
156 feed-in points. Since the hourly averages of forecasted GHI are practically the
157 same on such small spatial scales, the error introduced by this assumption is
158 negligible. The GHI forecast at a certain feed-in point has been estimated by
159 interpolating the GHI forecast at the grid points of the 20 km × 20 km grid, see
160 Figure 1. This is done by bilinear interpolation, see Hämmerlin and Hoffmann
161 (2012).

162 *2.2.2. Spatial compatibility for communities*

163 In a next step, the amounts of solar power supply and the GHI forecasts need
164 to be matched to communities. Therefore, we approximate the solar power sup-
165 ply generated in a community by summing the solar power supply measured at
166 the feed-in points over all feed-in points within the community. It is assumed

167 that there are no transmission losses, but alternatively the losses can be com-
168 puted either using explicit formulas or statistical estimations, see Dickert et al.
169 (2009), Council of European Energy Regulators (2017).

170 To compute the GHI forecasts for the corresponding communities, we use the
171 interpolated GHI forecasts at the feed-in points in a community as mentioned
172 in Section 2.2.1. By averaging the GHI forecasts over all feed-in points in a
173 community, we get the GHI forecast of the community.

174 2.2.3. Selection of data

175 Due to maintenance, repair work and risk of overloading, some solar plants
176 might have to be shut down for certain time periods. Since these actions are
177 not directly related with weather phenomena, corresponding time periods are
178 excluded by removing solar power supply being equal to zero from data. The
179 GHI forecasts for those locations and forecast times are also removed from data,
180 leaving only data pairs with matching time stamps.

181 Furthermore, the performance of solar plants is strongly influenced by many
182 factors apart from meteorological variables, e.g. nominal capacity, tilt angle and
183 composition of photovoltaic units. Most of these factors are constant over long
184 time periods, but their influence might largely depend on the time of day as it
185 is the case with the tilt angle. A simple way to remove such effects from data
186 is to consider each hour of the day and feed-in point separately.

187 Lastly, by matching both data sets, all hours which have no measurement of
188 solar power supply or GHI forecast are removed. This includes all night hours,
189 see Section 2.1.

190 2.2.4. Rescaling the data

191 For an easier comparison of the datasets with different scales, the interpo-
192 lated GHI forecasts and measured amounts of solar power are locally normal-
193 ized. More precisely, the datasets are rescaled for a certain feed-in point, or
194 community, by applying the transformation

$$195 \quad \phi_{a,b}(x) = (1 - 2c)(x - a)/(b - a) + c, \quad (1)$$

196 where $c > 0$ is close to zero. If a is the minimum and b the maximum of the
 197 dataset under consideration, then, ϕ maps onto the interval $[c, 1 - c]$. For $c > 0$
 198 close to zero, $[c, 1 - c]$ approximates the open interval $(0, 1)$. Thus, we use
 199 $c = 0.001$ in this paper.

200 2.3. Empirical data analysis

201 The rescaled and interpolated data of solar power supply and GHI forecast
 202 for feed-in points are analyzed in order to highlight their strong correlation and
 203 their spatial disparities. We consider the time frame May, June and July of
 204 the years 2015 till 2017 (11-12 UTC), denoted by T , and GHI forecasts with
 205 forecast lead time of one hour. The method, described in Section 2.2, is applied
 206 for each feed-in point and dataset separately, where the rescaling parameters a
 207 and b , in Section 2.2.4, are set to the minimum and maximum of each dataset
 208 in T .

209 Therefore, the local *empirical correlation coefficient* of GHI forecasts and
 210 amounts of solar power supply at each feed-in point ℓ , denoted by $\rho(\ell)$, is com-
 211 puted. Note that the empirical correlation coefficient $\rho(\ell)$ is defined as

$$212 \rho(\ell) = \frac{\sum_{t \in T} (r(\ell, t) - \bar{r}(\ell))(s(\ell, t) - \bar{s}(\ell))}{\sqrt{\sum_{t \in T} (r(\ell, t) - \bar{r}(\ell))^2 \sum_{t \in T} (s(\ell, t) - \bar{s}(\ell))^2}}, \quad (2)$$

213 where $r(\ell, t)$ and $s(\ell, t)$ are the preprocessed GHI forecast and solar power supply
 214 for forecast time t , respectively, and $\bar{r}(\ell)$ and $\bar{s}(\ell)$ are the corresponding time
 215 averages.

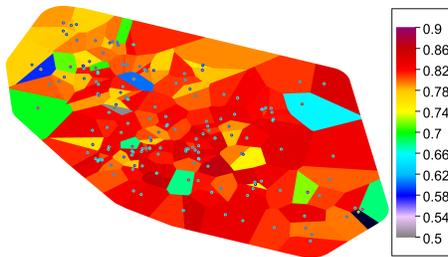


Figure 2: Empirical correlation coefficients of preprocessed GHI forecasts and preprocessed amounts of solar power supply for each feed-in point visualized by Voronoi tessellation.

216 In Figure 2, the results are visualized, which we obtained for the local em-
 217 pirical correlation coefficients, where the Region of Interest (ROI), i.e., the (ap-
 218 propriately dilated) convex hull of all feed-in points, is decomposed into the
 219 Voronoi tessellation generated by the feed-in points. The value computed for a
 220 feed-in point is assigned to the entire Voronoi cell of this point.

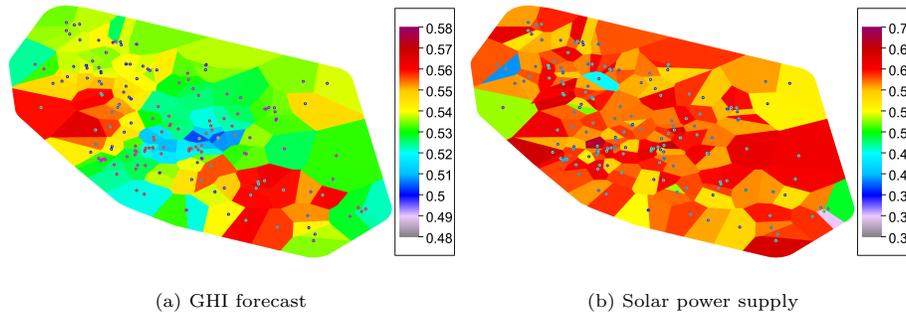


Figure 3: Local means over 273 days for each feed-in point visualized by Voronoi tessellation.

221 In general, Figure 2 shows rather high empirical correlation coefficients for
 222 all feed-in points, but also the existence of significant differences between feed-in
 223 points. By depicting the local means of the preprocessed datasets in Figure 3,
 224 which are obtained by averaging over the 273 days considered in this paper,
 225 this observation becomes even more evident. Indeed, the map of local means
 226 of the GHI forecasts does not show the same kind of pattern as the one of the
 227 solar power supply, see Figure 3. Thus, there have to be other factors except
 228 GHI forecasts, e.g. physical characteristics of the solar plants connected with
 229 the feed-in points, influencing the local means of solar power supply shown in
 230 Figure 3.

231 As a conclusion of this empirical analysis, our modeling approach has to con-
 232 sider the interdependence of GHI forecasts and solar power supply, see Figure
 233 2. Moreover, Figure 3 indicates that the parameters of the probability distri-
 234 butions considered in this paper should be determined for each feed-in point
 235 separately to take into account their spatial variability.

236 **3. Copula-based model for the prediction of solar power supply**

237 In many fields, where risk must be managed, probabilistic predictions are
 238 preferred as they allow to quantify the uncertainty. In our case, given some
 239 probability estimation for the occurrence of a critical feed-in event, distribu-
 240 tion network operators might make their decisions individually based on how
 241 much risk they want to take. Further advantages of probabilistic forecasts are
 242 discussed, e.g., in Antonanzas et al. (2016).

243 *3.1. Modeling approach*

244 For simplicity, we consider preprocessed solar power supply and preprocessed
 245 GHI forecast at a certain location in the distribution network for a single time of
 246 day and forecast lead time. Both are interpreted as realizations of some random
 247 variables R and S , which are strongly correlated as Figure 2 in Section 2.3
 248 depicts. The support of R and S is the interval $[0, 1]$, because of the rescaling
 249 transformation given in Eq. (1).

250 Given a GHI forecast r the conditional probability of solar power supply
 251 exceeding a certain threshold $v \in [0, 1]$ can be written in the following form:

$$252 \quad P(S \geq v \mid R = r) = \int_v^1 f_{S|R}(s \mid r) ds \quad (3)$$

$$253 \quad = \int_v^1 \frac{f_{(R,S)}(r, s)}{f_R(r)} ds, \quad (4)$$

254

255 where $f_{S|R}$ is the conditional density of the random variable S given R , $f_{(R,S)}$
 256 the joint density of the random variables S and R , and f_R the marginal density
 257 of R . Thus, our task lies in modeling the marginal and joint distributions of S
 258 and R .

259 The proposed method can be decomposed in the following two steps:

- 260 1. derive a parametric form of the marginal densities f_R and f_S of the random
 261 variables R and S ;
- 262 2. use a copula and the marginal densities f_R and f_S to model the joint
 263 density function $f_{(R,S)}$.

264 Then, the conditional level-crossing probabilities $P(S \geq v \mid R = r)$ can be
 265 computed using Eq. (4). In the following, these two steps are further explained
 266 in detail.

267 3.1.1. Model for the marginal densities

268 Vale (2015) investigated several types of parametric distributions for solar
 269 irradiation in Lisbon with respect to the day time and month, where it is con-
 270 cluded that the mixed beta distribution is the best fit for most months of the
 271 year. However, the difference in the considered time frame and location might
 272 influence the quality of fits for the tested distribution types.

273 In Section 4.1.1, we demonstrate that the mixed beta distribution is indeed
 274 suitable for GHI forecasts and also for solar power supply. Therefore, we choose
 275 the models of the marginal densities f_S and f_R to be a mixture of beta densities.
 276 The family of beta distributions allows for various shapes of probability density
 277 functions.

278 Given a mixing parameter $q \in (0, 1)$ and two densities of beta distributions
 279 $f_i : \mathbb{R} \rightarrow [0, \infty)$ with $i \in \{1, 2\}$, the probability density $f_X : \mathbb{R} \rightarrow [0, \infty)$ of the
 280 *mixed beta distribution* of a random variable X is given by

$$281 \quad f_X(x) = qf_1(x) + (1 - q)f_2(x) \quad (5)$$

282 for all $x \in \mathbb{R}$. For $i \in \{1, 2\}$, the *beta density* f_i has two shape parameters $a_i > 0$
 283 and $b_i > 0$, and is defined by

$$284 \quad f_i(x) = \frac{\Gamma(a_i + b_i)}{\Gamma(a_i)\Gamma(b_i)} x^{a_i-1}(1-x)^{b_i-1} \quad (6)$$

285 for $x \in (0, 1)$ and $f_i(x) = 0$ otherwise, where Γ denotes the gamma function.
 286 Then, the marginal densities f_R and f_S take the form of the mixed density given
 287 in Eq. (5) and are specified by five parameters each. The determination of these
 288 parameters is explained in Section 3.2.

289 3.1.2. The copula model

290 If R and S were independent random variables, the bivariate density function
 291 $f_{(R,S)}$ would turn out to be the product of the univariate densities f_R and f_S .

292 But, in fact, the random variables R and S are strongly correlated as depicted
 293 in Figure 2. Therefore, the design of a parametric joint density is more complex.

294 To calculate the joint density with the non-Gaussian marginal densities fit-
 295 ted in Section 3.2, a parametric modeling approach based on Sklar’s theorem is
 296 applied (Durante and Sempi, 2015). This fundamental result of copula theory
 297 allows us to represent the bivariate joint distribution function of two random
 298 variables by superposing a copula function upon the marginal distribution func-
 299 tions. Note that a *copula* is defined as the joint cumulative distribution function
 300 $C : [0, 1] \times [0, 1] \rightarrow [0, 1]$ of a two-dimensional random vector (U, V) with com-
 301 ponents U and V uniformly distributed on $[0, 1]$, see Nelsen (2006) for further
 302 details.

303 Let (R, S) be the two-dimensional random vector consisting of the random
 304 variables R and S introduced above with joint cumulative distribution function
 305 $F_{(R,S)} : \mathbb{R}^2 \rightarrow [0, 1]$ and marginal distribution functions F_R and F_S . Then,
 306 *Sklar’s theorem* says that a copula function $C : [0, 1] \times [0, 1] \rightarrow [0, 1]$ exists such
 307 that

$$308 \quad F_{(R,S)}(r, s) = C(F_R(r), F_S(s)) \quad (7)$$

309 for all $r, s \in \mathbb{R}$. Note that Eq. (7) can be written in the following differential
 310 form:

$$311 \quad f_{(R,S)}(r, s) = f_R(r) \cdot f_S(s) \cdot c(F_R(r), F_S(s)), \quad (8)$$

312 where $f_{(R,S)}$, f_R , f_S and c are the densities corresponding to the cumulative
 313 distribution functions $F_{(R,S)}$, F_R , F_S and C . Using (8), the conditional den-
 314 sity function $f_{S|R}(s | r)$ of solar power supply given a GHI forecast r can be
 315 computed by

$$316 \quad f_{S|R}(s | r) = \frac{f_{(R,S)}(r, s)}{f_R(r)} = f_S(s) \cdot c(F_R(r), F_S(s)). \quad (9)$$

317 Thus, to determine the conditional level-crossing probability considered in (4),
 318 we estimate the marginal densities f_S and f_R , which leads to estimates of the
 319 corresponding distribution functions F_S and F_R , and the copula density c in
 320 Eq. (9). For the estimation of c we apply *Archimedean copulas*. Archimedean

321 copulas are a commonly considered class of copulas which can be given by
 322 analytical formulas and are therefore especially easy to handle.

323 A function $g : [0, 1] \rightarrow [0, \infty]$ is called an *Archimedean generator* if g is
 324 continuous, strictly decreasing and solves $g(1) = 0$. The pseudo-inverse $g^{[-1]}$ of
 325 an Archimedean generator g is an extension of the inverse function $g^{(-1)}$ defined
 326 as

$$327 \quad g^{[-1]}(t) = \begin{cases} g^{(-1)}(t), & \text{if } 0 \leq t \leq g(0), \\ 0, & \text{if } g(0) < t \leq \infty. \end{cases} \quad (10)$$

328 The Archimedean copula generated by g is then given by

$$329 \quad C(u, v) = g^{[-1]}(g(u) + g(v)) \quad (11)$$

330 for $u, v \in [0, 1]$.

331 In this paper, we focus on four parametric types of Archimedean copulas,
 332 see Table 1, each of them having a single parameter $\theta \in \mathbb{R}$ to be fitted.

Type	Archimedean generator	Parameter
Joe	$g_\theta(t) = -\log(1 - (1 - t)^\theta)$	$\theta \in [1, \infty)$
Frank	$g_\theta(t) = (-\log(\frac{\exp(-\theta t) - 1}{\exp(-\theta) - 1}))^\theta$	$\theta \in \mathbb{R} \setminus \{0\}$
Clayton	$g_\theta(t) = \frac{1}{\theta}(t^{-\theta} - 1)$	$\theta \in [-1, \infty) \setminus \{0\}$
Gumbel	$g_\theta(t) = (-\log(t))^\theta$	$\theta \in [1, \infty)$

Table 1: Types of Archimedean copulas

333 3.2. Model fitting procedure

334 In this section, we describe the procedure to find the best model parameters
 335 of the marginal densities f_S and f_R , which are modeled by mixtures of beta
 336 densities, see Eq. (5). Next, we detail the method to search for the copula type
 337 and the copula parameter θ , which gives us the best fit to the data. Both fitting
 338 procedures are based on the *maximum likelihood estimation* (MLE) principle,
 339 see Wilks (2011).

340 *3.2.1. Fitting marginal densities*

341 There are five parameters to be determined for each marginal density f_R
 342 and f_S : the mixing parameter q , the shape parameters a_1 and b_1 of the first
 343 beta density f_1 , and the shape parameters a_2 and b_2 of the second beta density
 344 f_2 , see Eq. (5).

345 The main idea is to consider the maximum of a certain product of likelihood
 346 functions. Note that the likelihood function is defined, in the case of the random
 347 variable R , by

$$348 \quad L(\beta | r) = f_R(r), \quad (12)$$

349 where $\beta = (q, a_1, b_1, a_2, b_2)$ is the set of parameters and r is a GHI forecast. If
 350 we consider the dataset r_1, \dots, r_n and assume that the observations r_1, \dots, r_n
 351 are independently sampled realizations of the random variable R (n is the total
 352 number of observations), then the MLE consists in maximizing the function

$$353 \quad L(\beta | r_1, \dots, r_n) = \prod_{i=1}^n L(\beta | r_i) \quad (13)$$

$$354 \quad = \prod_{i=1}^n f_R(r_i) \quad (14)$$

$$355 \quad = \prod_{i=1}^n q f_1(r_i) + (1 - q) f_2(r_i), \quad (15)$$

356
 357 where the densities f_1 and f_2 depend on the parameters a_1 and b_1 , respectively
 358 a_2 and b_2 . The MLE expresses the fact that the realized observations occur
 359 with the highest possible probability. In general, it is more convenient to take
 360 the logarithmic form of the expression given in (15) in order to deal with a sum
 361 instead of product operations. Then, the best parameters q, a_1, b_1, a_2 and b_2
 362 for the marginal density of GHI forecasts are the solution of the maximization
 363 problem

$$364 \quad \arg \max_{q, a_1, b_1, a_2, b_2} \sum_{i=1}^n \log (q f_1(r_i) + (1 - q) f_2(r_i)). \quad (16)$$

365 The estimated parameters for the marginal density f_S of the solar power supply
 366 are the solution of an analogous maximization problem.

367 Note that the maximization problem stated in (16) is difficult to solve be-
368 cause of its high dimensionality. The iterative *expectation-maximization algo-*
369 *rithm* (EM algorithm) proposed in Dempster et al. (1977) is particularly ap-
370 propriate to deal with the maximization of such functions as it enables us to
371 decompose the maximization problem stated in (16) into easier sub-problems.
372 For further details regarding the EM algorithm, see Hastie et al. (2009), Leisch
373 (2004).

374 3.2.2. Fitting the copula function

375 To estimate the copula parameter θ , see Table 1, the *inference function for*
376 *margins method* proposed in Joe and Xu (1996) is applied. The first step of this
377 method is to estimate the marginal distributions, as described in Section 3.2.1.
378 In the second step, the copula parameter θ is fitted using the MLE principle
379 based on the previously estimated marginal distributions.

380 Indeed, for each copula type in Table 1, the copula density $c^{(\theta)}$ can be
381 obtained by differentiating Eq. (11) in dependence of the copula parameter θ .
382 Using Eq. (8) and the previously fitted marginal distributions we compute the
383 joint density $f_{(R,S)}^{(\theta)}$ for each copula type in dependence of θ . Then, we apply
384 the MLE method to fit the joint density $f_{(R,S)}^{(\theta)}$ for each copula type to the data.
385 As a result, we get an estimate of θ and the corresponding maximum of the
386 product of likelihood functions for each copula type. We choose the copula type
387 and its copula parameter, which gives the largest maximum.

388 Note that the inference function for margins method represents a certain
389 break with the classical MLE principle, where all model parameters, including
390 the parameters of the marginal distributions, are simultaneously estimated. As
391 an alternative, in Joe (2014) it is proposed to use the components of param-
392 eter vector β determined by the two-step procedure in Section 3.2.1 as initial
393 values for iterative numerical methods, which estimate all model parameters
394 simultaneously. But this has not been done in the present paper.

395 4. Results

396 To begin with, in Section 4.1, we present the results which we obtained
397 when fitting the model to the pre-processed data at the feed-in points. In
398 Section 4.2, the performance of the prediction model using various validation
399 scores are studied. Section 4.3 quantifies the economic value of the prediction
400 model based on the value score. In Section 4.4, the performance of the copula
401 model are compared to the ones of a quantile regression model. Finally, we
402 check the performance of the fitted models for communities in Section 4.5.

403 4.1. Fitted model characteristics

404 The idea of the model fitting procedure was described in Section 3.2. To
405 illustrate the results, which we obtained for the fitted model, we consider an
406 exemplary feed-in point ℓ_0 . The marginal distributions of GHI forecasts and
407 the copula parameter are fitted for each hour of the day and forecast lead time
408 separately, whereas the marginal distributions of solar power supply are fitted
409 for each hour of the day. Both datasets are spanning over May, June and July
410 of the years 2015 and 2016, 11-12 UTC. Thus, each of the datasets has about
411 180 timestamps.

412 4.1.1. Fitted marginal densities

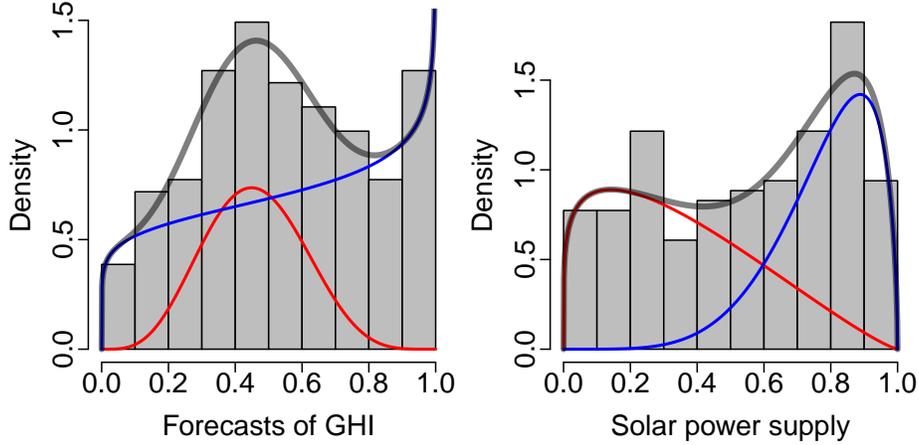


Figure 4: Histograms of rescaled GHI forecasts and rescaled solar power supply with fitted mixed beta density (grey) and weighted component distributions (blue and red) for feed-in point ℓ_0 .

413 The histograms in Figure 4 visualize the empirical distribution of GHI fore-
 414 casts and solar power supply for the feed-in point ℓ_0 . For data of length n the
 415 k equally distant bins of both histograms are determined by the *Sturges' rule*,
 416 see Scott (2011), i.e.,

$$417 \quad k = \lceil 1 + \log_2(n) \rceil. \quad (17)$$

418 If we sum up both suitably weighted curves we get the grey curve in Figure 4,
 419 which is the fitted mixed beta density.

420 From visual inspection of the empirical marginal distributions of the random
 421 variables R and S , we observe two modes in the histograms. We see that
 422 the fitted marginal densities approximate the shapes of both histograms quite
 423 accurately. Therefore, from a qualitative point of view a mixed distribution
 424 should be applied.

425 A quantitative assessment is also provided with the computation of the
 426 *Akaike information criterion* (AIC) for five different distribution types (for each
 427 feed-in point), see Figure 5. The AIC compares the fit for different distribution

428 types while penalizing choices with a larger number of fitted parameters. It is
 429 defined by

$$430 \quad AIC = 2k - 2\log(L), \quad (18)$$

431 where k is the number of fitted parameters. The smaller the AIC the better
 432 the fits of the distribution. Figure 5 shows that the mixed beta distribution fits
 433 the GHI and solar power supply better than the other considered distribution
 434 types.

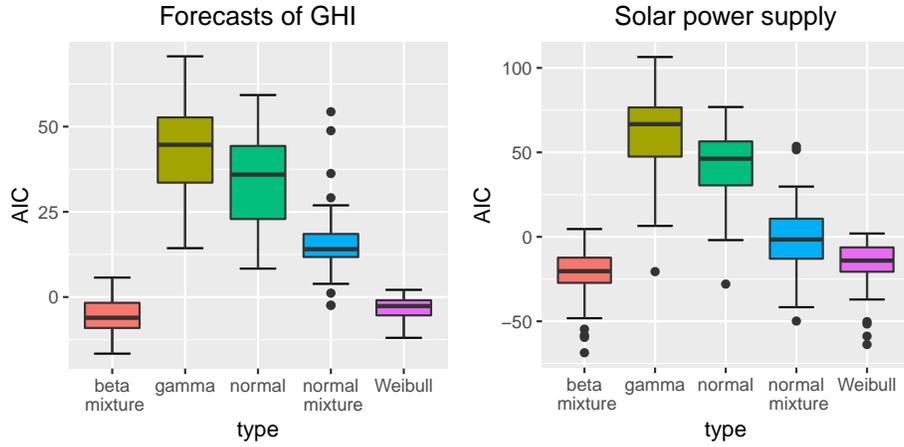


Figure 5: AIC computed for different distribution types fitted to GHI forecasts and solar power supply.

435 4.1.2. Variation in the marginal densities

436 Let R_1 and R_2 denote random variables with the beta densities f_1 and f_2 ,
 437 respectively, considered in Eq. (5), likewise for S_1 and S_2 . Thus, the random
 438 variables R_1 and R_2 are related to the two underlying beta densities of the
 439 random variable R , and S_1 and S_2 to the ones of the random variable S . For
 440 all feed-in points we computed the expectation and variance of the random
 441 variables R_i and S_i using the equations

$$442 \quad \mathbb{E}(Z) = \frac{a_i}{a_i + b_i}, \quad (19)$$

$$443 \quad \text{Var}(Z) = \frac{a_i b_i}{(a_i + b_i + 1)(a_i + b_i)^2}, \quad (20)$$

445 where Z is a random variable designating either R_i or S_i , and $a_i, b_i > 0$ are the
 446 shape parameters introduced in Section 3.1.1. The computed expectations and
 447 variances over the entire ROI are summarized in Figure 6.

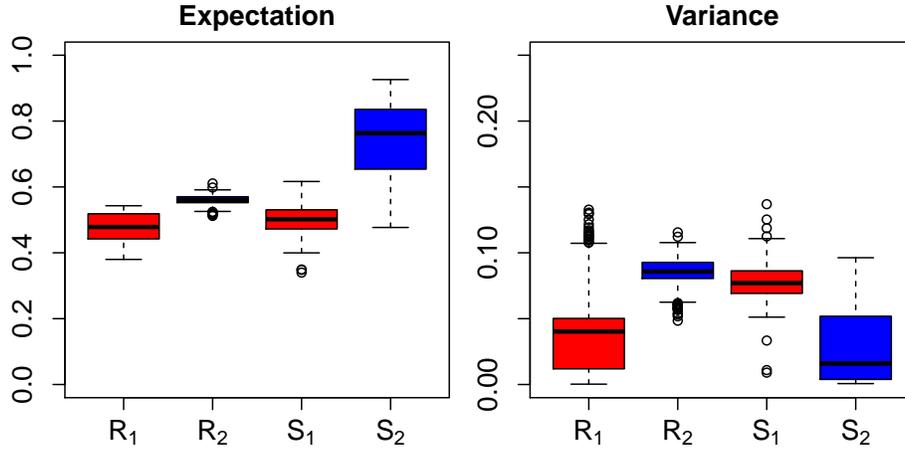


Figure 6: Expectation and variance of the fitted component beta distributions.

448 Note that the width of most of the boxplots is quite large. This indicates
 449 that applying averages of the density parameters over all feed-in points might
 450 introduce inaccuracies into the model. Thus, the fitting of the density param-
 451 eters should be done for each feed-in point separately as it was already mentioned
 452 in Section 3.2.1.

453 4.1.3. Fitted copula and two-dimensional density function

454 As described in Section 3.2.2, we computed the copula parameter θ and
 455 the maximum of the log-likelihood function for each copula type and feed-in
 456 point separately. The results are depicted in Table 2 for feed-in point ℓ_0 . Table
 457 2 shows that the Frank copula has the largest maximum of the log-likelihood
 458 function. Thus, for ℓ_0 the Frank copula performs better than the other copula
 459 types considered in Table 1.

copula type	Clayton	Frank	Gumbel	Joe
θ	3.50	7.26	2.94	3.21
$\log L(\theta)$	70.63	93.87	69.37	49.08

Table 2: Estimates of the parameter θ and maxima of the log-likelihood function for the exemplary feed-in point ℓ_0 .

460 Altogether, for 166 feed-in points out of the 168 considered feed-in points,
461 see Section 2.1, the Frank copula was determined as best fit while for only two
462 feed-in points a different copula type performed better. Therefore, in order to
463 reduce computational complexity, we decided to apply the Frank copula to all
464 feed-in points.

465 The copula parameter θ expresses how strongly the random variables R and
466 S are correlated. Note that for each feed-in point ℓ , *Spearman's rank correlation*
467 *coefficient* $rc(\ell)$ of R and S can be written as

$$468 \quad rc(\ell) = 12 \int_0^1 \int_0^1 C_\ell(u, v) du dv - 3, \quad (21)$$

469 where C_ℓ is the copula function fitted to the data at ℓ , see Schweizer and Wolff
470 (1981). In Figure 7 the local Spearman's rank coefficients are visualized. High
471 correlations, but also spatial disparities, can be observed.

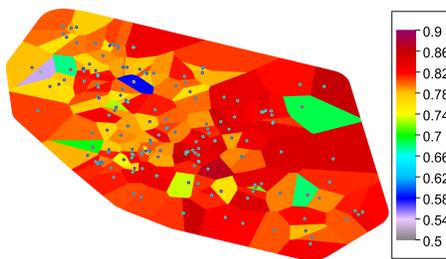


Figure 7: Local Spearman's rank coefficients computed based on fitted copula parameters.

472 Figure 8 visualizes the fitted joint density $f_{(R,S)}$ and the conditional densities
473 $f_{S|R}$ of solar power supply given a GHI forecast r at the exemplary feed-in point
474 ℓ_0 . As expected, with increasing GHI forecasts, the conditional densities assign

475 higher values to larger amounts of solar power supply.

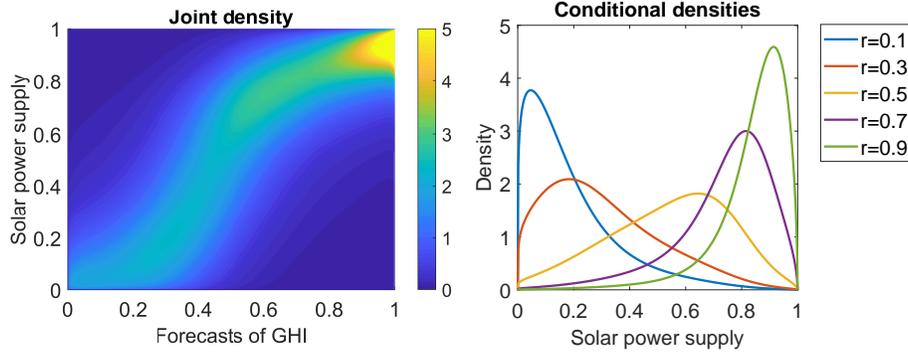


Figure 8: The joint density $f_{R,S}$ and conditional densities $f_{S|R}$ given certain GHI forecasts r for feed-in point ℓ_0 .

476 4.2. Validation based on scores

477 The proposed copula model was validated using various validation scores,
 478 such as the Brier skill score or the continuous ranked probability score, widely
 479 applied in the field of weather forecasting, see Wilks (2011). Recall that for
 480 model fitting we used datasets spanning over May, June and July of the years
 481 2015 and 2016. For validation we consider the data in the time frame May, June
 482 and July of the year 2017.

483 4.2.1. Notation

484 In the following we consider GHI forecasts r_i and corresponding measure-
 485 ments of solar power supply s_i where the index i belongs to the validation
 486 set $I_{val} = \{1, \dots, n\}$. We compute the conditional level-crossing probabilities
 487 $p_i(v) = P(S \geq v | R = r_i)$ for solar power supply exceeding the level $v \in [0, 1]$.
 488 The occurrence of the event that solar power supply s_i exceeds the level v is
 489 denoted by $o_i(v) = I(v, s_i)$, where I is the indicator function defined as

$$490 \quad I(v, s) = \begin{cases} 1, & \text{if } s \geq v, \\ 0, & \text{if } s < v. \end{cases} \quad (22)$$

491 Thus, for each level v and index $i \in I_{val}$, we consider the pair $(p_i(v), o_i(v))$ of
 492 the probability and occurrence of the event $\{S \geq v\}$ for a given feed-in point,
 493 hour of the day and forecast time.

494 4.2.2. Bias

495 The *bias* of our prediction model is defined as

$$496 \quad bias(v) = \frac{1}{n} \sum_{i=1}^n (p_i(v) - o_i(v)). \quad (23)$$

497 Note that the quantity $bias(v)$ takes values in $[-1, 1]$ for each $v \in [0, 1]$. In the
 498 ideal case, a good prediction model is unbiased, i.e. , $bias(v)$ is equal to zero.

499 However, an unbiased prediction model does not necessarily generate useful
 500 predictions. For instance, if all probabilities are equal to the relative frequency
 501 of the occurrences of the considered event, then the predictions are the same for
 502 all days. Such a prediction model would be unbiased, but it holds no information
 503 in regard to short term changes.

504 4.2.3. Brier score

505 Another measure of accuracy is the *Brier score* defined by

$$506 \quad bs(v) = \frac{1}{n} \sum_{i=1}^n (p_i(v) - o_i(v))^2. \quad (24)$$

507 The Brier score takes its values in the interval $[0, 1]$. It represents the mean
 508 squared differences between our predictions and the actual events. Thus, the
 509 Brier score of a good prediction model should be near zero.

510 The Brier score encompasses many important characteristics of a prediction
 511 model. Using the algebraic decomposition of the Brier score, more information
 512 on the model performance can be obtained (see Wilks (2011)). For a level v ,
 513 the Brier score $bs(v)$ can be expressed in terms of the reliability rel , resolution
 514 res and uncertainty unc , where

$$515 \quad bs(v) = rel(v) - res(v) + unc(v). \quad (25)$$

516 The definition of the reliability and resolution requires to partition the unit
 517 interval $[0,1]$ into sub-intervals B_1, \dots, B_J . Then, each sub-interval B_j contains

518 n_j values of forecasts $p_i(v)$ associated to the indicators of the occurring events
 519 $o_i(v)$. Furthermore, by $\bar{p}_j(v)$ and $\bar{o}_j(v)$ we denote the mean of the probabilities
 520 and the mean of the number of observations for each partition component B_j ,
 521 i.e.

$$522 \quad \bar{p}_j(v) = \frac{1}{n_j} \sum_{p_k(v) \in B_j} p_k(v), \quad (26)$$

$$523 \quad \bar{o}_j(v) = \frac{1}{n_j} \sum_{o_k(v) \in B_j} o_k(v). \quad (27)$$

524
 525 Moreover, by $\bar{o}(v)$ we denote the *climatological mean* for all observations, i.e.

$$526 \quad \bar{o}(v) = \frac{1}{n} \sum_{i=1}^n o_i(v). \quad (28)$$

527 Then, the *reliability* is defined as

$$528 \quad rel(v) = \frac{1}{n} \sum_{j=1}^J n_j (\bar{p}_j(v) - \bar{o}_j(v))^2. \quad (29)$$

529 A small reliability is typical for a well-calibrated prediction model. Besides, the
 530 *resolution* is defined as

$$531 \quad res(v) = \frac{1}{n} \sum_{j=1}^J n_j (\bar{o}_j(v) - \bar{o}(v))^2. \quad (30)$$

532 It measures how large the means of the number of observations for each sub-
 533 interval differ from the climatological mean for all observations. Higher resolu-
 534 tion means that the prediction model is able to distinguish between situations
 535 with different frequencies of occurrence. Last but not least, the *uncertainty* is
 536 given by

$$537 \quad unc(v) = \bar{o}(v)(1 - \bar{o}(v)), \quad (31)$$

538 which summarizes the variability of the observed events. The uncertainty does
 539 not depend on the prediction model. A low uncertainty value means that the
 540 observed events happen either with high or low frequency.

541 4.2.4. Brier skill score

542 The Brier score considered in Section 4.2.3 does not allow for a direct quan-
 543 titative comparison of the accuracy of prediction models as it depends on the

544 characteristics of the observed event. To draw a clear line between a good and
 545 a bad prediction model, the Brier skill score is used. The Brier skill score com-
 546 pares the Brier score $bs(v)$ of the prediction model with the Brier score $bs_r(v)$
 547 of some reference model. Thus, the *Brier skill score* is defined as

$$548 \quad bss_r(v) = 1 - \frac{bs(v)}{bs_r(v)}. \quad (32)$$

549 Note that $bss_r(v)$ takes its values in the interval $[-\infty, 1]$. If the reference model
 550 gives better results than the actually considered prediction model, the Brier
 551 skill score is negative and otherwise positive.

552 As reference model we consider the climatological model $\bar{o}(v)$ often used as
 553 a benchmark for weather forecast models. Since for the climatological model
 554 reliability and resolution are equal to zero, the Brier score of the climatological
 555 model is equal to the uncertainty unc leading to

$$556 \quad bss_c(v) = \frac{res(v) - rel(v)}{unc(v)}. \quad (33)$$

557 This implies that $bss_c < 0$ holds if the resolution is smaller than the reliability.
 558 This is clearly an undesirable outcome, as the resolution should be high and the
 559 reliability small. Therefore, prediction models with $bss_c < 0$ are commonly not
 560 considered.

561 4.2.5. Continuous ranked probability score

562 We consider the conditional cumulative distribution function of solar power
 563 supply $F_{S|R=r_i}$ given a GHI forecast r_i , and the corresponding measured solar
 564 power supply s_i . Then, we define the *continuous rank probability score* as

$$565 \quad crps_i = \int_{-\infty}^{\infty} [F_{S|R=r_i}(x) - F_i(x)]^2 dx \quad (34)$$

566 where F_i is the so-called cumulative-probability step function defined as

$$567 \quad F_i(x) = \begin{cases} 0 & \text{if } x \leq s_i, \\ 1 & \text{if } x > s_i. \end{cases} \quad (35)$$

568 Furthermore, we compute

$$569 \quad crps = \frac{1}{n} \sum_{i=1}^n crps_i \quad (36)$$

570 to quantify how concentrated around the corresponding observations are the
 571 computed conditional densities given GHI forecasts. As a general rule we can
 572 state: the lower the *crps*, the better.

573 *4.2.6. Validation for each feed-in point*

574 In this section we consider the forecast period 11-12 UTC and the forecast
 575 lead time of one hour. We compute the bias, Brier score and Brier skill score
 576 for each feed-in point separately. Furthermore, we compute the empirical cor-
 577 relation coefficient $\rho(v)$, see Eq. (2), of the observations $o_i(v)$ and probabilities
 578 $p_i(v)$.

579 The resulting validation scores are visualized in Figure 9 for the threshold
 580 $v = 0.8$. The biases and Brier scores are near zero for all feed-in points, whereas
 581 almost all computed Brier skill scores and empirical correlation coefficients are
 582 high. This indicates that the prediction model proposed in Section 3 works
 583 quite well regardless of the considered location.

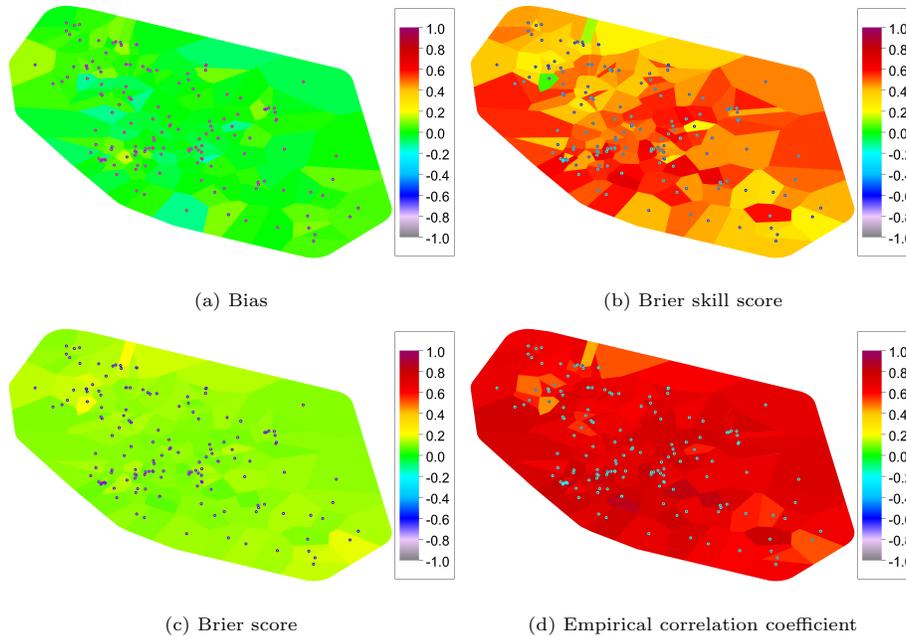


Figure 9: Validation scores computed at each feed-in point for the threshold $v = 0.8$

584 4.2.7. Validation for different lead times and hours of the day

585 It is commonly accepted that the longer the lead time is, the worse is the
 586 accuracy of forecasts. This effect is analyzed by merging the validation sets of
 587 all feed-in points and applying each validation score considered in Sections 4.2.2
 588 to 4.2.6 to the whole validation set. The analysis becomes then independent of
 589 the location and enables us to sum up the information to a single value for each
 590 score. Not surprisingly, Table 3 shows that our model yields better results for
 591 the shortest lead time of 1h than for longer lead times. However, the validation
 592 scores for longer forecast lead times are still highlighting good performance of
 593 the proposed model.

lead time (h)	<i>bias</i>	<i>bs</i>	<i>bss</i>	ρ	<i>rel</i>	<i>res</i>	<i>unc</i>	<i>crps</i>
1	0.006	0.124	0.437	0.661	0.001	0.097	0.221	0.088
4	-0.017	0.152	0.316	0.568	0.006	0.075	0.223	0.105
7	-0.022	0.151	0.315	0.565	0.004	0.073	0.221	0.108
10	-0.035	0.150	0.326	0.578	0.003	0.075	0.222	0.109
13	-0.037	0.141	0.363	0.608	0.003	0.083	0.221	0.106
16	-0.050	0.145	0.346	0.598	0.004	0.080	0.222	0.108
19	-0.069	0.151	0.327	0.590	0.006	0.079	0.224	0.110

Table 3: Validation scores of the combined validation sets for different forecast lead times and the threshold $v = 0.8$.

594 Finally, we applied our prediction model to GHI forecasts and solar power
 595 supply for different hours of the day. Table 4 shows very good validation scores
 596 for each 1h-period, regardless of the considered hour of the day. Regarding most
 597 scores we get slightly worse results for 13-14 UTC and 17-18 UTC, but even in
 598 these cases the brier skill score is clearly greater than zero and the bias is almost
 599 zero.

hour of the day	<i>bias</i>	<i>bs</i>	<i>bss</i>	ρ	<i>rel</i>	<i>res</i>	<i>unc</i>	<i>crps</i>
5-6	0.005	0.085	0.318	0.565	0.002	0.041	0.124	0.086
6-7	-0.019	0.125	0.415	0.646	0.002	0.090	0.214	0.092
7-8	-0.025	0.137	0.419	0.650	0.002	0.100	0.236	0.094
8-9	0.006	0.146	0.369	0.624	0.006	0.090	0.232	0.093
9-10	-0.002	0.155	0.308	0.570	0.005	0.072	0.224	0.097
10-11	-0.006	0.155	0.291	0.552	0.004	0.067	0.218	0.100
11-12	0.006	0.124	0.437	0.661	0.001	0.097	0.221	0.088
12-13	-0.026	0.147	0.347	0.593	0.002	0.078	0.225	0.098
13-14	-0.028	0.157	0.278	0.535	0.003	0.062	0.217	0.108
14-15	-0.026	0.126	0.397	0.633	0.001	0.084	0.209	0.097
15-16	-0.025	0.127	0.368	0.609	0.002	0.074	0.200	0.100
16-17	0.001	0.097	0.290	0.538	0.002	0.041	0.136	0.097
17-18	-0.060	0.104	0.318	0.601	0.007	0.053	0.153	0.168

Table 4: Validation scores of the combined validation sets for 1h forecast lead time, for different hours of the day (in UTC) and the threshold $v = 0.8$.

600 *4.3. Economic value of the forecast model*

601 Distribution network operators have high interest in the economic value of
602 prediction models. For that purpose, the value score VS is often used to analyze
603 the economic value of a forecast model compared to the climatological model,
604 see Wilks (2011).

605 The computation of the value score is based on the so-called cost-loss ratio.
606 In our case, the cost C is the cost of a curtailment which should be done if the
607 solar plant generates more solar power supply than the predefined threshold.
608 The loss L corresponds to the averaged economical damage caused by the over-
609 loading event when the predefined threshold is exceeded. Then, the *cost-loss*
610 *ratio* is given by the quotient $C/L \in [0, 1]$. It enables us to characterize various
611 possible scenarios for which the value score can be computed.

612 To formally define the value score, the relative joint frequencies $p_{0,0}, p_{1,0}, p_{0,1}$

and $p_{1,1}$ given in Table 5 have to be considered.

$p_{1,1} = \frac{\#\{i : p_i > C/L, o_i = 1\}}{m}$	$p_{1,0} = \frac{\#\{i : p_i > C/L, o_i = 0\}}{m}$
$p_{0,1} = \frac{\#\{i : p_i \leq C/L, o_i = 1\}}{m}$	$p_{0,0} = \frac{\#\{i : p_i \leq C/L, o_i = 0\}}{m}$

Table 5: Relative joint frequencies considered for the computation of the value score.

613

614 Then, the *value score* VS is defined as

$$VS = \begin{cases} \frac{(C/L)(p_{1,1} + p_{1,0} - 1) + p_{0,1}}{(C/L)(\bar{o} - 1)}, & \text{if } C/L < \bar{o}, \\ \frac{(C/L)(p_{1,1} + p_{1,0}) + p_{0,1} - \bar{o}}{\bar{o}((C/L) - 1)}, & \text{if } C/L > \bar{o}. \end{cases} \quad (37)$$

616 For more details regarding the interpretation of the value score, see Wilks (2011).

617 Following the approach of Wilks (2001), we compute the value score for the
618 cost/loss ratio values equal to $(p_i + p_{i+1})/2$, where $i \in \{1, \dots, m - 1\}$. The
619 evolution of the value score with respect to the cost-loss ratio is illustrated in
620 Figure 10. The value score is positive for all considered cost/loss ratios which
621 means that the proposed model is more valuable than the climatological model.
622 Moreover, the proposed model shows a greater economic utility for decision
623 making with value scores larger than 0.5 if $C/L \in [0.12, 0.55]$.

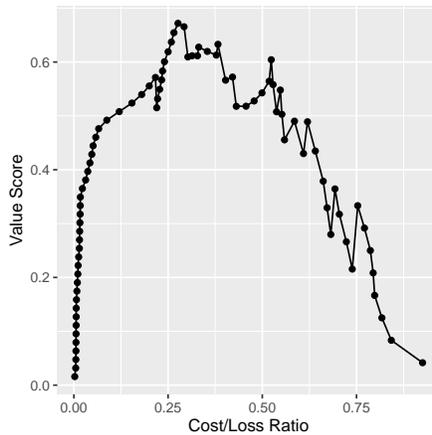


Figure 10: Value score with respect to cost-loss ratio for the exemplary feed-in point ℓ_0 .

624 *4.4. Comparison with the quantile regression model*

625 The accuracy of the copula model proposed in this paper is finally compared
 626 to the one of the quantile regression (QR) model. QR is a non-parametric ap-
 627 proach to estimate conditional quantiles of a random variable, called predictand,
 628 given some independent random variables, called predictors. Linear QR assumes
 629 a linear relationship between the conditional quantiles of the predictand and the
 630 predictors, see Davino et al. (2013) for more details.

631 We consider the time period 11 – 12 UTC and a forecast lead time of one
 632 hour. For each feed-in point, the solar power supply S is the predictand and
 633 the global horizontal irradiation R the unique predictor. Then, the conditional
 634 α -quantile $\hat{q}_\alpha(r) = \hat{b}_\alpha r$ is computed by minimizing the expression

$$635 \quad \sum_{i=1}^n \rho_\alpha(s_i - \hat{b}_\alpha r_i), \quad (38)$$

636 where the quantile loss function is defined as

$$637 \quad \rho_\alpha(u) = \begin{cases} \alpha u, & \text{if } u \geq 0, \\ (\alpha - 1)u, & \text{otherwise.} \end{cases} \quad (39)$$

638 By applying 101 linear quantile regressions, i.e. with $\alpha \in \{0, 1/100, \dots, 1\}$,
 639 we compute conditional α -quantiles for solar power supply given GHI forecasts.
 640 Next, for a threshold v we approximate the conditional level-crossing probability
 641 of S given the GHI forecast r under the quantile regression model by

$$642 \quad \pi(v) = 1 - \min_{\alpha} (\hat{q}_\alpha(r) - v). \quad (40)$$

643 The reliability diagrams are then computed for the copula model and the
 644 quantile regression. Note that the corresponding reliability diagrams are the
 645 pairs of points $(\bar{p}_j(v), \bar{o}_j(v))$ and $(\bar{\pi}_j(v), \bar{d}_j(v))$ for $j \in \{1, \dots, J\}$, where $\bar{p}_j(v)$
 646 and $\bar{\pi}_j(v)$ are the mean conditional level-crossing probabilities of S given R
 647 under the copula model and the QR model, respectively, for the sub-interval B_j
 648 (using the same notation as in Section 4.2.3).

649 For a perfectly reliable model, the difference between the mean conditional
 650 probability and the corresponding frequency of occurrence of the considered

651 event is zero, i.e., the reliability diagram coincides with the diagonal (grey line
 652 in Figure 11). Vice versa, the more the reliability diagram deviates from the
 653 diagonal, the less reliable the model, see Wilks (2011) for further details.

654 Figure 11 compares the reliability diagrams of the two models. The numbers
 655 of samples n_j in the sub-intervals B_j are visualized in bars below the diagrams.
 656 We observe that the copula model is more reliable than the quantile regression
 657 model which tends to underestimate the probability that the solar power supply
 658 exceeds the predefined threshold $v = 0.8$.

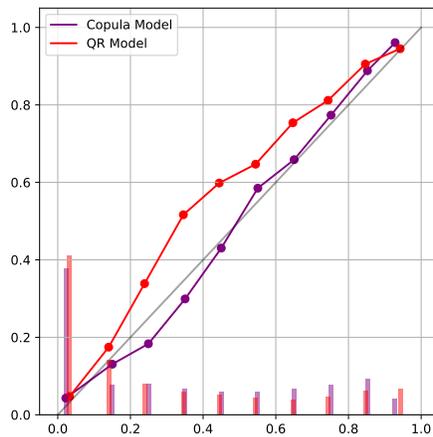


Figure 11: Reliability diagrams of the copula model (in violet) and the QR model (in red), for the predefined threshold $v = 0.8$.

659 4.5. Comparison between different hierarchy levels

660 Since overloading problems can occur at each hierarchy level of a distribution
 661 network, it is crucial that probabilistic predictions quantifying the risk of over-
 662 loading can be generated for several hierarchy levels. In fact, the consequences of
 663 a critical event are usually even larger the higher the hierarchy level is, where it
 664 occurs. Indeed, in such a case the solar power supply of a whole region might be
 665 interrupted. Consequently, it is convenient to have a flexible prediction model,
 666 which can be applied to different hierarchy levels of a distribution network.

667 Using the methods stated in Section 3.1 we compute and visualize condi-
 668 tional level-crossing probabilities for a certain threshold at feed-in points and

669 communities. It turns out that there are less local disparities for communities
 670 than for feed-in points. In fact the aggregation of solar power supply causes a
 671 smoothing effect on the conditional probabilities, see Figure 12.

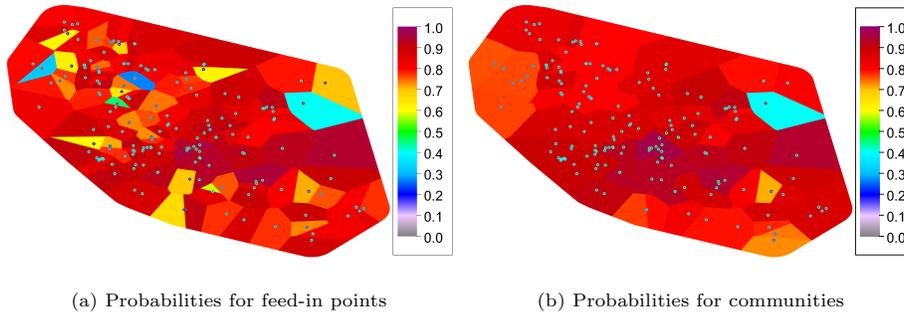


Figure 12: Conditional level-crossing probabilities for solar power supply exceeding the threshold of 0.8 for May 10, 2017, 11-12 UTC.

672 Table 6 clearly shows the flexibility of our prediction model by comparing the
 673 averaged validation scores computed for communities with the averaged scores
 674 we get for individual feed-in points. Some scores for communities are even
 675 slightly better than for feed-in points. This stands to reason, since averaging
 676 over a certain number of feed-in points eliminates noise.

hierarchy level	<i>bias</i>	<i>bs</i>	<i>bss</i>	ρ	<i>rel</i>	<i>res</i>	<i>unc</i>	<i>crps</i>
feed-in points	0.006	0.124	0.437	0.661	0.001	0.097	0.221	0.088
communities	0.004	0.113	0.491	0.701	0.002	0.110	0.223	0.080

Table 6: Validation scores of the combined validation sets for feed-in points and communities with lead time one hour and forecast period 11-12 UTC.

677 5. Conclusion

678 In this paper a probabilistic prediction model to quantify the risk of overload-
 679 ing at feed-in points was proposed. The model is based on hourly deterministic
 680 GHI forecasts and applies copula theory to compute the joint distribution of GHI
 681 forecast and solar power supply for each feed-in point. Based on marginal and

682 joint distributions, conditional probabilities for solar power supply exceeding a
683 predefined threshold are computed.

684 The model was validated using prediction scores, such as bias, Brier score
685 (decomposed into reliability, resolution and uncertainty), Brier skill score, con-
686 tinuous rank probability score and the empirical correlation coefficient. The
687 scores were validated for forecast lead times ranging from 1 to 19 hours and 1h
688 periods of the day ranging from 5 to 18 UTC. These validations showed a high
689 accuracy of the proposed copula-based model, regardless of the considered hour
690 of the day or forecast lead time. Moreover, a comparison of the copula model
691 reliability diagram with the one of the quantile regression model emphasized
692 higher reliability than a state-of-the-art model. Besides, we also showed that
693 the model can be applied to higher hierarchy levels in the distribution network,
694 such as communities. Finally, the value score was computed to quantify the
695 economic value of the proposed model for an exemplary feed-in point. The cop-
696 ula model demonstrated a higher economic utility than the climatological model
697 for the selected feed-in point, regardless of the cost-loss ratio values. Particu-
698 larly, the economic utility was outstanding for cost-loss ratios in the interval
699 $[0.12, 0.55]$.

700 The fitting and the validation of the copula model have been undertaken
701 on three typical months with high GHI in Germany (May, June and July). If
702 more than three months have to be considered, it may be necessary to split the
703 fitting and validation period in order to capture information relevant to seasonal
704 variation. The way to split the data, if necessary, is not straightforward and
705 may have to be undertaken by trial-and-error.

706 Besides, it may also be possible to integrate other factors, except GHI, that
707 are correlated to the solar power supply. This will be investigated in a forth-
708 coming paper using more advanced tools of copula theory such as nesting or
709 vine copulas (Joe, 2014).

710 **Acknowledgments**

711 We would like to thank “Bundesministerium für Bildung und Forschung”
712 (BMBF) for financially supporting this research project (grant 05M18VUB).

713 **References**

714 Alessandrini, S., Delle Monache, L., Sperati, S., Cervone, G., 2015. An analog
715 ensemble for short-term probabilistic solar power forecast. *Applied Energy*
716 157, 95–110.

717 Almeida, M.P., Perpignan, O., Narvarte, L., 2015. PV power forecast using a
718 nonparametric PV model. *Solar Energy* 115, 354–368.

719 Antonanzas, J., Osorio, N., Escobar, R., Urraca, R., Martinez-de Pison, F.,
720 Antonanzas-Torres, F., 2016. Review of photovoltaic power forecasting. *Solar*
721 *Energy* 136, 78–111.

722 Bacher, P., Madsen, H., Nielsen, H.A., 2009. Online short-term solar power
723 forecasting. *Solar Energy* 83, 1772–1783.

724 Baldauf, M., Seifert, A., Förstner, J., Majewski, D., Raschendorfer, M., Rein-
725 hardt, T., 2011. Operational convective-scale numerical weather prediction
726 with the COSMO model: Description and sensitivities. *Monthly Weather*
727 *Review* 139, 3887–3905.

728 Bessa, R.J., Trindade, A., Silva, C.S., Miranda, V., 2015. Probabilistic solar
729 power forecasting in smart grids using distributed information. *International*
730 *Journal of Electrical Power & Energy Systems* 72, 16–23.

731 Bird, L., Cochran, J., Wang, X., 2014. Wind and solar energy curtailment:
732 Experience and practices in the United States. Technical Report. National
733 Renewable Energy Laboratory.

734 Coiffier, J., 2011. *Fundamentals of Numerical Weather Prediction*. Cambridge
735 University Press.

- 736 Council of European Energy Regulators, 2017. Ceer report on
737 power losses. URL: [https://www.ceer.eu/documents/104400/-/-/
738 09ecee88-e877-3305-6767-e75404637087](https://www.ceer.eu/documents/104400/-/-/09ecee88-e877-3305-6767-e75404637087).
- 739 Davino, C., Furno, M., Vistocco, D., 2013. Quantile Regression: Theory and
740 Applications. volume 988. J. Wiley & Sons.
- 741 Dempster, A., Laird, N., Rubin, D., 1977. Maximum likelihood from incomplete
742 data via the EM algorithm. *Journal of the Royal Statistical Society: Series
743 B (Methodological)* 39, 1–22.
- 744 Dickert, J., Hable, M., Schegner, P., 2009. Energy loss estimation in distribution
745 networks for planning purposes, in: *IEEE Bucharest PowerTech*, pp. 1–6.
- 746 Durante, F., Sempi, C., 2015. *Principles of Copula Theory*. Chapman and
747 Hall/CRC.
- 748 Golestaneh, F., Gooi, H.B., 2017. Multivariate prediction intervals for photo-
749 voltaic power generation, in: *2017 IEEE Innovative Smart Grid Technologies-
750 Asia (ISGT-Asia)*, IEEE. pp. 1–5.
- 751 Golestaneh, F., Gooi, H.B., Pinson, P., 2016a. Generation and evaluation of
752 space–time trajectories of photovoltaic power. *Applied Energy* 176, 80–91.
- 753 Golestaneh, F., Pinson, P., Gooi, H.B., 2016b. Very short-term nonparametric
754 probabilistic forecasting of renewable energy generation—with application to
755 solar energy. *IEEE Transactions on Power Systems* 31, 3850–3863.
- 756 Hastie, T., Tibshirani, R., Friedman, J., 2009. *The Elements of Statistical
757 Learning: Data Mining, Inference and Prediction*. Springer.
- 758 Heinemann, D., Lorenz, E., Girodo, M., 2006. *Solar Energy Resource Man-
759 agement for Electricity Generation from Local Level to Global Scale*. Nova
760 Science Publishers.
- 761 Hess, R., Glashoff, J., Reichert, B., 2015. The Ensemble-MOS of Deutscher
762 Wetterdienst, in: *EMS Annual Meeting Abstracts*, Sofia.

- 763 Huang, J., Perry, M., 2016. A semi-empirical approach using gradient boosting
764 and k-nearest neighbors regression for gefcom2014 probabilistic solar power
765 forecasting. *International Journal of Forecasting* 32, 1081–1086.
- 766 Hämmerlin, G., Hoffmann, K., 2012. *Numerical Mathematics*. Springer.
- 767 Jobson, J., 1991. *Applied Multivariate Data Analysis: Regression and Experi-*
768 *mental Design*. Springer.
- 769 Joe, H., 2014. *Dependence Modeling with Copulas*. Chapman and Hall/CRC.
- 770 Joe, H., Xu, J., 1996. The estimation method of inference functions for
771 margins for multivariate models. URL: <https://open.library.ubc.ca/collections/facultyresearchandpublications/52383/items/1.0225985>.
772
773
- 774 Kaldellis, J., Kapsali, M., Kavadias, K., 2014. Temperature and wind speed
775 impact on the efficiency of PV installations. Experience obtained from outdoor
776 measurements in greece. *Renewable Energy* 66, 612–624.
- 777 Karimi, M., Mokhlis, H., Naidu, K., Uddin, S., Bakar, A., 2016. Photovoltaic
778 penetration issues and impacts in distribution network - A review. *Renewable*
779 *and Sustainable Energy Reviews* 53, 594–605.
- 780 Lauret, P., David, M., Pedro, H., 2017. Probabilistic solar forecasting using
781 quantile regression models. *Energies* 10, 1591.
- 782 Leisch, F., 2004. A general framework for finite mixture models and latent class
783 regression in R. *Journal of Statistical Software*, 11 (8) , 1–18.
- 784 Lu, Q., Hu, W., Min, Y., Yuan, F., Gao, Z., 2014. Wind power uncertainty
785 modeling considering spatial dependence based on pair-copula theory, in: *PES*
786 *General Meeting— Conference & Exposition*, IEEE. pp. 1–5.
- 787 Massidda, L., Marrocu, M., 2018. Quantile regression post-processing of weather
788 forecast for short-term solar power probabilistic forecasting. *Energies* 11,
789 1763.

790 Mekhilef, S., Saidur, R., Kamalisarvestani, M., 2012. Effect of dust, humidity
791 and air velocity on efficiency of photovoltaic cells. *Renewable and Sustainable*
792 *Energy Reviews* 16, 2920–2925.

793 Nelsen, R., 2006. *An Introduction to Copulas*. Springer.

794 Panamtaash, H., Zhou, Q., Hong, T., Qu, Z., Davis, K.O., 2020. A copula-based
795 Bayesian method for probabilistic solar power forecasting. *Solar Energy* 196,
796 336–345.

797 Papaefthymiou, G., Kurowicka, D., 2009. Using copulas for modeling stochastic
798 dependence in power system uncertainty analysis. *IEEE Transactions on*
799 *Power Systems* 24, 40 – 49.

800 Schweizer, B., Wolff, E., 1981. On nonparametric measures of dependence for
801 random variables. *The Annals of Statistics* 9, 879–885.

802 Scott, D.W., 2011. *Sturges’ and Scott’s Rules*. Springer. pp. 1563–1566.
803 URL: [https://doi.org/10.1007/](https://doi.org/10.1007/978-3-642-04898-2_578)
804 [978-3-642-04898-2_578](https://doi.org/10.1007/978-3-642-04898-2_578), doi:10.1007/
[978-3-642-04898-2_578](https://doi.org/10.1007/978-3-642-04898-2_578).

805 SolarPower Europe, 2017. *Global market outlook 2018-2022*. URL:
806 [http://www.solarpowereurope.org/wp-content/uploads/2018/09/](http://www.solarpowereurope.org/wp-content/uploads/2018/09/Global-Market-Outlook-2018-2022.pdf)
807 [Global-Market-Outlook-2018-2022.pdf](http://www.solarpowereurope.org/wp-content/uploads/2018/09/Global-Market-Outlook-2018-2022.pdf).

808 Vale, P., 2015. *Energy assessment of photovoltaic conversion systems*. Technical
809 Report. Instituto Superior Técnico.

810 Wang, Y., Infield, D., Stephen, B., Galloway, S., 2014. Copula-based model for
811 wind turbine power curve outlier rejection. *Wind Energy* 17, 1677–1688.

812 Wilks, D., 2001. A skill score based on economic value for probability forecasts.
813 *Meteorological Applications* 8, 209–219.

814 Wilks, D.S., 2011. *Statistical Methods in the Atmospheric Sciences*. Academic
815 Press.

- 816 Zamo, M., Mestre, O., Arbogast, P., Pannekoucke, O., 2014. A benchmark
817 of statistical regression methods for short-term forecasting of photovoltaic
818 electricity production. Part II: probabilistic forecast of daily production. *Solar*
819 *Energy* 105, 804–816.
- 820 Zhang, B., Dehghanian, P., Kezunovic, M., 2016. Spatial-temporal solar power
821 forecast through use of gaussian conditional random fields, in: *IEEE Power*
822 *and Energy Society General Meeting (PESGM)*, pp. 1–5.