

Document downloaded from:

<http://hdl.handle.net/10251/40653>

This paper must be cited as:

Zarzo Castelló, M.; Martí Pérez, PC. (2011). Modeling the variability of solar radiation data among weather stations by means of principal components analysis. *Applied Energy*. 88(8):2775-2784. doi:10.1016/j.apenergy.2011.01.070.



The final publication is available at

<http://dx.doi.org/10.1016/j.apenergy.2011.01.070>

Copyright Elsevier

1 **Modeling the variability of solar radiation data among weather**  
2 **stations by means of principal components analysis**

3

4 Manuel Zarzo <sup>a</sup>, Pau Martí <sup>b\*</sup>

5

6 <sup>a</sup>Departamento de Estadística e Investigación Operativa Aplicadas y Calidad.

7 Universidad Politécnica de Valencia. Camino de Vera, s/n. 46022 Valencia. Spain.

8

9 <sup>b</sup>Departamento de Ingeniería Rural y Agroalimentaria. Universidad Politécnica de

10 Valencia. Camino de Vera, s/n. 46022 Valencia. Spain.

11

12

13 \*Corresponding author. Tel.: +34 963877490; Fax.: +34 963877499.

14 E-mail addresses: mazarcas@eio.upv.es (M. Zarzo), paumarpe@hotmail.es  
15 (P. Martí).

16

17

18

19 ***Abstract***

20 Measurements of global terrestrial solar radiation ( $R_s$ ) are commonly recorded in  
21 meteorological stations. Daily variability of  $R_s$  has to be taken into account for the  
22 design of photovoltaic systems and energy efficient buildings. Principal components  
23 analysis (PCA) was applied to  $R_s$  data recorded at 30 stations in the Mediterranean  
24 coast of Spain. Due to equipment failures and site operation problems, time series of  $R_s$   
25 often present data gaps or discontinuities. The PCA approach copes with this problem  
26 and allows estimation of present and past values by taking advantage of  $R_s$  records from  
27 nearby stations. The gap infilling performance of this methodology is compared with  
28 alternative conventional approaches. A new method was also developed for  $R_s$   
29 estimation if previous measurements are not available. Four principal components  
30 explain 66% of the data variability with respect to the average trajectory. By means of  
31 multiple linear regression, it was found that this variability can be fitted according to the  
32 latitude, longitude and altitude of the station where data were recorded from. Additional  
33 geographical or climatic variables did not increase the predictive goodness-of-fit. The  
34 resulting models allow the estimation of daily  $R_s$  values at any location in the area  
35 under study. The proposed methodology for estimating  $R_s$  based on geographical  
36 parameters would be of interest to design solar energy systems and to select their best  
37 location.

38

39

40 **Keywords:** solar radiation, missing data estimation, PCA, multivariate statistical  
41 monitoring

42

43

44 **Introduction**

45 Solar radiation plays a key role in evaporation, plant photosynthesis or crop growth and  
46 productivity [1-4]. In architecture, accurate estimates of long-term global solar radiation  
47 are required for the design and development of energy efficient buildings [5,6]. Solar  
48 radiation is also essential in biophysical models for risk assessment of forest fires, in  
49 hydrological simulation models of natural processes [7], in environmental and  
50 agrometeorological research as well as in atmospheric physics [8]. Nonetheless, the  
51 major interest of measuring solar radiation is for the simulation and design of solar  
52 energy systems [9].

53

54 The most common solar radiation measurements recorded in meteorological stations  
55 correspond to total radiation on a horizontal surface,  $R_s$ , also called global terrestrial  
56 solar radiation, which is normally given on an hourly or daily basis [6]. These data are  
57 required for the design of photovoltaic applications in remote or isolated areas where no  
58 connection to an electrical supply grid is available, for instance in rural or mountainous  
59 areas, natural parks, small islands and developing countries in general [10-14]. In these  
60 places, the daily  $R_s$  variability has to be taken into account for the design of  
61 photovoltaic systems in order to guarantee enough power generation for the essential  
62 electric devices [15]. In developed countries, solar energy systems are frequently  
63 implemented in buildings for water heating or electric power supply. The design of  
64 these systems is also based on  $R_s$  measurements. However, in many applications of  
65 solar energy, especially in the aforementioned isolated areas, projects are not supported  
66 by the required  $R_s$  data at the place of interest. Actually, solar radiation is measured at  
67 relatively few weather stations in comparison to other variables such as temperature or

68 relative humidity (RH). This is generally due to the high cost, maintenance and  
69 calibration requirements of the measuring equipment [9,10,12,16,17].

70

71 Although suitable weather records have become more and more available in recent  
72 years, data reliability and quality is another problem. Even in automatic stations where  
73 solar radiation is measured, data records often lie outside the expected range [3-5] and  
74 are erroneous because of sensor calibration problems. According to Muneer et al. [9],  
75 another cause of errors is site operation problems such as instrument proximity to  
76 shading elements, electrical and magnetic fields, weather elements as well as bird or  
77 insect activity. Erroneous data need to be discarded, and equipment failures also cause  
78 missing data. As a result, time series of  $R_s$  often present data gaps or discontinuities.

79

80 One alternative to cope with the lack of accurate  $R_s$  measurements is to use  
81 mathematical predictive models relying on climatic inputs. Several empirical, numerical  
82 and physically-based models have been proposed for  $R_s$  estimation based on different  
83 input combinations. They differ in sophistication from simple empirical equations based  
84 on common climatic data to more complex numerical models involving high  
85 computational costs and relying on numerous inputs. The most frequent inputs are  
86 sunshine duration, extraterrestrial radiation, mean temperature, maximum temperature,  
87 soil temperature, RH, number of rainy days, altitude, latitude, total precipitation,  
88 cloudiness, and evaporation [16]. Among the simplest methods for estimating solar  
89 radiation data, Hargreaves and Samani [18], Bristow and Campbell [19] and Allen [20]  
90 propose equations relying on maximum and minimum temperatures as well as on  
91 extraterrestrial radiation. These approaches, modified by other authors [17], take into  
92 account implicitly the geographical information of the studied locations by including

93 theoretical extraterrestrial radiation values based on latitude, day of the year, sunset  
94 hour angle and relative distance earth-sun.

95  
96 As an alternative to conventional approaches, artificial neural networks (ANNs) have  
97 been successfully applied for solar radiation estimation [1,2,5,8,10-13,21-31].

98 Techniques based on artificial intelligence have also been proposed particularly for  
99 isolated areas [14]. However, only a small part of these works present models fed by  
100 few easily measurable inputs such as temperature and/or RH records [11,28,29].

101  
102 The development of  $R_s$  estimation methods not relying on local climatic records turns  
103 into a task of great relevance because even the simplest climatic parameters are not  
104 available in many cases given the limited number of available automatic weather  
105 stations. One approach hardly tackled in literature would be to develop models relying  
106 exclusively on exogenous  $R_s$  inputs from nearby locations with similar climatic  
107 conditions. These models are of relevant interest given the aforementioned ubiquitous  
108 problems such as data scarcity, equipment failures, maintenance and calibration as well  
109 as physical and biological constraints.

110  
111 Data of  $R_s$  recorded daily at different stations can be regarded as a multivariate time  
112 series. Such type of data is common in the monitoring and control of industrial  
113 processes. For example, chemical reactors are usually monitored by means of electronic  
114 sensors that record the temperature at different points of the process. In this context of  
115 multivariate statistical process control (MSPC), principal components analysis (PCA) is  
116 a useful technique for process monitoring and diagnosis because it allows data  
117 estimation in case of faulty sensors [32]. PCA is one of the multivariate techniques

118 wider spread [33]. In a recent study, PCA was used for modeling the spatial data  
119 variability from a set of RH sensors located at different positions [34]. The same PCA  
120 approach was applied here for the analysis of  $R_s$  values measured daily at 30 weather  
121 stations. PCA copes with gap infilling by taking advantage of  $R_s$  records from nearby  
122 stations. Multiple linear regression (MLR) was used to identify which geographical or  
123 climatic parameters are the ones that best explain the differences in  $R_s$  measurements  
124 recorded at the 30 stations. Once identified the key variables, a new methodology is  
125 proposed to estimate  $R_s$  when, apart from exogenous measurements, these parameters  
126 are also available.

127

## 128 **Materials and methods**

129

### 130 1. Data characterization

131 The database analyzed here consisted of daily  $R_s$  values from 30 weather stations  
132 located on the Mediterranean coast of Spain (Table 1). Data were obtained from the  
133 Valencian Institute of Agricultural Research (IVIA). The dataset was structured as a  
134 matrix of 30 rows (stations) by 2920 variables (in columns). Each variable corresponds  
135 to one day in the 8-year period under study (January 2000 to December 2007). The  
136 original  $R_s$  series contained missing data. In order to assess different procedures for gap  
137 infilling, it is convenient to work with a complete data matrix. Thus, all variables in the  
138 initial  $R_s$  dataset containing missing data were discarded, resulting a complete matrix  
139 with 1203 variables. Fig. 1 displays the average  $R_s$  for these 1203 days. A clear periodic  
140 trend can be observed due to  $R_s$  annual seasonability.

141

142

**[FIGURE 1 NEAR HERE]**

143 2. PCA models

144 2.1. *PCA configuration and data pretreatment*

145 Principal components are directions of maximum data variance obtained as  
146 linear combinations of the original variables. The projections of observations (weather  
147 stations, in this case) over these directions are called scores. The variable containing  
148 these projections over the first principal component (PC1) is called score vector ( $\mathbf{t}_1$ ).  
149 Similarly,  $\mathbf{t}_2$  contains the projections over PC2, and so on. The contributions of  
150 variables in the formation of a given component are called loadings, being  $\mathbf{p}_1$  the  
151 loadings in the formation of PC1.

152

153 The software SIMCA-P 10.0 (Umetrics AB, Malmö, Sweden) was used to carry out all  
154 PCA models. It uses the NIPALS algorithm [35] which extracts components one by  
155 one. Given a matrix  $\mathbf{X}$ , this algorithm calculates  $\mathbf{t}_1$  and  $\mathbf{p}_1$ , resulting a residual matrix  
156  $\mathbf{E}_1$  (eq. 1). PC2 is obtained by applying the NIPALS algorithm to  $\mathbf{E}_1$ . Thus, PC2 is the  
157 direction that explains the maximum data variability of  $\mathbf{E}_1$  and remains orthogonal to  
158 PC1. Next, a residual matrix  $\mathbf{E}_2$  is calculated (eq. 1). This procedure can be conducted  
159 sequentially until  $\mathbf{E}=0$ .

160

$$\left. \begin{array}{l} \mathbf{E}_1 = \mathbf{X} - \mathbf{t}_1 \mathbf{p}_1^T \\ \mathbf{E}_2 = \mathbf{E}_1 - \mathbf{t}_2 \mathbf{p}_2^T \\ \dots \\ \mathbf{E}_k = \mathbf{E}_{k-1} - \mathbf{t}_k \mathbf{p}_k^T \end{array} \right\} \Rightarrow \mathbf{X} = \mathbf{t}_1 \mathbf{p}_1^T + \dots + \mathbf{t}_k \mathbf{p}_k^T + \mathbf{E}_{k-1} \quad (1)$$

161 Each row of the  $R_s$  dataset is a time series that reflects the evolution of  $R_s$  recorded at  
162 one station. In the context of MSPC, such time series is often called ‘trajectory’ because  
163 of the trend observed when the parameter is plotted versus time. If a new row is  
164 obtained by averaging the values of each column, it could be regarded as the mean  
165 trajectory. In order to highlight the relationships (i.e. similarities and dissimilarities)



166 among stations, data were mean-centered prior to PCA by subtracting the mean value of  
167 each column. As a result, the average of all centered variables becomes null. In the  
168 MSPC of batch chemical processes, the idea of subtracting the mean trajectory prior to  
169 PCA was first proposed by Nomikos and MacGregor [36]. The same methodology has  
170 also been successfully applied by other works [37,38]. Variables are also scaled to  
171 unitary variance prior to PCA when the variance among them is very different [38], but  
172 this is not the case here.

173

## 174 *2.2. Number of relevant PCs and outlier detection*

175 Different methods can be applied to decide how many components should be extracted  
176 (i.e. the value of  $k$  in eq. 1) for the purpose of modeling the systematic data variability  
177 of  $\mathbf{X}$  [33]. Further components not calculated are included in a matrix of errors ( $\mathbf{E}_{k-1}$  in  
178 eq. 1) that is assumed to account for random variation. One criterion implemented in the  
179 software SIMCA-P 10.0 is cross-validation [39]. It considers that one PC does not  
180 provide relevant information if it changes significantly when several observations are  
181 randomly removed.

182

183         Applying PCA to all  $R_s$  data assumes that the relationships among stations are  
184 basically maintained all the year round. In order to test this hypothesis, the  $R_s$  dataset  
185 was split in two subsets of about equal size, one containing those variables with an  
186 average value higher than 200 and another one containing the remaining variables.  
187 These subsets will be referred to as  $R_{s>200}$  and  $R_{s<200}$ , respectively. The value of 200  
188 is approximately the average value of all data in the  $R_s$  matrix (horizontal dotted line in  
189 Fig. 1). Next, two new PCA models were fitted, one with each submatrix. Results from  
190 both models were compared.

191

192           A scatterplot of the scores corresponding to two different components is referred  
193 to as a score plot. The score plot corresponding to PC1 and PC2 (i.e.,  $\mathbf{t}_2$  vs.  $\mathbf{t}_1$ ), referred  
194 to here as the PC1/PC2 plot, usually highlights the basic similarities and dissimilarities  
195 among observations. Score plots with different combinations of PCs were visually  
196 inspected in order to detect outliers as well as to identify stations with a similar  
197 performance. The distance of observations to the PCA model was also checked.

198

### 199 *2.3. PCA infilling approach for $R_s$ estimation*

200 Missing data due to sensor failures is a problem often encountered in MSPC. Different  
201 approaches have been proposed for PCA to deal with incomplete observations [40,41].  
202 One of these algorithms is implemented in the software SIMCA-P 10.0 [42]. Starting  
203 from the complete  $R_s$  matrix, three new ones were obtained containing 5%, 10% and  
204 15% randomly distributed gaps. In the four cases, data were mean-centered prior to  
205 PCA. After obtaining the score and loading vectors for each PC, they were used to  
206 reconstruct the  $\mathbf{X}$  matrix (eq. 1). This procedure was conducted using Matlab version  
207 7.4.0 (MathWorks Inc., Natick, MA, USA), considering an increasing number of PCs.  
208 Next, in order to assess the accuracy of the gap infilling method, the estimated missing  
209 values were compared with the original ones.

210

211           Four additional methods were tested for gap infilling: (i) by adopting as  $R_s$   
212 estimations for a given station, the  $R_s$  records from the nearest station (1-neighbor); (ii)  
213 by obtaining the  $R_s$  average of two nearest stations with a similar altitude (2-neighbor);  
214 (iii) by adopting the  $R_s$  values from the nearest station in the score plot for PC1/PC2 (1-  
215 neighbor-SP); (iv) by assigning the  $R_s$  average of two neighboring stations in the score

216 plot for PC1/PC2 (2-neighbor-SP). These methods will be referred to hereafter with the  
217 name indicated within brackets. In order to assess their efficiency for gap infilling, they  
218 were applied to the matrices with 5%, 10% and 15% of gaps. The neighboring stations  
219 in the 1- and 2-neighbor-SP methods were established only according to the score plot  
220 of the complete matrix.

221

#### 222 2.4. $R_s$ estimation from geographic parameters

223 Principal component regression (PCR) was used to study if score vectors can be  
224 predicted according to geographic and climatic parameters. The proposed methodology  
225 comprises two steps. First, score and loading vectors of the  $k$  relevant PCs were  
226 extracted from the  $R_s$  matrix. Second, step-wise MLR was applied to fit each score  
227 vector according to the following independent variables: *latitude, longitude, altitude,*  
228 *minimum distance to the sea, temperature (average, maximum, minimum), RH, wind*  
229 *speed, Gorezynski continentality index [43] and cumulated rain.* Climatic parameters  
230 correspond to the average values for the 1203 days of study. The software Statgraphics  
231 plus 5.1 (StatPoint Technologies Inc., Warrenton, VA, USA) was used to conduct all  
232 regression models.

233

234 Once obtained the  $k$  predictive equations, they might be applied to estimate the  
235  $t_1, t_2, \dots, t_k$  scores of a new station according to its geographic and climatic data. Next,  
236 the  $R_s$  estimation for the  $j$ -th day would be obtained based on these predicted scores and  
237 the loadings calculated in the first step for the  $j$ -th day (eq. 2), being  $k$  the number of  
238 relevant PCs and  $\mu_j$  the  $j$ -th column average of the  $R_s$  matrix.

239

$$240 \quad (R_s)_j = \mu_j + \hat{t}_1 p_{1j} + \hat{t}_2 p_{2j} + \dots + \hat{t}_k p_{kj} \quad (2)$$

241

242 In order to assess the performance of the proposed method, MLR equations were  
243 applied using the geographic and climatic data of each station and, next, the  $R_s$  matrix  
244 was reconstructed by applying eq. 2 for the 1203 days. Predicted values were compared  
245 with the original ones. The 1- and 2- neighbor-SP approaches described in the previous  
246 section were also tested. The first one consisted of adopting for a given station, the  $R_s$   
247 records from the station with most similar  $t_1$  and  $t_2$  scores. Similarly, the estimation  
248 according to the 2-neighbor-SP method was obtained as the  $R_s$  average of the two  
249 nearest stations in the PC1/PC2 score plot. The  $t_1$  and  $t_2$  scores of the target station were  
250 previously estimated by applying the MLR equations based on its geographic and  
251 climatic parameters.

252

### 253 3 Performance indicators

254 Several error parameters were calculated to assess the performance accuracy of the  
255 proposed estimation methods. The average absolute relative error (AARE), the mean  
256 absolute error (MAE) and the mean squared error (MSE), which are commonly used in  
257 time series analysis, were obtained according to eqs. 3, 4 and 5, respectively, being  $x_i$   
258 the observed  $R_s$  value,  $\hat{x}_i$  the prediction, and  $n$  the number of missing data randomly  
259 created in the  $R_s$  matrix.

$$260 \quad AARE = \frac{1}{n} \cdot \sum_{i=1}^n \left| \frac{x_i - \hat{x}_i}{x_i} \right| \quad (3)$$

$$261 \quad MAE = \frac{1}{n} \cdot \sum_{i=1}^n |x_i - \hat{x}_i| \quad (4)$$

$$262 \quad MSE = \frac{1}{n} \cdot \sum_{i=1}^n (x_i - \hat{x}_i)^2 \quad (5)$$

263

264 **Results and discussion**

265 1. PCA of the  $R_s$  matrix: relevant PCs and outlier detection

266 Three PCA models were carried out, one with all 1203 variables of the  $R_s$  dataset,  
267 another using the set of  $R_s > 200$  variables and a third one with  $R_s < 200$ . In the three  
268 models, the score plot for PC3/PC4 reveals that station s19 presents abnormal values in  
269 both components (figures not shown). This station is the most northern one, which  
270 might explain its different performance. However, the position of s19 in the PC1/PC2  
271 score plot is not abnormal (Fig. 2). Taking into account that PC3 and PC4 provide  
272 relevant information, station s19 was discarded and the three models were repeated. A  
273 summary overview of these six models is shown in Table 2.

274

275 The software SIMCA-P 10.0 considers that a certain component explains  
276 systematic data variability if the goodness-of-fit for that component obtained by cross-  
277 validation ( $Q^2$ ) is higher than a certain threshold [42]. In five of the six models, the  
278 cross-validation criterion is satisfied up to PC4 (Table 2). In order to further investigate  
279 the number of relevant PCs, it was checked that the  $t_1$  score vector obtained from the  
280  $R_s < 200$  model with 29 stations is strongly correlated with that from the  $R_s > 200$  model  
281 ( $r = 0.964$ ,  $p < 0.0001$ ). The correlation is also statistically significant for the  $t_2$ ,  $t_3$  and  
282  $t_4$  vectors ( $p < 0.0001$ ) but not in the case of  $t_5$  ( $r = 0.267$ ,  $p = 0.161$ ) nor  $t_6$  ( $r = 0.291$ ,  
283  $p = 0.125$ ). Again, this result suggests that 4 components should be used to describe the  
284 systematic variability of the  $R_s$  matrix.

285

286 In the  $R_s > 200$  model with 29 observations, station s27 has abnormal values of  
287 PC5 and PC7 and appears as an outlier in the PC5/PC7 score plot (Fig. not shown). The  
288 same result was obtained with the  $R_s < 200$  model and the one with 1203 variables. The

289  $R_s$  pattern of station s27 is slightly different to the rest probably because it has the  
290 highest distance to the sea and the most continental climate. Actually, it presents the  
291 lowest average and minimum temperature among the 30 stations. Nonetheless, s27 was  
292 not discarded because its performance is not abnormal in the four relevant PCs.  
293 Different score plots were visually inspected, but no additional outliers were identified.

294

## 295 2. Similarities among stations based on the score plots

296  $R^2_X$  is usually called goodness-of-fit because it indicates how good is a given PC  
297 to fit the observed values. PC1 explains about 37% of the mean-centered data variability  
298 (Table 2). The coordinate position of stations in the PC1/PC2 score plot, if properly  
299 rotated, is strikingly similar to their geographic position (Fig. 2). The rotation was  
300 achieved by plotting  $(-2 t_1 + t_2)$  vs  $(- t_1 - 2t_2)$ .

301

302 **[FIGURE 2 NEAR HERE]**

303

304 In order to assess if the differences among stations are relevant in practice, the  
305 average  $R_s$  value was calculated for all stations. Averages follow approximately a  
306 normal distribution, being 225.3 the maximum value (station s29) and 187.4, the  
307 minimum value (station s6). Thus, the average  $R_s$  of s29 is 20.2% higher than in the  
308 case of station s6, which highlights the importance of choosing correctly the location for  
309 a solar energy system. It was found that average  $R_s$  values were correlated with  $t_1$   
310 scores ( $r = 0,756$ ). Thus, PC1 will highlight which stations provide higher or lower  
311 values than the average trajectory. Further PCs will describe changes in the shape with  
312 respect to the mean trajectory. Moreover,  $t_1$  scores are also correlated with latitude ( $r =$   
313  $- 0,939$ ), as reflected in Fig. 2. The PC3/PC4 score plot for the  $R_s < 200$  and  $R_s > 200$

314 models are quite similar (Fig. 3), which again indicates that PC3 and PC4 provide  
315 relevant information.

316

317 **[FIGURE 3 NEAR HERE]**

318

319 The PCA model with 4 components accounts for about 66% of the mean-  
320 centered data variability (Table 2). Thus, the PC1/PC2 and PC3/PC4 score plots will  
321 highlight the most relevant similarities and dissimilarities. Stations close to each other  
322 in both plots will present a similar performance, i.e., a trajectory of  $R_s$  recordings with a  
323 similar average value and shape. After visually inspecting both score plots (Figs. 2 and  
324 3), four clusters of stations with a similar  $R_s$  pattern were established: cluster A  
325 (stations s1, s4, s6, and s15), B (s2, s5, s9, and s10), C (s26, s29, and s30) and D (s11,  
326 s12, and s14). They basically differ in latitude (Fig. 2) while C is the cluster with  
327 highest altitude, which implies a more continental climate. Trajectories of stations  
328 belonging to the same clusters were averaged, centered with respect to the mean  
329 trajectory and smoothed using a moving average of order 50 (Fig. 4). Cluster C yields  
330 the trajectory with highest average values and, moreover, its pattern is somewhat  
331 different. This distinctive performance is basically explained by PC2. Clusters C and D  
332 correspond to southern stations and their  $R_s$  values are higher than clusters A and B,  
333 which reflects the negative correlation between latitude and  $R_s$ . Fig. 3 shows that  
334 stations in clusters A and B are discriminated by PC4 which implies that their  
335 trajectories are somewhat different, as reflected by Fig. 4.

336

337 **[FIGURE 4 NEAR HERE]**

338

339 3. Gap infilling results

340 Taking into account that station s19 is an outlier, it was disregarded for the gap infilling  
341 study. PCA was applied to the complete  $R_s$  matrix (29 stations by 1203 variables) and  
342 to the three matrices containing 5%, 10% and 15% of missing data (i.e., 1744, 3489 and  
343 5233 gaps, respectively). Table 3 shows that  $R^2_X$  and  $Q^2$  values of these models are  
344 nearly the same regardless of the amount of missing data. Despite the presence of gaps,  
345 PC4 satisfies the cross-validation criterion. MSE, MAE and AARE (eqs. 3 to 5) were  
346 calculated by comparing discarded data with the predictions obtained by eq. 1 with an  
347 increasing number of PCs. In the complete  $R_s$  matrix, error parameters become null  
348 using 28 components (see Fig. 5), which implies a perfect fit if the maximum number of  
349 possible PCs is used. Only a slight increase of MAE and AARE is observed as the  
350 amount of missing data increases. Fig. 5 also shows that PCA models built with 4 or 5  
351 PCs lead to similar errors. Additional components just provide a slight decrease of the  
352 error indicators, which is consistent with the cross-validation results suggesting that  
353 only four PCs are relevant.

354

355 **[FIGURE 5 NEAR HERE]**

356

357 Missing  $R_s$  data were also estimated according to four alternative methods based  
358 on neighbor assignment. This assignment presents in some cases several valid options  
359 because the number of available stations is limited and there is not always a single  
360 optimum choice. Therefore, only the complete matrix was considered to provide the  $t_1$   
361 and  $t_2$  scores used to establish the neighbor assignment, which was the same for the 3  
362 gap sizes. As observed in Table 4, error indicators are higher if the percentage of  
363 missing data increases. After PCA and 2-neighbor-SP, the best results correspond to the  
364 2-neighbor procedure (Table 4) which presents average indicators slightly better than



365 the 1-neighbor and 1-neighbor-SP methods. As could be expected, the approaches based  
 366 only in the information of one neighboring station are worse than those taking into  
 367 account data from two neighbors.

368

369 4. R<sub>s</sub> estimation from geographical data

370 Score and loading vectors of PC3 and PC4 were extracted from the R<sub>s</sub> matrix with 29  
 371 stations and 1203 variables. Station s19 was disregarded because it becomes an outlier  
 372 in both PCs but not for the previous components. Thus, all stations were considered to  
 373 obtain score and loading vectors of PC1 and PC2.

374

375 MLR was applied next to determine if  $\mathbf{t}_1$  is correlated with *latitude* ( $\varphi$ ),  
 376 *longitude* ( $\tau$ ), *altitude* ( $z$ ), and *distance* to the sea. The same study was conducted with  
 377  $\mathbf{t}_2$ ,  $\mathbf{t}_3$ , and  $\mathbf{t}_4$ . After trying several alternative models, it was decided to consider also two  
 378 indicator variables and their interactions.  $I_{z>400}$  takes the value 1 for the 5 stations with  
 379 an altitude higher than 400 m and zero otherwise. The indicator variable  $I_{\varphi<38.7}$  takes the  
 380 value 1 for stations that satisfy the condition  $\varphi < 38.7$ .

381

$$382 \quad \mathbf{t}_1 = 22403 - 579.1 \varphi + 363.1 I_{\varphi<38.7} - 185.6 \tau + 0.58 z \quad (6)$$

$$383 \quad \mathbf{t}_4 = 7932 - 208.0 \varphi - 333.3 I_{\varphi<38.7} + 676.9 \tau - 341.8 I_{z>400} \quad (7)$$

$$384 \quad \mathbf{t}_2 = -16846 + 421.4 \varphi - 540.7 I_{\varphi<38.7} (\varphi - 39) + 1.3 z \quad (8)$$

$$385 \quad \mathbf{t}_3 = -38125 + 968.1 \varphi - 1611.9 I_{\varphi<38.7} (\varphi - 39.1) - 181 \tau \quad (9)$$

386

387 All regression coefficients of the best predictive equations are statistically  
 388 significant ( $p \leq 0.002$  for eq. 6,  $p \leq 0.0002$  for eq. 7,  $p < 0.003$  for eq. 8, and  $p < 0.009$   
 389 for eq. 9). It was also checked that residuals followed approximately a normal

390 distribution and no outliers were detected. *Longitude* is significantly correlated with  
391 *altitude* ( $r = 0.492, p = 0.006$ ) as well as with *latitude* ( $r = -0.494, p = 0.0055$ ). Given  
392 the correlation among predictive variables, trying to interpret the effect of each  
393 parameter in eqs. 6 to 9 might be misleading.

394

395       Coefficients of determination are the following: 0.983 (eq. 6), 0.901 (eq. 7),  
396 0.785 (eq. 8) and 0.900 (eq. 9). These high values suggest that a considerable amount of  
397 the centered  $R_s$  data variability depends on the geographical position of the station. Fig.  
398 6 shows that  $t_1$  and  $t_2$  scores predicted from eqs. 6 and 8 based on geographical  
399 information are similar to those originally obtained from the  $R_s$  matrix. Predictive MLR  
400 equations for  $t_5$  and  $t_6$  were also tried using geographic and climatic variables, but a  
401 very poor goodness-of-fit was obtained.

402

403                                   **[FIGURE 6 NEAR HERE]**

404

405       It was found that none of the climatic variables entered in the MLR models.  
406 Thus, *latitude*, *longitude* and *altitude*, which are parameters readily available for any  
407 location in the region under study, are enough to predict the four relevant scores. Eqs. 6  
408 to 9 are only valid for the  $R_s$  estimation in weather stations located on the  
409 Mediterranean coast of Spain with similar geographical characteristics as those in Table  
410 1. Nevertheless, the proposed methodology could be applied to any kind of climatic  
411 conditions.

412

413       The PCR approach was applied to reconstruct the  $R_s$  matrix. PCR using four  
414 components provides more accurate estimations than the 1- and 2-neighbor-SP

415 approaches, which were also tested (Table 5). Error indicators of 1-neighbor-SP are  
416 similar as those of PCR using only one component, which suggests that this method  
417 would not be recommended. By contrast, errors of 2-neighbor-SP and PCR with 3  
418 components are similar. Again, these results indicate that it would be better to use data  
419 from two neighboring stations instead of just one.

420

## 421 **Conclusions**

422 Choosing the right location for a solar energy system is a key factor for  
423 maximizing the power generation. Moreover, the estimation of daily  $R_s$  values is of  
424 interest for the design of photovoltaic systems and energy efficient buildings. Available  
425  $R_s$  data are useful for this purpose, particularly if they are recorded from nearby weather  
426 stations. PCA was applied to  $R_s$  data recorded at 30 stations in the Mediterranean coast  
427 of Spain. Four principal components account for the systematic data variation and  
428 explain about 66% of the mean-centered  $R_s$  variability. By means of MLR, it was found  
429 that the latent variables associated to the four relevant PCs can be predicted according to  
430 latitude, longitude and altitude. Climatic variables did not increase the predictive  
431 goodness-of-fit. Based on the results, a new methodology is proposed to estimate daily  
432  $R_s$  values at any location in the region under study when only local geographical  
433 parameters are available. The proposed method exhibits a higher accuracy than simpler  
434 procedures using data from neighboring stations.

435

436 Time series of  $R_s$  often present data gaps or discontinuities. In practice, this  
437 problem is often solved by adopting the measurements from neighboring stations. The  
438 PCA approach characterizes the similarities among weather stations and also allows  
439 estimation of present and past  $R_s$  values. The proposed method for gap infilling is more

440 accurate than four alternative procedures also tested. The statistical methodology  
441 applied here is commonly used in the context of MSPC, particularly in chemical  
442 processes, but as far as we know this is the first work that applies such methodology for  
443 the estimation of solar radiation.

444

#### 445 **Mathematical notation**

**X** matrix (upper case, bold)

**t** vector, i.e. column matrix (lower case, bold)

**p<sup>T</sup>** transposed vector, i.e. row matrix (lower case, bold)

*k* scalar (lower case, italicized)

446

447

#### 448 **Acknowledgements**

449 We are grateful to the Valencian Institute of Agricultural Research (IVIA) for providing  
450 the dataset used in the present work.

451

452 **References**

- 453 [1] López G, Rubio MA, Martínez M, Batlles FJ. Estimation of hourly global  
454 photosynthetically active radiation using artificial neural network models. *Agric  
455 Forest Meteorol* 2001;107:279-91.
- 456 [2] Reddy KS, Ranjan M. Solar resource estimation using artificial neural networks and  
457 comparison with other correlation models. *Energy Convers Manage* 2003;44:2519–  
458 30.
- 459 [3] Hunt LA, Kuchar L, Swanton CJ. Estimation of solar radiation for use in crop  
460 modelling. *Agric Forest Meteorol* 1998;91:293-300.
- 461 [4] Abraha MG, Savage MJ. Comparison of estimates of daily solar radiation from air  
462 temperature range for application in crop simulations. *Agric Forest Meteorol*  
463 2008;148:401-16.
- 464 [5] Hontoria L, Aguilera J, Zufiria P. Generation of hourly irradiation synthetic series  
465 using the neural network multilayer perceptron. *Sol Energy* 2002;72(5):441-6.
- 466 [6] Zekai Ş. *Solar energy fundamentals and modeling techniques: atmosphere,  
467 environment, climate change and renewable energy*. London: Springer; 2008.
- 468 [7] Meza F, Varas E. Estimation of mean monthly solar global radiation as a function of  
469 temperature. *Agric Forest Meteorol* 2000;100:231-41.
- 470 [8] Tymvios FS, Jacovides CP, Michaelides SC, Scouteli C. Comparative study of  
471 Ångström's and artificial neural networks' methodologies in estimating global solar  
472 radiation. *Sol Energy* 2005;78:752-62.
- 473 [9] Muneer T, Younes S, Munawwar S. Discourses on solar radiation modeling. *Renew  
474 Sustain Energy Rev* 2005;11:551-602
- 475 [10] Mohandes M, Rehman S, Halawani TO. Estimation of global solar radiation using  
476 artificial neural networks. *Renew Energy* 1998;14(1-4):179-84.

- 477 [11] Rehman S, Mohandes M. Artificial neural network estimation of global solar  
478 radiation using air temperature and relative humidity. *Energy Policy* 2008;36:571-6.
- 479 [12] Bosch JL, López G, Batlles FJ. Daily solar irradiation estimation over a  
480 mountainous area using artificial neural networks. *Renew Energy* 2008;33:1622-8.
- 481 [13] Azadeh A, Maghsoudi A, Sohrabkhani S. An integrated artificial neural networks  
482 approach for predicting global radiation. *Energy Convers Manage* 2009;50:1497-  
483 1505.
- 484 [14] Mellit A, Kalogirou SA, Hontoria L, Shaari S. Artificial intelligence techniques for  
485 sizing photovoltaic systems: A review. *Renew Sustain Energy Rev* 2009;13:406-19.
- 486 [15] Mellit A, Benghanem M, Hadj Arab A, Guessoum A. An adaptive artificial neural  
487 network model for sizing stand-alone photovoltaic systems: application for isolated  
488 sited in Algeria. *Renew Energy* 2005;30:1501-24.
- 489 [16] Bakirci K. Models of solar radiation with hours of bright sunshine: A review.  
490 *Renew Sustain Energy Rev* 2009;13:2580-8.
- 491 [17] Liu X, Mei X, Li Y, Wang Q, Jensen JR, Zhang Y, Porter JR. Evaluation of  
492 temperature-based global solar radiation models in China. *Agric Forest Meteorol*  
493 2009;149:1433–46.
- 494 [18] Hargreaves GH, Samani ZA. Estimating potential evapotranspiration. *J Irrig Drain*  
495 *Eng* 1982;108(3):225-30.
- 496 [19] Bristow KL, Campbell GS. On the relationship between incoming solar radiation  
497 and daily maximum and minimum temperature. *Agric Forest Meteorol*  
498 1984;31:159-66.
- 499 [20] Allen RG. Self-calibrating method for estimating solar radiation from air  
500 temperature. *J Hydrol Eng* 1997;2(2):56-67.

- 501 [21] Fadare DA. Modelling of solar energy potential in Nigeria using an artificial neural  
502 network model. *Appl Energy* 2009;86:1410–22.
- 503 [22] Al-Alawi SM, Al-Hinai HA. An ANN-based approach for predicting global  
504 radiation in locations with no direct measurement instrumentation. *Renew Energy*  
505 1998;14(1-4):199-204.
- 506 [23] Mohandes M, Balghonaim A, Kassas M, Rehman S, Halawani TO. Use of radial  
507 basis functions for estimating monthly mean daily solar radiation. *Sol Energy*  
508 2000;68(2):161-8.
- 509 [24] Sfetsos A, Coonick AH. Univariate and multivariate forecasting of hourly solar  
510 radiation with artificial intelligence techniques. *Sol Energy* 2000;68(2):169-78.
- 511 [25] Dorvlo ASS, Jervase JA, Al-Lawati A. Solar radiation estimation using artificial  
512 neural networks. *Appl Energy* 2002;71:307–19.
- 513 [26] Sözen A, Arcaklioğlu E, Özalp M. Estimation of solar potential in Turkey by  
514 artificial neural networks using meteorological and geographical data. *Energy*  
515 *Convers Manage* 2004;45:3033–52.
- 516 [27] Lam JC, Wan KKW, Yang L. Solar radiation modelling using ANNs for different  
517 climates in China. *Energy Convers Manage* 2008;49:1080–90.
- 518 [28] Benghane M, Mellit A, Alamri SN. ANN-based modelling and estimation of  
519 daily global solar radiation data: A case study. *Energy Convers Manage*  
520 2009;50:1644-55.
- 521 [29] Kalogirou S, Michaelides S, Tymvios F. Prediction of maximum solar radiation  
522 using artificial neural networks. *Proceedings of World Renewable Energy*  
523 *Congress VII, Cologne, Germany; 2002.*
- 524 [30] López G, Batlles FJ, Tovar-Pescador J. Selection of input parameters to model  
525 direct solar irradiance by using artificial neural networks. *Energy* 2005;30:1675-84.

- 526 [31] Hontoria L, Aguilera J, Zufiria P. An application of the multilayer perceptron:  
527 Solar radiation maps in Spain. *Sol Energy* 2005;79:523-30.
- 528 [32] Dunia R, Qin SJ, Edgar TF, McAvoy TJ. Identification of faulty sensors using  
529 principal component analysis. *AIChE J* 1996;42:2797-2812.
- 530 [33] Jolliffe IT. *Principal component analysis*. 2nd ed. New York: Springer-Verlag;  
531 2002.
- 532 [34] García-Diego FJ, Zarzo M. Microclimate monitoring by multivariate statistical  
533 control: The renaissance frescoes of the Cathedral of Valencia (Spain). *J Cult*  
534 *Heritage* 2010;11:339-44.
- 535 [35] Wold H. Estimation of principal components and related models by iterative least  
536 squares. In: Krishnaiah PR, editor. *Multivariate analysis*, New York: Academic  
537 Press; 1966, p. 391-420.
- 538 [36] Nomikos P, MacGregor JF. Monitoring of batch processes using multiway  
539 principal component analysis. *AIChE J* 1994;40:1361-75.
- 540 [37] Rännar S, MacGregor JF, Wold S. Adaptive batch monitoring using hierarchical  
541 PCA. *Chemom Intell Lab Syst* 1998;41:73-81.
- 542 [38] Zarzo M, Ferrer A. Batch process diagnosis: PLS with variable selection versus  
543 block-wise PCR. *Chemom Intell Lab Syst* 2004;73:15-27.
- 544 [39] Diana G, Tommasi C. Cross-validation methods in principal components: a  
545 comparison. *Statist Methods Applicat* 2002;11(1):71-82.
- 546 [40] Nelson PRC, Taylor PA, MacGregor JF. Missing data methods in PCA and PLS:  
547 Score calculations with incomplete observations. *Chemom Intell Lab Syst*  
548 1996;35:45-65.
- 549 [41] Arteaga F, Ferrer A. Dealing with missing data in MSPC: several methods,  
550 different interpretations, some examples. *J Chemometrics* 2002;16:408-18.



- 551 [42] Eriksson L, Johansson E, Kettaneh-Wold N, Wold S. Introduction to multi- and  
552 megavariate data analysis using projection methods (PCA & PLS). Umea, Sweden:  
553 Umetrics AB; 1999.
- 554 [43] Martí P, M. Gasque M. Improvement of temperature-based ANN models for solar  
555 radiation estimation through exogenous data assistance. Energy Convers Manage.  
556 doi:10.1016/j.enconman.2010.08.027
- 557

558 **Figure captions**

559

560 **Fig. 1.**  $R_s$  values averaged for the 30 stations. Gaps correspond to days with missing  
561 data for at least one station.

562

563 **Fig. 2.** *Left:* Map of the eastern coast of Spain (provinces of Alicante, Valencia and  
564 Castellón) indicating the location of the 30 weather stations (codes as in Table 1). *Right:*  
565 rotated score plot of PC1/PC2 obtained from the initial  $R_s$  matrix (30 stations by 1203  
566 variables).

567

568 **Fig. 3.** Score plot of PC3/PC4 from the  $R_s > 200$  model (filled triangles) and  $R_s < 200$   
569 (empty triangles). Station s19 was disregarded. Both plots were slightly rotated to  
570 achieve a better fit between scores corresponding to the same station

571

572 **Fig. 4.**  $R_s$  centered trajectories (i.e. difference with respect to the mean trajectory)  
573 averaged for stations with a similar performance (A: s1-s4-s6-s15, B: s2-s5-s9-s10, C:  
574 s26-s29-s30, D: s11-s12-s14).

575

576 **Fig. 5.** Error parameters showing the gap infilling performance of PCA for different gap  
577 sizes with an increasing number of PCs.

578

579 **Fig. 6.** Comparison between  $t_1$  and  $t_2$  scores from the  $R_s$  matrix with those obtained by  
580 applying eqs. 6 and 8 taking into account the geographical data of stations.

581

582

583 **Table 1**

584 Geographical parameters of the 30 weather stations.  $z$  : altitude (m) with respect to sea  
 585 level;  $\varphi$  : latitude (degrees);  $\tau$  : longitude (degrees).

586

Station	$c^a$	$z$	$\varphi$	$\tau^b$		$c^a$	$z$	$\varphi$	$\tau^b$
Benavites	1	8	39.7333	0.2150	Dénia-Gata	16	102	38.7939	-0.0836
Tavernes de Valldigna	2	15	39.0964	0.2367	Vila Joiosa	17	138	38.5294	0.2553
Catral	3	27	38.1544	0.8042	Pedralba	18	200	39.5678	0.7164
Sagunt	4	33	39.6492	0.2925	San Rafel del Riu	19	205	40.5956	-0.3703
Carcaixent	5	35	39.1167	0.5047	Altea	20	210	38.6056	0.0775
Vila Real	6	42	39.9333	0.1000	Monforte del Cid	21	244	38.3997	0.7289
Ondara	7	49	38.8197	-0.0075	Llíria	22	250	39.6919	0.6253
Moncada	8	58	39.5877	0.3992	Turís	23	299	39.4006	0.6836
Vilanova de Castelló	9	58	39.0667	0.5228	Cheste	24	323	39.5217	0.7417
Carlet	10	66	39.2264	0.5459	Agost	25	345	38.4278	0.6433
Almoradí	11	74	38.0908	0.7714	Villena	26	495	38.5967	0.8733
Pilar de la Horadada	12	77	37.8700	0.8103	Campo Arcís	27	584	39.4344	1.1608
Elx	13	86	38.2667	0.7000	El Pinós	28	606	38.4286	1.0594
Orihuela	14	99	38.1828	0.9536	Camp de Mirra	29	627	38.6803	0.7717
Vall d'Uixó	15	100	39.7975	0.2272	Castalla	30	708	38.6053	0.6728

587

588 <sup>a</sup>Station code, which was assigned according to an increasing altitude.

589 <sup>b</sup>Positive values: West; negative values: East, with respect to the Greenwich meridian.

590

591

592 **Table 2**

593 Summary overview of 3 PCA models: (i)  $R_s$  matrix with 1203 variables ('all'), (ii)  
 594 subset of 604 variables with an average  $R_s < 200$ , and (iii) subset of 599 variables with  
 595 an average  $R_s > 200$ . These models were repeated after discarding station s19.  
 596 Goodness-of-fit ( $R^2_x$ ), eigenvalue ( $\lambda$ ), goodness-of-fit by cross-validation ( $Q^2$ ), and  
 597 threshold value ( $Q^2_{limit}$ ).

598

model	PC	PCA with 30 stations				PCA with 29 stations			
		$R^2_x$	$\lambda$	$Q^2$	$Q^2_{limit}$	$R^2_x$	$\lambda$	$Q^2$	$Q^2_{limit}$
All	1	0.366	11.0	0.300	0.034	0.376	10.9	0.312	0.035
All	2	0.129	3.86	0.097	0.035	0.138	4.01	0.120	0.036
All	3	0.104	3.12	0.111	0.036	0.100	2.90	0.140	0.038
All	4	0.061	1.83	0.041	0.038	0.057	1.67	0.070	0.039
All	5	0.045	1.34	0.011	0.039	0.043	1.25	-0.038	0.041
All	6	0.039	1.18	0.012	0.041	0.038	1.10	0.004	0.042
$R_s < 200$	1	0.380	11.4	0.322	0.035	0.389	11.3	0.334	0.036
$R_s < 200$	2	0.101	3.04	0.061	0.036	0.107	3.11	0.038	0.037
$R_s < 200$	3	0.096	2.87	0.076	0.037	0.098	2.84	0.125	0.039
$R_s < 200$	4	0.072	2.15	0.061	0.039	0.065	1.88	0.044	0.040
$R_s < 200$	5	0.047	1.41	-0.020	0.040	0.049	1.42	-0.059	0.042
$R_s < 200$	6	0.043	1.29	-0.007	0.042	0.045	1.30	0.028	0.043
$R_s > 200$	1	0.363	10.9	0.281	0.035	0.375	10.9	0.294	0.036
$R_s > 200$	2	0.156	4.69	0.153	0.036	0.165	4.79	0.192	0.037
$R_s > 200$	3	0.107	3.22	0.130	0.037	0.105	3.04	0.156	0.039
$R_s > 200$	4	0.057	1.70	-0.003	0.039	0.049	1.41	0.047	0.040
$R_s > 200$	5	0.042	1.26	0.026	0.040	0.042	1.21	-0.021	0.042
$R_s > 200$	6	0.038	1.14	-0.017	0.042	0.037	1.08	0.035	0.043

599

600

601 **Table 3**

602 PCA of the  $R_s$  matrix (29 stations, 1203 variables): goodness-of-fit ( $R^2_x$ ) and goodness-  
 603 of-fit by cross-validation ( $Q^2$ ) considering different gap sizes.

604

PC	0% gaps		5% gaps		10% gaps		15% gaps	
	$R^2_x$	$Q^2$	$R^2_x$	$Q^2$	$R^2_x$	$Q^2$	$R^2_x$	$Q^2$
1	0.376	0.312	0.375	0.324	0.376	0.321	0.385	0.327
2	0.138	0.120	0.139	0.123	0.142	0.126	0.140	0.120
3	0.100	0.140	0.100	0.133	0.101	0.136	0.099	0.129
4	0.057	0.070	0.058	0.042	0.058	0.043	0.059	0.048
5	0.043	-0.038	0.044	0.005	0.044	-0.017	0.043	-0.013
6	0.038	0.004	0.037	0.004	0.037	0.011	0.036	0.006

605

606 **Table 4**

607 Error parameters (eqs. 3 to 5) as performance indicators for several infilling methods  
 608 and gap sizes: PCA approach based on 4 components (A), 1-neighbor (B), 2-neighbor  
 609 (C), 1-neighbor-SP (D), and 2-neighbor-SP (E).

610

611

gap size	method	MSE	MAE	AARE
5%	A	297.54	11.840	0.0814
10%	A	343.31	12.505	0.0931
15%	A	391.21	13.306	0.0949
5%	B	373.64	13.736	0.0857
10%	B	435.44	14.388	0.0959
15%	B	459.40	14.775	0.0944
5%	C	340.14	12.965	0.0834
10%	C	393.27	13.290	0.0910
15%	C	414.81	13.858	0.0912
5%	D	457.63	13.941	0.1070
10%	D	435.51	14.480	0.0941
15%	D	461.63	14.342	0.0994
5%	E	309.96	12.459	0.0861
10%	E	365.77	12.983	0.0969
15%	E	373.31	13.119	0.0925

612

613

614 **Table 5**

615 Error parameters (eqs. 3 to 5) as indicators of the goodness-of-fit for the reconstruction  
 616 of the  $R_s$  matrix according to four different methods.

617

method	N	MSE	MAE	AARE
PCA <sup>a</sup>	1	422.76	14.283	0.1135
	2	336.16	12.425	0.1015
	3	266.71	11.170	0.0870
	4	225.72	10.295	0.0767
PCR <sup>b</sup>	1	427.03	14.409	0.1142
	2	359.78	13.104	0.1049
	3	311.93	12.252	0.0923
	4	286.04	11.784	0.0853
1-neighbor-SP	-	425.19	14.044	0.0990
2-neighbor-SP	-	309.45	11.853	0.0842

618

619 <sup>a</sup> $R_s$  matrix (centered values) reconstructed according to eq. 1 with an increasing number  
 620 of components (N).

621 <sup>b</sup>Same as the PCA method but  $t_i$  scores were obtained using eqs. 6 to 9.

622

623