

1 **Methods for interpolating missing data in aerobiological databases**

2 Picornell, A.^{a,*}; Oteros, J.^{b,c}; Ruiz-Mata, R.^a; Recio, M.^a; Trigo, M.M.^a; Martínez-
3 Bracero, M.^{b,c,d}; Lara, B.^e; Serrano-García, A.^e; Galán, C.^{b,c}; García-Mozo, H.^{b,c};
4 Alcázar, P.^{b,c}; Pérez-Badía, R.^e; Cabezudo, B.^a; Romero-Morte, J.^e; Rojo, J.^{e,f}

5 a. Department of Botany and Plant Physiology. University of Malaga.
6 Campus de Teatinos s/n E-29071. Malaga (Spain).

7 b. Department of Botany, Ecology and Plant Physiology. Agrifood Campus
8 of International Excellence CeiA3, University of Cordoba. Cordoba
9 (Spain).

10 c. Andalusian Inter-University Institute for Earth System IISTA, University of
11 Cordoba, Spain.

12 d. School of Chemical and Pharmaceutical Sciences. Technological
13 University Dublin. Dublin (Ireland).

14 e. University of Castilla-La Mancha. Institute of Environmental Sciences
15 (Botany). Toledo (Spain).

16 f. Department of Pharmacology, Pharmacognosy and Botany, Complutense
17 University. Madrid (Spain).

18 * Corresponding author: Antonio Picornell

19 Department of Botany and Plant Physiology, University of Malaga.

20 Campus de Teatinos s/n, Malaga, E-29071, Spain.

21 E-mail address: picornell@uma.es

22 +34 952131912

23 **Abstract**

24 Missing data is a common problem in scientific research. The availability of
25 extensive environmental time series is usually laborious and difficult, and
26 sometimes unexpected failures are not detected until samples are processed.
27 Consequently, environmental databases frequently have some gaps with missing
28 data in it. Applying an interpolation method before starting the data analysis can
29 be a good solution in order to complete this missing information. Nevertheless,
30 there are several different approaches whose accuracy should be considered and
31 compared. In this study, data from 6 aerobiological sampling stations were used
32 as an example of environmental data series to assess the accuracy of different
33 interpolation methods. For that, observed daily pollen/spore concentration data
34 series were randomly removed, interpolated by using different methods and then,
35 compared with the observed data to measure the errors produced. Different
36 periods, gap sizes, interpolation methods and bioaerosols were considered in
37 order to check their influence in the interpolation accuracy. The moving mean
38 interpolation method obtained the highest success rate as average. By using this
39 method, a success rate of the 70% was obtained when the risk classes used in
40 the alert systems of the pollen information platforms were taken into account. In
41 general, errors were mostly greater when there were high oscillations in the
42 concentrations of biotic particles during consecutive days. That is the reason why
43 the pre-peak and peak periods showed the highest interpolation errors. The
44 errors were also higher when gaps longer than 5 days were considered. So, for
45 completing long periods of missing data, it would be advisable to test other
46 methodological approaches. A new Variation Index based on the behaviour of the
47 pollen/spore season (measurement of the variability of the concentrations every

48 2 consecutive days) was elaborated, which allows to estimate the potential error
49 before the interpolation is applied.

50 **Keywords**

51 Missing data; aerobiology; time-series; modelling; interpolation; environmental
52 sampling; bioaerosols

53 **1. Introduction**

54 Environmental time series databases require continuous and reliable monitoring
55 systems which may be affected by technical breakdowns and human factors that
56 can interrupt the sampling process (Oteros et al., 2013). Thus, the presence of
57 gaps in time series data is a very widespread problem in scientific research
58 (Junger and Ponce de Leon, 2015; Navares and Aznarte, 2019; Orlandi et al.,
59 2014; Rubin, 1976; Schouten et al., 2018). This is why in many scientific
60 disciplines, interpolation is commonly used to complete missing data or to
61 increase its resolution (Lehmann et al., 1999; Luedeling et al., 2013; J. Oteros et
62 al., 2013).

63 Many different methodologies have been developed for completing missing
64 information depending on the nature of the data and the required accuracy. Some
65 of the most extended methods are multiple imputation-based, and multiple
66 likelihood-based estimations (Junger and Ponce de Leon, 2015). These methods
67 are widely implemented in most statistical softwares, in particular, in several
68 statistical R packages (e.g. “chillR”, “MICE”, “missForest”, “rrcovNA”, “mtsdi”,
69 “mi”). Some of them use either parametric and non-parametric statistics methods
70 (Junger and Ponce de Leon, 2015; Luedeling et al., 2013; Stekhoven and
71 Buhlmann, 2012; Su et al., 2011; Todorov, 2020; van Buuren and Groothuis-

72 Oudshoorn, 2011). In general, these methods analyse the nature of the data in
73 order to create new data that can replace the missing observations when they
74 are randomly distributed (i.e. data are missing independently of their value or the
75 value of other related variables) (Rubin, 1976).

76 Aerobiology is the scientific discipline based on the study of the atmospheric
77 bioaerosol dynamics (pollen, spore, bacteria, virus...) (Fröhlich-Nowoisky et al.,
78 2016). In this context, in case of gaps detection, it is not enough to create data
79 which are statistically coherent with the rest of the database. Aerobiological data
80 are sequential and missing data estimation must be linked with the previous and
81 subsequent observations, given the stochastic nature of the bioaerosols in the
82 atmosphere. In addition, aerobiological data follow a time series evolution that
83 does not fit the stationary criterion, which increases the difficulty of predictions
84 (Ritenberga et al., 2016).

85 Additionally, aerobiological samplings are complex and sometimes pollen traps
86 are installed in non-easily accessible locations (García-Mozo et al., 2007; Oteros
87 et al., 2019; Picornell et al., 2019c). In the case of the Hirst-type volumetric traps,
88 the proper operation of the traps is checked once a week, but unexpected failures
89 such as power outages or device breakdowns may happen in between, resulting
90 in a few days period of missing data (Navares and Aznarte, 2019). Even the
91 development of new real-time automatic sampling devices also requires
92 interpolation methods to complete gaps during the phase of processing the
93 database (Oteros et al., 2020). Such missing data events are produced
94 completely at random (Missing Completely At Random; MCAR) since they are
95 not conditioned, a priori, by any other variable or by their concentrations values
96 (Junger and Ponce de Leon, 2015). In some cases, the gaps may may not

97 hamper proper data analysis, but in other cases it can seriously affect the
98 establishment of the principal dates of the main pollen season (MPS) or the main
99 spore season (MSS) (Navares and Aznarte, 2017; Picornell et al., 2019a).
100 Alternatively, missing concentration values might be considered as 0 pollen
101 grains or spores/m³ of air for most MPS/MSS definitions, which in many cases
102 would produce more errors than estimated data.

103 The ideal method to complete missing concentration data in aerobiological
104 databases, in terms of usability, should be independent of other variables and
105 directly applicable. Linear interpolation is the most commonly used method to
106 complete missing data, but its use is not so extended in Aerobiology as in other
107 disciplines (Belmonte et al., 1999; Gabarra et al., 2002; Navares and Aznarte,
108 2019, 2017; Picornell et al., 2019a; Skjøth et al., 2016). Other methods are rarely
109 applied (e.g. moving mean interpolation or interpolation by using data of nearby
110 location) and their accuracy have never been measured (Jesús Rojo et al., 2019;
111 Skjøth et al., 2016).

112 For all the aforementioned, the main aim of this study was to comparatively
113 evaluate different methods which allow to interpolate aerobiological data, as well
114 as to check their effectiveness and accuracy depending on the circumstances
115 based on real life observation data.

116 **2. Material and methods**

117 To carry out this study, the databases of 6 aerobiological stations, situated in
118 different localities of the Iberian Peninsula, have been used: Cordoba,
119 Hornachuelos Natural Park, Malaga, Ronda, Sierra de las Nieves Natural Park,
120 and Toledo (Fig. 1). The altitudinal range of the sampling stations varied from 58

121 to 1073 m a.s.l., with an average annual total precipitation between 382 mm and
122 996 mm, and an annual mean temperature of between 11.9 °C and 18.4 °C (Table
123 1). All the sampling stations are within the Mediterranean macroclimate (Rivas-
124 Martínez et al., 2017).



125
126 **Fig. 1.** Map of the pollen and spores sampling stations used in this study. Spatial
127 information obtained from REDIAM (Junta de Andalucía, 2011). NP: Natural
128 Park.

129 **Table 1.** Climatic parameters, sampling years and coordinates of the sampling
130 sites included in this study. Data extracted from García-Mozo et al., 2006;
131 Hernández-Ceballos et al., 2015; Picornell et al., 2020, 2019b.

Location	Annual total rainfall (mm)	Annual average temperature (°C)	Altitude (m a.s.l.)	Coordinates	Years of sampling
Cordoba	621	17.8	138	37°54' N 4°43' W	2006-2018
Hornachuelos NP	700	16.8	225	38°4' N 5°24' W	1998-2019
Malaga	540	18.4	58	36°42' N 4°28' W	1991-2019
Ronda	681	16.4	768	36°44' N 5°10' W	2017-2019
Sierra de las Nieves NP	996	11.9	1073	36°39' N 5°5' W	2018-2019
Toledo	342	15.8	450	39°51' N 4°2' W	2003-2019

132 NP: Natural Park.

133 2.1. Pollen and spore data

134 Airborne pollen and fungal spores were collected by means of 6 Hirst-type
135 volumetric traps, one per location (Hirst, 1952). The air flow was adjusted in all of
136 them to 10 l/min. The aerobiological samples obtained were processed and
137 analysed following the recommendations of both the Spanish Aerobiology
138 Network (REA) (Galán et al., 2007) and the European Aerobiology Society (EAS)
139 (Galán et al., 2014). More than the 10% surface of each daily sample were
140 analysed for pollen identification and counted by light microscopy at a
141 magnification of 400X. In the case of *Alternaria* spores, at least the 5% of each
142 daily sample were analysed at the same magnification (Galán et al., 2021). Pollen
143 and spore concentrations were expressed as pollen grains/m³ of air and

144 spores/m³ of air, respectively, according to the international recommendations
145 (Galán et al., 2017, 2014).

146 Daily pollen concentrations of *Amaranthaceae*, *Cupressaceae*, *Olea*, *Pinus*,
147 *Plantago*, *Platanus*, *Poaceae*, *Quercus*, and *Urticaceae* were used to test the
148 accuracy of the interpolation methods at all sampling stations, while *Arecaceae*
149 and *Casuarina* pollen concentrations only were used in Malaga, Ronda, and
150 Sierra de las Nieves, due to its scarcity in the atmosphere of the other localities.
151 Regarding fungal spores, *Alternaria* concentrations registered in Cordoba,
152 Malaga, Ronda, Sierra de las Nieves, and Toledo were also included in the study.

153 For each pollen/spore type and year, the main pollen/spore season was
154 calculated, this being defined as the period between the first day of the year in
155 which the 5% of the annual pollen/spore integral is reached and the first day in
156 which the 95% annual is accumulated (Nilsson and Persson, 1981). In the case
157 of *Cupressaceae* and *Alternaria* two different pollen/spore curves were detected
158 within a year in all sampling sites. Therefore, each curve was studied separately
159 by dividing the year in two periods: January-July (winter *Cupressaceae* and
160 spring *Alternaria*) and August-December (autumn *Cupressaceae* and *Alternaria*).
161 Since *Urticaceae* pollen type was abundantly detected during the whole year, the
162 start and end dates of the MPS were defined by adjusting the cumulative pollen
163 concentrations to a logistic curve and selecting the dates in which the fourth
164 derivative of the logistic curve crossed the x-axis (Cunha et al., 2015; Ribeiro et
165 al., 2007). The MPS/MSS were calculated with the “AeRobiology” R package
166 (Jesús Rojo et al., 2019). The optimal definition method was applied in each case.
167 The defined seasons, independently of the method applied, helped to categorize
168 the time series into different periods in order to analyse any effect of the seasonal

169 stages on the missing data estimation. Therefore, the method used to define the
170 pollen seasons was not a crucial point in this study, and the results are not
171 affected by them.

172 **2.2. Interpolation methods tested**

173 For this study, the different interpolation methods integrated in the “AeRobiology”
174 R package were tested, i.e. linear interpolation, moving mean interpolation, spline
175 interpolation, interpolation by using time series analysis, and interpolation by
176 using data from nearby locations (Jesús Rojo et al., 2019). In all cases,
177 calculations are based on the daily mean pollen/spore concentrations.

178 2.2.1. Linear interpolation

179 A linear regression is calculated by taking the first data previous and subsequent
180 to the gap (i.e., first day before and after the gap), so, the missing data are
181 estimated by using the regression equation (Fig. 2).

182 2.2.2. Moving mean interpolation

183 Each missing data is replaced by the mean value of a certain number of data
184 placed on both sides of the gap. The number of days took for calculating the
185 mean is the double of the gap size, and it is centred in the missing value (Fig. 2).

186 2.2.3. Spline interpolation

187 A spline regression is calculated by taking the first 3 data on both sides of the
188 gap. Then, the missing data are estimated by using the regression equation (Fig.
189 2).

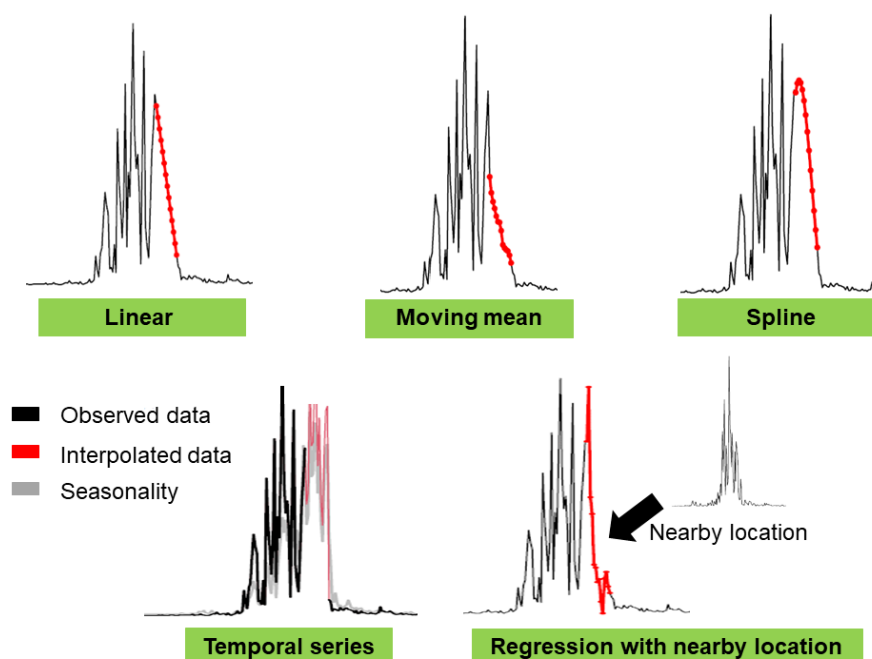
190 2.2.4. Interpolation by using temporal series analysis

191 For each pollen/spore type, the seasonality is calculated by taking all the daily
192 data available for several years by performing a seasonal trend decomposition
193 based on LOESS (Cleveland et al., 1990). Then, a linear regression is calculated
194 between the pollen/spore curve of the year in which there are missing data and
195 the seasonality curve in order to regulate the curve intensity based on the known
196 data of the target pollen/spore season (Fig. 2). Missing data are estimated by
197 using the regression equation.

198 2.2.5. Interpolation with data from nearby locations

199 In this case, the database of a nearby sampling station with complete data for the
200 missing period is used to complete the gaps. For the year in which the loss of
201 data occurred, a linear regression is calculated between the pollen/spore curve
202 of the nearby locality (independent variable), and the pollen/spore curve of the
203 target locality (dependent variable). If the regression is significant ($p\text{-value} \leq 0.05$)
204 and the regression coefficient is higher than 0.6, the data from the nearby location
205 are transformed by applying the regression equation, and the missing data are
206 replaced by the calculated values (Fig. 2). Regression coefficients under 0.6 have
207 not been considered high enough to reflect a direct relationship between the
208 concentration values of both sampling locations. In this method, it is possible to
209 include more than one nearby location simultaneously. In such cases, the data of
210 each nearby sampling station is included as an independent variable in a multiple
211 linear regression and the missing data (dependant variable) is calculated by
212 applying the regression equation.

213 This method was tested in Cordoba by using Hornachuelos National Park as
214 nearby sampling station, as well as in Ronda by using Malaga, and Sierra de las
215 Nieves databases as nearby localities, both individually and simultaneously.



216

217 **Fig. 2.** Graphical examples to visualise the application of the different methods
218 of data interpolation applied in this study. Example elaborated with *O/ea* pollen
219 data in Ronda during 2018.

220 **2.3. Relative error calculation**

221 To check the effectiveness of each interpolation method, some observation data
222 were removed from the original databases of each sampling site in order to create
223 artificial gaps. After that, the data missed were interpolated by the methods
224 explained above and the estimated data were compared with those removed.

225 To avoid bias in removing the original data, an algorithm that performs random
226 cuts in the data series was developed. This algorithm made random cuts in
227 different periods of the pollen seasons for the different pollen/spore types, being
228 these periods: pre-season, pre-peak, peak, post-peak and post-season. The pre-
229 and post-season periods are those outside the MPS/MSS. The peak cut was

230 obtained by centring the peak day in the centre of the removed data. Cuts of 3,
231 5, 7 and 10 consecutive days were tested.

232 Relative errors (RE) were calculated by means of the formula according to
233 equation 1. In such of the mathematical formula, the error values range between
234 0 and 2. Cases whose observed concentrations were zero pollen grains or
235 spores/m³ caused mathematical indeterminacy when the estimated value was
236 zero too (0/0), and relative errors of 2 when the estimated value was non-zero.
237 Therefore, they were excluded since these concentrations were not frequent and
238 they have scarce relevance.

239 (1)

240
$$\text{Relative error} = \frac{|e - o|}{\frac{|e| + |o|}{2}} \text{ if } o \neq 0$$

241 where *e* is the estimated pollen/spore concentration, and *o* the observed
242 concentration.

243 Besides this, estimated and observed pollen/spore concentrations were classified
244 into the Spanish Aerobiology Network pollen classes (nil, low, moderate and high)
245 (Galán et al., 2007), but nil class was modified to concentrations ≤ 1 pollen grain
246 or spore/m³. Due to there are no stablished classes for *Alternaria* spores, the
247 thresholds for the moderate and high categories were set in 30 and 50 spores/m³
248 respectively, according to the most frequent concentrations detected in the
249 sampling sites. After classifying the observed and estimated pollen/spore
250 concentrations, the percentage of correct classification, i.e. the success rate
251 (observed category = estimated category), was calculated by means of equation

252 2, in which, observed concentrations of 0 pollen grains or spores/m³ were not
253 excluded since they do not induce mathematical artefacts in the formula.

254 (2)

$$255 \quad \text{Success rate} = \frac{N^{\circ} \text{ correct classifications}}{N^{\circ} \text{ total classifications}} * 100$$

256 Differences in the relative errors and in the success rates were tested with
257 pairwise Mann-Whitney-Wilcoxon tests with Bonferroni post-hoc corrections
258 since data did not fit a normal distribution according to Kolmogorov-Smirnov tests
259 with Lilliefors corrections ($\alpha=0.05$).

260 **2.4. Variation Index of each pollen/spore type**

261 The same pollen/spore type can show different curve profiles in different locations
262 depending on the abundance of the emission sources, wind dynamics, climate,
263 meteorological conditions, and phenology of the species present in the territory
264 (Grinn-Gofroń and Rapiejko, 2009; Picornell et al., 2019b; Velasco-Jiménez et
265 al., 2013). Therefore, results of a certain pollen/spore type may not be directly
266 comparable among sampling locations. Moreover, interpolation success may be
267 strongly related with the curve shape and the variation coefficients in the
268 concentrations of consecutive days (noise). In order to characterise the daily
269 variations of a given pollen/spore type in the different sampling stations, we
270 developed the so defined “Variation Index” (VIn). It consisted on calculating the
271 average of the variation coefficients (CV, equation 3) of every two consecutive
272 days for the main pollen/spore season (equation 4). The average VIn is then
273 calculated for the years included in the study. This index measures the average
274 variations during consecutive days that a certain pollen/spore type shows at a

275 certain locality (i.e. the more variations in the concentrations among consecutive
276 days, the highest VIn is obtained).

277 (3)

$$278 \quad CV = \frac{\sigma}{\bar{X}}$$

279 where CV is the coefficient of variation, σ the standard deviation and \bar{X} the
280 average.

281 (4)

$$282 \quad \text{Variation Index} = \frac{\sum_{i=1}^n CV_i}{n}$$

283 where CV is the coefficient of variation, and n the number of days within the main
284 pollen/spore season.

285 To check if there is any relationship between the VIn of a certain pollen/spore
286 type and the errors obtained when interpolating, a linear regression was
287 calculated between these two variables.

288 **3. Results and discussion**

289 **3.1. Relative errors and Variation Index**

290 Once gaps of 3, 5, 7 and 10 days were artificially created in the data series of
291 different pollen/spore types, years and localities, and the gaps filled with the
292 estimated data, the results obtained were analysed by comparing them with the
293 observed data. In general, regarding the different interpolation methods, the one
294 that obtained the lowest relative error was the moving mean (0.77 as average),
295 followed by the linear interpolation (0.80 as average) (Fig. 3A). Their relative
296 errors showed significant differences between them as well as with the other

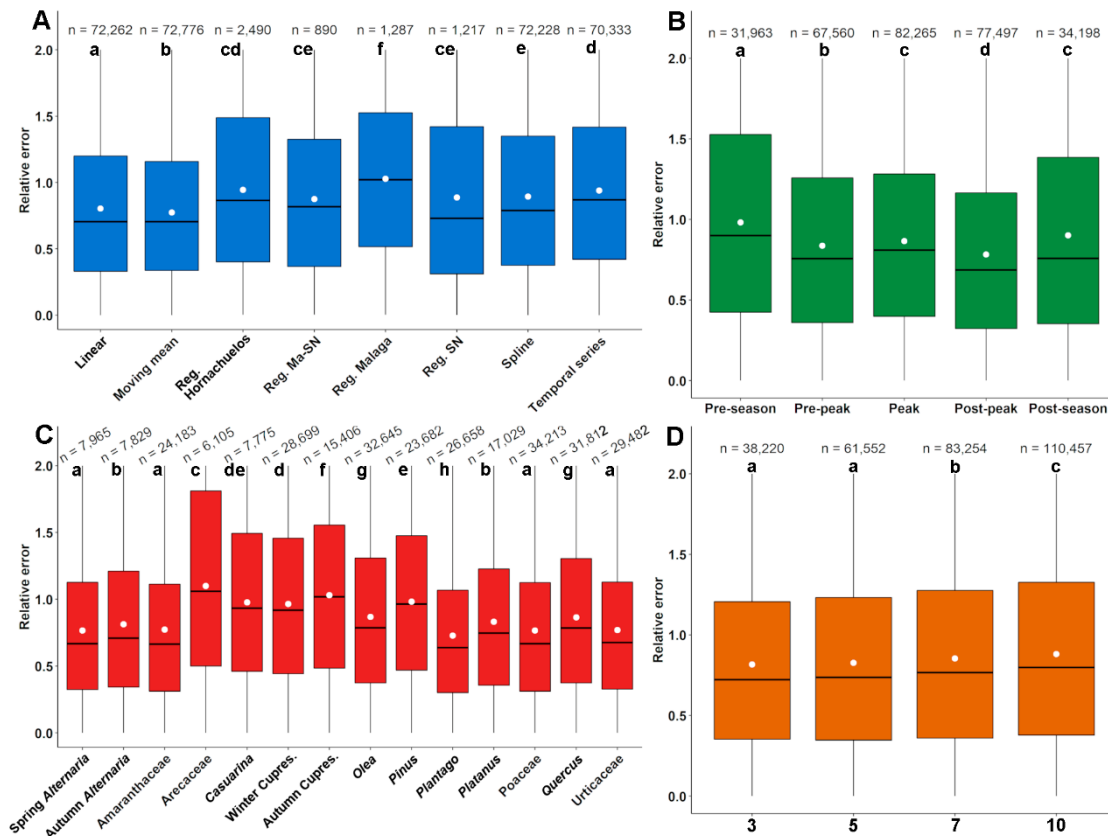
297 methods. These results are due to, despite its mathematical simplicity, the
298 moving mean takes into account the curve trend and the pollen/spore
299 concentrations immediately before and after the gaps, what provides a more
300 accurate adjustment to the pollen/spore curve. Although abrupt changes may
301 happen during the missing period, this interpolation method follows the general
302 trend of the serie and, in general, it is less likely that predictions contain major
303 errors. The linear interpolation generally takes into account the curve trend too,
304 but in some cases it is oversimplified and the new data obtained may be affected
305 by punctual concentrations that do not fit the general trend.

306 The regressions with nearby locations, the spline regresions, and the regressions
307 with nearby locations obtained significant higher relative errors according to
308 Mann-Whitney-Wilcoxon tests (Fig. 3A).

309 The interpolation with nearby location obtained very different relative errors
310 depending on the nearby sampling station considered. In the case of Ronda, the
311 regressions with Sierra de las Nieves sampling station (14 km away) obtained
312 lower relative errors (RE=0.89 as average) than the regressions with Malaga
313 sampling station (RE=1.03), situated 62 km away (Fig. 3A). However, when
314 Sierra de las Nieves and Malaga databases were simultaneously taken into
315 account, the error rates were estatistically similar to those obtained when
316 considering only Sierra de las Nieves. These last errors were also similar to the
317 ones obtained in the case of Cordoba when using Hornachuelos Natural Park
318 sampling station as neighbour location, which is 64 km apart, with a RE of 0.94,
319 as average. This interpolation method is the only one of all tested that, indirectly,
320 takes into account variables such as meteorological conditions, the effect of the
321 vegetation or land use. The effects of these variables are also reflected in the

322 pollen/spore daily concentrations of the nearby location, so they are indirectly
323 integrated in the regression with the target location. Therefore, more accurate
324 interpolations could be expected if the nearby locations had similar climatic
325 conditions, vegetation and land use, which would be also reflected in the
326 pollen/spore timing (phenology) and in the airborne pollen/spore load (intensity)
327 (El-Moslimany, 2019; García-Mozo, 2017; Ruiz-Valenzuela and Aguilera, 2018).
328 For these reasons, closer localities are generally more likely to have similar
329 concentration curves and, therefore, regressions analysis resulting more
330 accurate than when using further away locations, as suggested in previous
331 studies (Hjort et al., 2016; Lara et al., 2020; Navares and Aznarte, 2019).
332 Nevertheless, it is possible that further sampling sites with similar conditions to
333 the target station or with similar ornamental taxa in the vicinity of the pollen trap
334 obtain lower errors than geographically closer sampling sites. Consequently, the
335 errors obtained for the method of the nearby locations should be cautelously
336 considered because the accuracy of the interpolation depends on the factors
337 aforementioned. When applying this method, it would be interesting to select the
338 nearby location by studing its similarity to the target location as proposed by
339 Oteros et al. (2019).

340 Usually, new aerobiological sampling sites are selected in order to cover areas
341 with different environmental conditions than other previously settled stations,
342 including meteorological conditions, land use and vegetal coverage.
343 Consequently, the results of the interpolation with nearby locations, as observed
344 in the results, are expected to produce high relative errors as average, given that
345 the different sampling stations are installed to cover the geographical
346 heterogeneity of a territory.



347

348 **Fig. 3.** Relative errors obtained by the different interpolation methods (A), periods
 349 of the year (B), pollen/spore types (C), and gap sizes (days; D) with all the
 350 available data. n: number of observations. Each box includes the interquartile
 351 range (Q1-Q3), bold lines indicate the median and white dots indicate the mean.
 352 Groups which share the same letter above have not any significant differences
 353 ($\alpha=0.05$) between them according to Mann-Whitney-Wilcoxon tests with
 354 Bonferroni *post-hoc* corrections. Reg.: regression with; Ma-SN: Malaga and
 355 Sierra de las Nieves; SN: Sierra de las Nieves.

356 The spline interpolation, as well as the moving mean and the linear interpolations,
 357 also takes into account the curve trend, but they produce more pronounced curve
 358 trends than the other methods as a consequence of spline approximation. This
 359 can lead to very accurate fits or to big errors in the predictions, which, in general,

360 gives higher errors (0.89) than the linear and moving mean interpolations (Fig.
361 3A).

362 The temporal series analysis is highly dependent on the extension of the historical
363 database (number of years in this case). Therefore, if the number of years were
364 not enough (as occurred in Ronda or in Sierra de las Nieves), the obtained
365 seasonality curve might not be representative of the regular behaviour of a given
366 pollen/spore type. Also, data series of uncommon years (according to phenology
367 and flowering intensity) deviated from the standard behaviour may result in non-
368 accurate interpolations when using this method (Belda et al., 2020). In wind
369 pollinated trees, remarkable differences between pollen seasons of consecutive
370 years are frequent due to mast seeding cycles (Bogdziewicz et al., 2017).
371 Additionally, when performing the linear regression between the seasonality
372 curve and the curve of the target year, more errors are accumulated (0.94 as
373 average relative error).

374 Regarding the main pollen/spore season (Fig. 3B), in general, the period that
375 obtained the lowest relative errors was the post-peak (RE average=0.78),
376 followed by the pre-peak period (RE average=0.84). Post-peak periods usually
377 present smoother curve shapes than the pre-peaks and the peaks periods due to
378 these last ones are more conditioned by plant phenology and flowering intensity
379 than the post-peaks, in which the plants progressively reduce the pollen emission
380 intensity (Cunha et al., 2016; Kasprzyk and Walanus, 2014; Picornell et al.,
381 2019a). Therefore, as abrupt changes in the data series are less likely during the
382 post-peaks, fewer errors are also expected in the interpolation.

383 On the other hand, the peak-day concentration is difficult to predict, since it is
384 usually an abrupt change caused by the interaction of both meteorological and

385 biological parameters which are not easily predictable. Moreover, the peak-day
386 concentration usually varies widely from one year to another, what makes more
387 difficult to successfully apply interpolation methods (Devadas et al., 2018; García-
388 Mozo et al., 2009; Picornell et al., 2019a; Valencia et al., 2019). Consequently,
389 the relative errors obtained were higher than for the rest of the MPS/MSS period
390 (RE average=0.86).

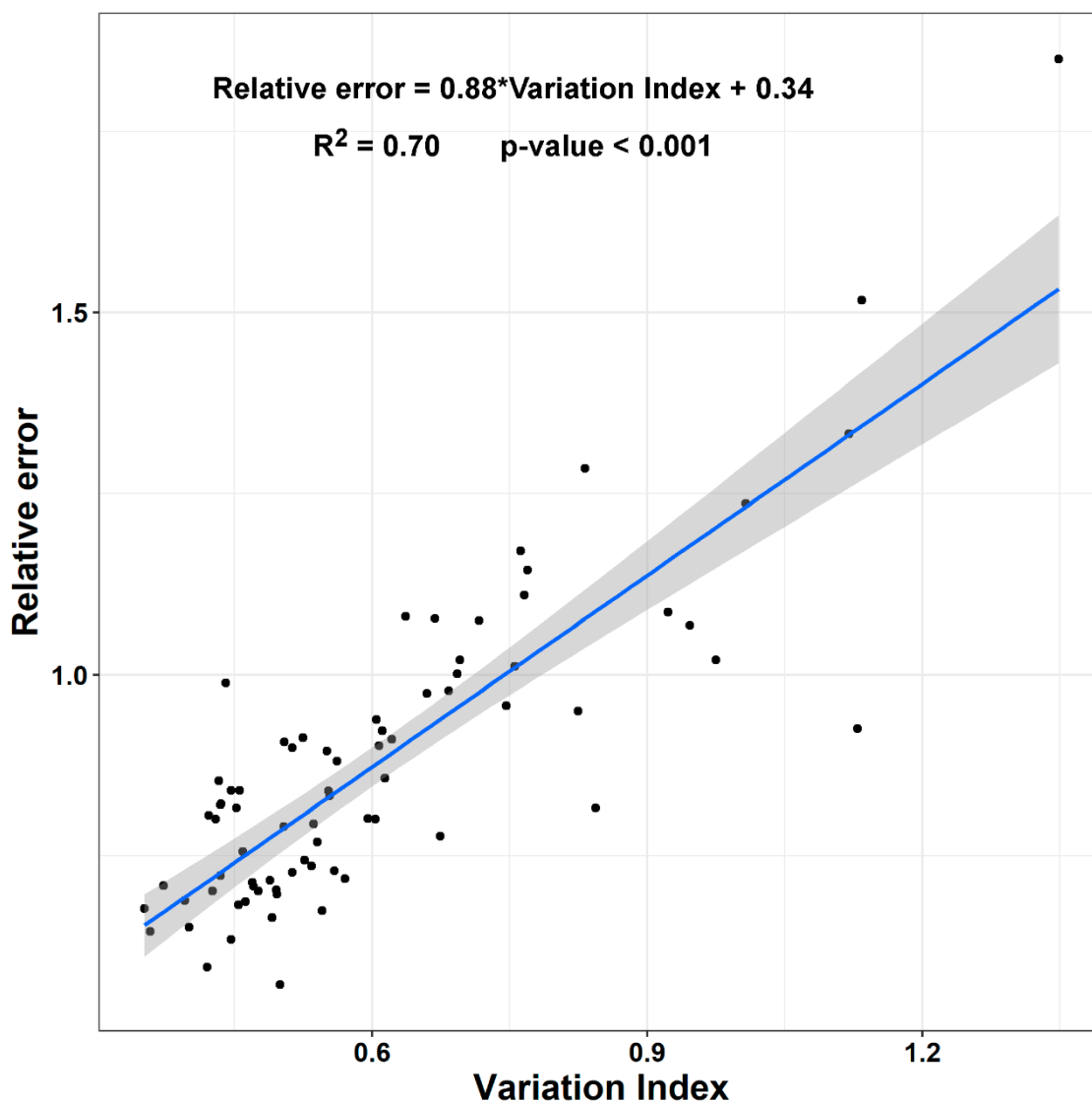
391 Outside the MPS/MSS, pre- and post-seasonal periods obtained significant
392 higher relative errors (0.98 and 0.90 as average respectively) (Fig. 3B). During
393 these periods, days with null value in pollen/spore concentrations are frequent,
394 interspersed with small rises and falls, what makes it more difficult to predict or
395 estimate the daily values. However, these errors might not be relevant for defining
396 the MPS/MSS unless they are located near to the start or the end dates. Given
397 that the concentrations outside the MPS/MSS usually are very low, such errors
398 are less relevant for allergy alerts.

399 Gap sizes longer than 5 days obtained significantly higher relative errors (Fig.
400 3D). In these gaps, abrupt changes or trend changes are more likely to happen
401 than in smaller gaps, what may lead to higher errors. Accordingly to this, the
402 lowest error rate was obtained for gaps of 3 days (0.82), but with non-significant
403 differences with the errors for gaps of 5 days (0.83). As expected, longer gaps
404 produced higher errors since the uncertainty increases when it comes to
405 estimating longer periods. In fact, relative errors of 0.85 and 0.88 were obtained
406 for gaps of 7 and 10 days respectively. However, these errors are expected to
407 induce less changes in the MPS/MSS calculation than leaving the gaps without
408 data, so it is still recommendable to interpolate them.

409 Regarding the results obtained by pollen/spore types (Fig. 3C), spring *Alternaria*
410 (0.76), Amaranthaceae (0.77), Poaceae (0.77) and Urticaceae (0.77) were the
411 pollen and spore types that obtained the lowest relative errors, as average. In
412 general, these pollen and spore types are integrated by several species that,
413 usually, have wide distribution areas. These pollen/spore types are detected
414 during a relatively long period of the year and it probably makes their seasonal
415 trends be smoother than the other pollen types such as Arecaceae (in which the
416 highest relative errors were obtained, RE=1.10), *Casuarina* (0.98) or *Platanus*
417 (0.83). Other pollen types, such as Cupressaceae, *Pinus*, *Quercus*, and *Olea*,
418 have several concentration peaks within their MPS, which may correspond to the
419 flowering of the different species or varieties that integrate the pollen types. These
420 peaks are difficult to predict, and it may increase the relative errors obtained. For
421 the same reason, autumn *Alternaria* obtained higher relative errors than spring
422 *Alternaria* because autumn MSS is usually shorter and contains more
423 pronounced peaks.

424 Despite the pollen type has been considered, it is more interesting to consider
425 the behaviour of the pollen curve in general, which have been characterized in
426 this case by the Variation Index (VIn, see methods for the definition). The same
427 pollen/spore type showed different relative errors at different sampling locations
428 (data not shown). Therefore, it would be pointless to establish the average
429 relative error by pollen/spore type if it is going to vary when considering a new
430 sampling site. As observed in Fig. 4, the more variations during consecutive days
431 (higher VIn), the more relative errors are obtained during interpolation, a linear
432 and direct relationship existing between the VIn and the relative errors. So, by
433 means of the regression equation, included in Fig. 5, it is possible to estimate the

434 average error rate when interpolating values of a pollen/spore type by calculating
435 the VIn. This error estimation is independent of the pollen/spore type, and only
436 relies in the behaviour of its daily concentration curve. Furthermore, this
437 regression equation has been elaborated with pollen and spore data of 6
438 sampling sites located at different environmental conditions and so, it can be used
439 as a calibration curve for estimating the errors at new locations. Anyway, we
440 recommend taking interpolation results with great caution when the VIn is higher
441 than 0.75, since relative errors greater than 1 are expected, as average, above
442 this value.

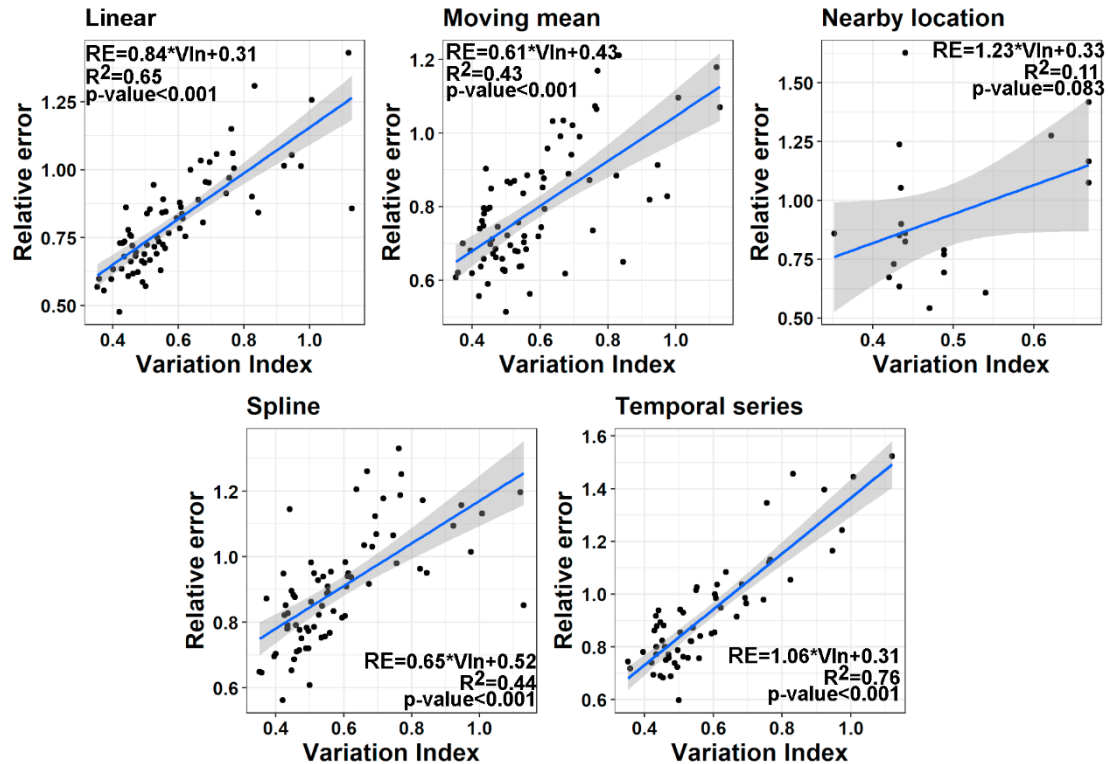


443

444 **Fig. 4.** Linear regression between the Variation Index and the relative error of the
445 different pollen/spore type in the different sampling localities during the
446 MPS/MSS. The grey area marks the 95% confidence interval.

447 Additionally, in Fig. 5 we have represented, separately, the regression lines
448 between the Variation Indexes and the relative errors for each interpolation
449 method. It would allow to roughly estimate the average relative error that this
450 interpolation method would produce for a given pollen/spore type. The methods
451 that obtained the highest coefficient of regression were the temporal series, linear
452 and spline interpolation. The regression with nearby location is not statistically
453 significant since the relative errors obtained depend on the similarity between
454 sampling sites, and not on the VIn of each pollen/spore type.

455 According to the obtained results, some interpolation methods, such as temporal
456 series and linear interpolations, are more sensitive to pollen/spore types with high
457 variations in their concentrations during consecutive days than the others (Fig.
458 5). This can be observed in the regression equations slopes that, when
459 significant, are higher than in the other methods. Linear interpolation may obtain
460 lower relative errors than moving mean interpolation if the pollen/spore type had
461 a low VIn, but the errors would be higher if the pollen type presented a higher
462 VIn. In the case of the interpolation with data from nearby location, the points did
463 not fit a linear regression ($p\text{-value} > 0.05$). It can be explained, as commented
464 above, because the errors obtained during the interpolation are related to the
465 similarity of both sampling locations, rather than to the characteristics of the
466 pollen/spore type.



467

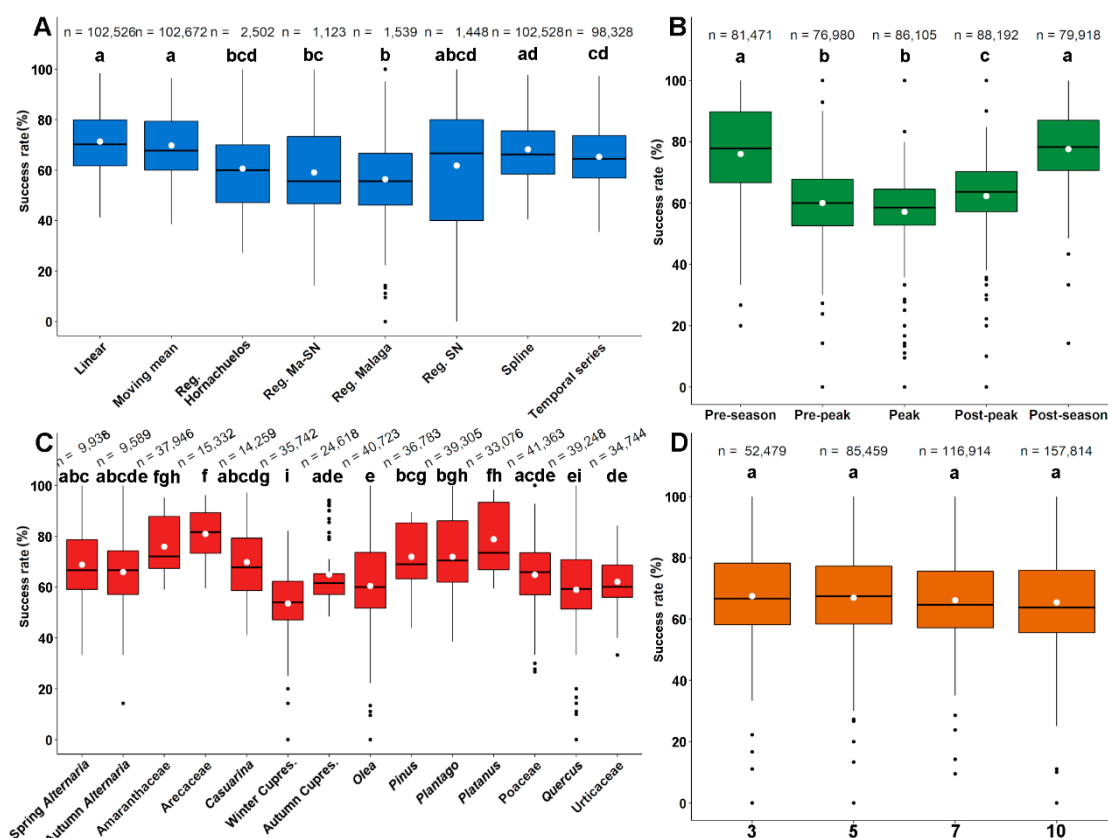
468 **Fig. 5.** Linear regressions between the Variation Index and the relative error of
469 each pollen/spore type during its MPS/MSS sorted by interpolation method. The
470 grey area marks the 95% confidence interval. RE: relative error; VIn: Variation
471 Index.

472 3.2. Success rates

473 Observed and estimated pollen/spore concentrations were categorized by
474 thresholds following the criteria of the Spanish Aerobiology Network (REA) for
475 each pollen/spore type. Then, observed and estimated categories were
476 compared and, in general, success rates above 60% were obtained for all the
477 studied bioaerosols (Fig. 6). The REA and other pollen information platforms use
478 the categories null/nil, low, moderate and high for releasing the pollen risk
479 information to the population (Galán et al., 2007; Pérez-Badía et al., 2010).

480 Therefore, many of the concentrations that showed relative errors in the
 481 continuous variable, now are classified in the same risk category.

482 The highest average success rate (i.e. lowest errors) were obtained, once more,
 483 for the linear interpolation (71%) and the moving mean interpolation (70%),
 484 without significant differences between them, but neither with the spline
 485 interpolation (68%), or the regression with Sierra de las Nieves (62%; in the case
 486 of Malaga sampling location) (Fig. 6A). The lowest average success rates (i.e.
 487 highest errors) were detected when the interpolation was performed by using the
 488 nearby location of Malaga (56% in the case of Ronda). In general, these results
 489 were similar to the obtained by calculating the relative errors (Fig. 3A) with the
 490 exception that, in this case, using levels instead daily concentrations, the spline
 491 interpolation and the regression with Sierra de las Nieves did not show significant
 492 differences when compared to linear and moving mean methods.



493

494 **Fig. 6.** Success rates obtained by the different interpolation methods (A), periods
495 of the year (B), pollen/spore types (C), and gap sizes (days; D) when comparing
496 the observed and predicted pollen/spore levels. n: number of observations. Each
497 box includes the interquartile range (Q1-Q3), bold lines indicate the median and
498 white dots indicate the mean. Groups which share the same letter above have
499 not any significant differences ($\alpha=0.05$) among them according to Mann-Whitney-
500 Wilcoxon tests with Bonferroni *post-hoc* corrections. Reg.: regression with; Ma-
501 SN: Malaga and Sierra de las Nieves Natural Park; SN: Sierra de las Nieves
502 Natural Park; Hornachuelos: Hornachuelos Natural Park.

503 Regarding the periods of the pollen season (Fig. 6B), those outside the MPS/MSS
504 (i.e. pre-season and post-season) obtained the highest success rates (76 and
505 78% respectively). As previously commented, pollen and spore concentrations
506 during these periods are generally low and, although interpolations produced high
507 relative errors, the variations between expected and observed concentrations
508 imply very little changes in the pollen/spore categories. Therefore, these errors
509 are less relevant for the allergy alert systems since they do not imply high
510 changes in the information of the atmospheric allergenic potential. The highest
511 success rates inside the MPS/MPS were obtained for the post-peak period (62%
512 of success), what matches the results of the relative errors, while the lowest
513 success rates were obtained for the peak (57%), and pre-peak (60%) periods,
514 once again the peak being the most unpredictable period.

515 As observed with the relative errors, there is a decrease in the average success
516 rate (i.e. higher errors) when gaps of more than 5 days are interpolated (Fig. 6D).
517 However, in this case, these errors did not involve any significant difference
518 between any gap size since all average success rates were between 65 and 67%.

519 Finally, as can be seen in Fig. 6C, the different pollen/spore types obtained
520 different success rates, ranged from 54 to 81%. The highest average success
521 rate was obtained for Arecaceae (81%). This pollen type did not present
522 significant differences with Amaranthaceae (76%) or *Platanus* (79%). Usually,
523 Arecaceae pollen concentrations detected are low, as occurred also with
524 *Platanus* in most sampling locations. The errors when interpolating such low
525 concentrations may involve high relative errors but little changes in the established
526 categories, a similar effect that the observed outside the MPS/MSS periods and
527 commented above (Fig. 3B). The lowest average success rate was obtained for
528 winter Cupressaceae (54%), followed by *Quercus* (59%), *Olea* (60%), Urticaceae
529 (62%), and autumn Cupressaceae (65%). These pollen types are usually
530 detected in a wider range of concentrations than in the other pollen types and so,
531 errors during the interpolation are more likely to entail errors in the categories
532 and, therefore, lower accuracy rates (Fig. 6C).

533 Although the results obtained in some cases have not been the most favourable,
534 we consider that they have been accurate enough (relative errors are generally
535 under 0,8) for not leaving blank the gaps in an aerobiological database, without
536 assigning a value, due to it would lead to take these concentrations as 0 pollen
537 grains or spores/m³. This would introduce higher errors in the annual spore/pollen
538 integral and in the MPS/MSS definition than with the interpolated data. These
539 errors are potentially greater when a percentage definition of the MPS/MSS is
540 applied. Moreover, when working with pollen/spore levels during these missing
541 days, the accuracy can reach to the 70-71% of the cases (using moving mean or
542 linear interpolation, respectively), which would allow to use these data to give

543 pollen/spore information to the population or to make comparisons between
544 pollen data and allergic symptoms (Karatzas, 2009).

545 Data quality and errors involved in the aerobiological sampling method may also
546 play an important role in the measurement of interpolation accuracy (Oteros et
547 al., 2015; Rojo et al., 2019). If these errors increase the variability of the
548 bioaerosol concentrations during consecutive days, they might compromise the
549 measurement of the interpolation accuracy (i.e., they will increase the VIn).
550 However, this effect is not easily measurable since the data used for validating
551 the interpolation provides from the same pollen trap and would have the same
552 potential sampling error. In these terms, the interpolation methods that does not
553 only rely on the data in both sides of the gap, such as the interpolation by using
554 temporal series analysis or the interpolation with nearby locations, would be less
555 affected by sampling errors.

556 Apart from the methods proposed in this study, additional methods to complete
557 missing data may be considered, such as elaborating regional forecast models
558 based on meteorological variables and emission maps or dispersal models (Lara
559 et al., 2019; Verstraeten et al., 2021). Nevertheless, such models should be
560 elaborated separately for each pollen/spore type, and it might be necessary to
561 elaborate individual models for different climate areas (García-Mozo et al., 2008).
562 Hence, such methods would not be easy to automatize, they require individual
563 validation, depend on the availability of meteorological data for the target location
564 and, in most cases, a long time series of data is required to train and validate the
565 models.

566 Due to recent movility restrictions caused by the COVID-19 pandemic, many
567 aerobiological samplings have been interrupted to a lesser or greater extend. This

568 has caused missing data for some weeks and even months in several
569 aerobiological sampling stations. Most of these monitoring gaps occurred during
570 spring, which affected the MPS/MSS data collection. As observed in the results,
571 when the gap is longer than 5 consecutive days, the error rates increase
572 significantly. For these long gaps, most of the presented interpolation methods
573 might not be appropriate, so, it would be interesting for further studies to test the
574 performance of time series analysis or the regressions with nearby sampling
575 stations when a great part of the data of the MPS/MSS is missing. Nevertheless,
576 predicting the temporality and intensity of the MPS/MSS is not an easy task, and
577 often requires the adjustment of the models to the local conditions (Picornell et
578 al., 2019a; Rojo et al., 2021, 2016).

579 This work is an approach to perform interpolations in order to fill in the gaps that,
580 due to different reasons, are generated in pollen/spore databases. Despite the
581 most accurate method was generally the moving mean, for each specific case it
582 would be necessary to select the method according to the particularities of each
583 sampling station and pollen/spore type. Although this study was conducted with
584 aerobiological databases, these results may be useful for interpolating missing
585 data in other environmental databases.

586 **4. Conclusions**

- 587 • The moving mean interpolation is the method that generates the lowest
588 relative errors, as average. This method is independent of the availability of
589 additional data and of the length of the database, and it is also less sensible
590 to variations in the pollen/spore concentrations during consecutive days than
591 the other methods considered in this study. In addition, this method showed

592 a success rate of the 70% when assigning the risk classes that are frequently
593 used in the allergy alert systems.

594 • Periods with high variation indexes (VIn) make the pollen concentrations
595 difficult to predict and, generally, cause high errors when interpolating data.
596 Probably that be the reason why the pre-peak and peak periods present
597 higher error rates.

598 • The Variation Index proposed, based on the pollen/spore season behaviour,
599 is a good indicator of the success rate. Therefore, it is advisable to take this
600 index into consideration since it allows to estimate the relative error before
601 applying interpolation methods.

602 • The errors during the interpolation generally increase when gaps of more than
603 5 days are considered. For that reason, alternative methods should be
604 considered for interpolating longer gaps.

605 **Funding:** This work was supported by the Spanish Ministry of Economy and
606 Competitiveness [project CGL2014-54731-R]; by the Ministry of Science and
607 Innovation [projects RTI2018-096392-B-C22]; by the Junta de Andalucía
608 [contract 8.06/503.4764]; and by the Area of Environment and Sustainability of
609 the Malaga City Council [contracts 8.06/5.03.4721 and 8.07/5.03.5159], and the
610 Junta Comunidades de Castilla-La Mancha, which provides financial support for
611 the Castilla-La Mancha Aerobiology Network (AEROCAM). Antonio Picornell was
612 supported by a predoctoral grant financed by the Spanish Ministry of Education,
613 Culture and Sport, in the Program for the Promotion of Talent and its
614 Employability [grant number FPU15/01668]. The pollen trap installed in Sierra de
615 las Nieves was funded by the Herbarium MGC of the SCAI (Central Services of
616 Research Support) of the University of Malaga under the agreement signed

617 between the Junta de Andalucía and the University of Malaga [contract
618 8.07/5.034764].

619 **Acknowledgments:** The authors specially want to thanks the SCAI (Central
620 Service for Research Support) of the University of Malaga for supporting the
621 acquisition of the pollen trap installed in Sierra de las Nieves; the Parauta City
622 Council, the direction of Sierra de las Nieves Natural Park, Las Conejeras
623 campsite for facilitating the installation of the pollen trap in Sierra de las Nieves;
624 and the staff of Pérez de Guzmán High School for providing support to install and
625 maintain the pollen trap in Ronda, and to Enresa for facilitating the installation
626 and maintenance of the pollen trap in Hornachuelos Natural Park.

627 **Conflicts of Interest:** The authors declare no conflict of interest. The funders
628 had no role in the design of the study, collection, analyses, or interpretation of
629 data; in the writing of the manuscript, or in the decision to publish the results.

630 **References**

631 Belda, S., Pipia, L., Morcillo-Pallarés, P., Rivera-Caicedo, J.P., Amin, E., De
632 Grave, C., Verrelst, J., 2020. DAtimeS: A machine learning time series GUI
633 toolbox for gap-filling and vegetation phenology trends detection. *Environ.*
634 *Model. Softw.* 127, 104666. <https://doi.org/10.1016/j.envsoft.2020.104666>

635 Belmonte, J., Canela, M., Guàrdia, R.A., Sbai, L., Vendrell, M., Cariñanos, P.,
636 Díaz de la Guardia, C., Dopazo, A., Fernández, D., Gutiérrez, M., Trigo,
637 M.M., Guàrdia, R.A., Sbai, L., Vendrell, M., Cariñanos, P., Díaz de la
638 Guardia, C., Dopazo, A., Fernández, D., Gutiérrez, M., Trigo, M.M., 1999.
639 *Aerobiological dynamics of the Cupressaceae pollen in Spain, 1992-98.*
640 *Polen* 10, 27–38.

- 641 Bogdziewicz, M., Szymkowiak, J., Kasprzyk, I., Grewling, Ł., Borowski, Z.,
642 Borycka, K., Kantorowicz, W., Myszkowska, D., Piotrowicz, K., Ziemianin,
643 M., Pesendorfer, M.B., 2017. Masting in wind-pollinated trees: System-
644 specific roles of weather and pollination dynamics in driving seed production.
645 *Ecology* 98, 2615–2625. <https://doi.org/10.1002/ecy.1951>
- 646 Cleveland, R.B., Cleveland, W.S., McRae, J.E., Terpenning, I., 1990. STL: A
647 Seasonal-Trend Decomposition Procedure based on Loess. *J. Off. Stat.* 6,
648 3–73.
- 649 Cunha, M., Ribeiro, H., Abreu, I., 2016. Pollen-based predictive modelling of wine
650 production: Application to an arid region. *Eur. J. Agron.* 73, 42–54.
651 <https://doi.org/10.1016/j.eja.2015.10.008>
- 652 Cunha, M., Ribeiro, H., Costa, P., Abreu, I., 2015. A comparative study of
653 vineyard phenology and pollen metrics extracted from airborne pollen time
654 series. *Aerobiologia*. 31, 45–56. <https://doi.org/10.1007/s10453-014-9345-3>
- 655 Devadas, R., Huete, A.R., Vicendese, D., Erbas, B., Beggs, P.J., Medek, D.,
656 Haberle, S.G., Newnham, R.M., Johnston, F.H., Jaggard, A.K., Campbell,
657 B., Burton, P.K., Katelaris, C.H., Newbigin, E., Thibaudon, M., Davies, J.M.,
658 2018. Dynamic ecological observations from satellites inform aerobiology of
659 allergenic grass pollen. *Sci. Total Environ.* 633, 441–451.
660 <https://doi.org/10.1016/j.scitotenv.2018.03.191>
- 661 El-Moslimany, A., 2019. Reduced Poaceae pollen under conditions of severe
662 summer drought in the Middle East: Implications for rainfall seasonality in
663 pollen diagrams. *Rev. Palaeobot. Palynol.* 271, 104068.
664 <https://doi.org/10.1016/j.revpalbo.2019.04.007>

- 665 Fröhlich-Nowoisky, J., Kampf, C.J., Weber, B., Huffman, J.A., Pöhlker, C.,
666 Andreae, M.O., Lang-Yona, N., Burrows, S.M., Gunthe, S.S., Elbert, W., Su,
667 H., Hoor, P., Thines, E., Hoffmann, T., Després, V.R., Pöschl, U., 2016.
668 Bioaerosols in the Earth system: Climate, health, and ecosystem
669 interactions. *Atmos. Res.* 182, 346–376.
670 <https://doi.org/10.1016/j.atmosres.2016.07.018>
- 671 Gabarra, E., Belmonte, J., Canela, M., 2002. Aerobiological behaviour of
672 *Platanus L.* pollen in Catalonia (North-East Spain). *Aerobiologia*. 18, 185–
673 193. <https://doi.org/10.1023/A:1021370724043>
- 674 Galán, C., Ariatti, A., Bonini, M., Clot, B., Crouzy, B., Dahl, A., Fernández-
675 González, D., Frenguelli, G., Gehrig, R., Isard, S., Levetin, E., Li, D.W.,
676 Mandrioli, P., Rogers, C.A., Thibaudon, M., Sauliene, I., Skjoth, C., Smith,
677 M., Sofiev, M., 2017. Recommended terminology for aerobiological studies.
678 *Aerobiologia*. 33, 293–295. <https://doi.org/10.1007/s10453-017-9496-0>
- 679 Galán, C., Cariñanos, P., Alcázar, P., Domínguez-Vilches, E., 2007. Spanish
680 Aerobiology Network (REA): Management and Quality Manual. Servicio de
681 Publicaciones Universidad de Córdoba, Córdoba.
- 682 Galán, C., Smith, M., Damialis, A., Frenguelli, G., Gehrig, R., Grinn-Gofroñ, A.,
683 Kasprzyk, I., Magyar, D., Oteros, J., Šaulienė, I., Thibaudon, M., Sikoparija,
684 B., 2021. Airborne fungal spore monitoring: between analyst proficiency
685 testing. *Aerobiologia*. 1–11. <https://doi.org/10.1007/s10453-021-09698-4>
- 686 Galán, C., Smith, M., Thibaudon, M., Frenguelli, G., Oteros, J., Gehrig, R.,
687 Berger, U., Clot, B., Brandao, R., 2014. Pollen monitoring: minimum
688 requirements and reproducibility of analysis. *Aerobiologia*. 30, 385–395.

- 689 <https://doi.org/10.1007/s10453-014-9335-5>
- 690 García-Mozo, H., 2017. Poaceae pollen as the leading aeroallergen worldwide:
691 A review. *Allergy* 72, 1849–1858. <https://doi.org/10.1111/all.13210>
- 692 García-Mozo, H., Chuine, I., Aira, M.J., Belmonte, J., Bermejo, D., Díaz de la
693 Guardia, C., Elvira, B., Gutiérrez, M., Rodríguez-Rajo, J., Ruiz, L., Trigo,
694 M.M., Tormo-Molina, R., Valencia, R., Galán, C., 2008. Regional
695 phenological models for forecasting the start and peak of the Quercus pollen
696 season in Spain. *Agric. For. Meteorol.* 148, 372–380.
697 <https://doi.org/10.1016/j.agrformet.2007.09.013>
- 698 García-Mozo, H., Dominguez-Vilches, E., Galan, C., 2007. Airborne allergenic
699 pollen in natural areas: Hornachuelos Natural Park, Cordoba, Southern
700 Spain. *Ann. Agric. Environ. Med.* 14, 63–69.
- 701 García-Mozo, H., Galán, C., Belmonte, J., Bermejo, D., Candau, P., Díaz de la
702 Guardia, C., Elvira, B., Gutiérrez, M., Jato, V., Silva, I., Trigo, M.M., Valencia,
703 R., Chuine, I., 2009. Predicting the start and peak dates of the Poaceae
704 pollen season in Spain using process-based models. *Agric. For. Meteorol.*
705 149, 256–262. <https://doi.org/10.1016/j.agrformet.2008.08.013>
- 706 García-Mozo, H., Pérez-Badía, R., Fernández-González, F., Galán, C., 2006.
707 Airborne pollen sampling in Toledo, Central Spain. *Aerobiologia.* 22, 55–66.
708 <https://doi.org/10.1007/s10453-005-9015-6>
- 709 Grinn-Gofroń, A., Rapiejko, P., 2009. Occurrence of *Cladosporium* spp. and
710 *Alternaria* spp. spores in Western, Northern and Central-Eastern Poland in
711 2004–2006 and relation to some meteorological factors. *Atmos. Res.* 93,
712 747–758. <https://doi.org/10.1016/J.ATMOSRES.2009.02.014>

- 713 Hernández-Ceballos, M.A., García-Mozo, H., Galán, C., 2015. Cluster analysis
714 of intradiurnal holm oak pollen cycles at peri-urban and rural sampling sites
715 in southwestern Spain. *Int. J. Biometeorol.* 59, 971–982.
716 <https://doi.org/10.1007/s00484-014-0910-9>
- 717 Hirst, J.M., 1952. An automatic volumetric spore trap. *Ann. Appl. Biol.* 39, 257–
718 265. <https://doi.org/10.1111/j.1744-7348.1952.tb00904.x>
- 719 Hjort, J., Hugg, T.T., Antikainen, H., Rusanen, J., Sofiev, M., Kukkonen, J.,
720 Jaakkola, M.S., Jaakkola, J.J.K., 2016. Fine-Scale Exposure to Allergenic
721 Pollen in the Urban Environment: Evaluation of Land Use Regression
722 Approach. *Environ. Health Perspect.* 124, 619–626.
723 <https://doi.org/10.1289/ehp.1509761>
- 724 Junger, W.L., Ponce de Leon, A., 2015. Imputation of missing data in time series
725 for air pollutants. *Atmos. Environ.* 102, 96–104.
726 <https://doi.org/10.1016/j.atmosenv.2014.11.049>
- 727 Junta de Andalucía, 2011. Red de Información Ambiental de Andalucía
728 (REDIAM) [WWW Document]. URL
729 <https://www.juntadeandalucia.es/medioambiente/site/rediam> (accessed
730 2.6.20).
- 731 Karatzas, K.D., 2009. Informing the public about atmospheric quality: Air pollution
732 and pollen. *Allergo J.* 18, 212–217. <https://doi.org/10.1007/BF03362059>
- 733 Kasprzyk, I., Walanus, A., 2014. Gamma, Gaussian and logistic distribution
734 models for airborne pollen grains and fungal spore season dynamics.
735 *Aerobiologia.* 30, 369–383. <https://doi.org/10.1007/s10453-014-9332-8>

- 736 Lara, B., Rojo, J., Fernández-González, F., González-García-Saavedra, A.,
737 Serrano-Bravo, M.D., Pérez-Badía, R., 2020. Impact of Plane Tree
738 Abundance on Temporal and Spatial Variations in Pollen Concentration.
739 *Forests* 11, 817. <https://doi.org/10.3390/f11080817>
- 740 Lara, B., Rojo, J., Fernández-González, F., Pérez-Badía, R., 2019. Prediction of
741 airborne pollen concentrations for the plane tree as a tool for evaluating
742 allergy risk in urban green areas. *Landsc. Urban Plan.* 189, 285–295.
743 <https://doi.org/10.1016/j.landurbplan.2019.05.002>
- 744 Lehmann, T.M., Gönner, C., Spitzer, K., 1999. Survey: Interpolation methods in
745 medical image processing. *IEEE Trans. Med. Imaging* 18, 1049–1075.
746 <https://doi.org/10.1109/42.816070>
- 747 Luedeling, E., Kunz, A., Blanke, M.M., 2013. Identification of chilling and heat
748 requirements of cherry trees-a statistical approach. *Int. J. Biometeorol.* 57,
749 679–689. <https://doi.org/10.1007/s00484-012-0594-y>
- 750 Navares, R., Aznarte, J.L., 2019. Geographical imputation of missing poaceae
751 pollen data via convolutional neural networks. *Atmosphere (Basel)*. 10, 717–
752 727. <https://doi.org/10.3390/atmos10110717>
- 753 Navares, R., Aznarte, J.L., 2017. Predicting the Poaceae pollen season: six
754 month-ahead forecasting and identification of relevant features. *Int. J.*
755 *Biometeorol.* 61, 647–656. <https://doi.org/10.1007/s00484-016-1242-8>
- 756 Nilsson, S., Persson, S., 1981. Tree pollen spectra in the Stockholm region
757 (Sweden), 1973-1980. *Grana* 20, 179–182.
758 <https://doi.org/10.1080/00173138109427661>

- 759 Orlandi, F., Oteros, J., Aguilera, F., Ben Dhiab, A., Msallem, M., Fornaciari, M.,
760 2014. Design of a downscaling method to estimate continuous data from
761 discrete pollen monitoring in Tunisia. *Environ. Sci. Process. Impacts* 16,
762 1716–1725. <https://doi.org/10.1039/c4em00153b>
- 763 Oteros, Jose, Galán, C., Alcázar, P., Domínguez-Vilches, E., 2013. Quality
764 control in bio-monitoring networks, Spanish Aerobiology Network. *Sci. Total*
765 *Environ.* 443, 559–565. <https://doi.org/10.1016/J.SCITOTENV.2012.11.040>
- 766 Oteros, J., García-Mozo, H., Vázquez, L., Mestre, A., Domínguez-Vilches, E.,
767 Galán, C., 2013. Modelling olive phenological response to weather and
768 topography. *Agric. Ecosyst. Environ.* 179, 62–68.
769 <https://doi.org/10.1016/j.agee.2013.07.008>
- 770 Oteros, J., Pusch, G., Weichenmeier, I., Heimann, U., Möller, R., Röseler, S.,
771 Traidl-Hoffmann, C., Schmidt-Weber, C., Buters, J.T.M., 2015. Automatic
772 and Online Pollen Monitoring. *Int. Arch. Allergy Immunol.* 167, 158–166.
773 <https://doi.org/10.1159/000436968>
- 774 Oteros, J., Sofiev, M., Smith, M., Clot, B., Damialis, A., Prank, M., Werchan, M.,
775 Wachter, R., Weber, A., Kutzora, S., Heinze, S., Herr, C.E.W., Menzel, A.,
776 Bergmann, K.-C., Traidl-Hoffmann, C., Schmidt-Weber, C.B., Buters, J.T.M.,
777 2019. Building an automatic pollen monitoring network (ePIN): Selection of
778 optimal sites by clustering pollen stations. *Sci. Total Environ.* 688, 1263–
779 1274. <https://doi.org/10.1016/J.SCITOTENV.2019.06.131>
- 780 Oteros, J., Weber, A., Kutzora, S., Rojo, J., Heinze, S., Herr, C., Gebauer, R.,
781 Schmidt-Weber, C.B., Buters, J.T.M., 2020. An operational robotic pollen
782 monitoring network based on automatic image recognition. *Environ. Res.*

- 783 191, 110031. <https://doi.org/10.1016/j.envres.2020.110031>
- 784 Pérez-Badía, R., Rapp, A., Morales, C., Sardinero, S., Galán, C., García-Mozo,
785 H., 2010. Pollen spectrum and risk of pollen allergy in central Spain. *Ann.*
786 *Agric. Environ. Med.* 17, 139–151.
- 787 Picornell, A., Buters, J., Rojo, J., Traidl-Hoffmann, C., Damialis, A., Menzel, A.,
788 Bergmann, K.C., Werchan, M., Schmidt-Weber, C., Oteros, J., 2019a.
789 Predicting the start, peak and end of the *Betula* pollen season in Bavaria,
790 Germany. *Sci. Total Environ.* 690, 1299–1309.
791 <https://doi.org/10.1016/J.SCITOTENV.2019.06.485>
- 792 Picornell, A., Oteros, J., Trigo, M.M., Gharbi, D., Docampo, S., Melgar, M., Toro,
793 F.J., García-Sánchez, J., Ruiz-Mata, R., Cabezudo, B., Recio, M., 2019b.
794 Increasing resolution of airborne pollen forecasting at a discrete sampled
795 area in the southwest Mediterranean Basin. *Chemosphere* 234, 668–681.
796 <https://doi.org/10.1016/j.chemosphere.2019.06.019>
- 797 Picornell, A., Recio, M., Ruiz-Mata, R., García-Sánchez, J., Cabezudo, B., Trigo,
798 M.M., 2020. Medium- and long-range transport events of *Alnus* pollen in
799 western Mediterranean. *Int. J. Biometeorol.* 64, 1637–1647.
800 <https://doi.org/10.1007/s00484-020-01944-7>
- 801 Picornell, A., Recio, M., Trigo, M.M., Cabezudo, B., 2019c. Preliminary study of
802 the atmospheric pollen in Sierra de las Nieves Natural Park (Southern
803 Spain). *Aerobiologia.* 35, 571–576. [https://doi.org/10.1007/s10453-019-](https://doi.org/10.1007/s10453-019-09591-1)
804 [09591-1](https://doi.org/10.1007/s10453-019-09591-1)
- 805 Ribeiro, H., Cunha, M., Abreu, I., 2007. Definition of main pollen season using a
806 logistic model. *Ann. Agric. Environ. Med.* 14, 259–264.

- 807 Ritenberga, O., Sofiev, M., Kirillova, V., Kalnina, L., Genikhovich, E., 2016.
808 Statistical modelling of non-stationary processes of atmospheric pollution
809 from natural sources: Example of birch pollen. *Agric. For. Meteorol.* 226–
810 227, 96–107. <https://doi.org/10.1016/j.agrformet.2016.05.016>
- 811 Rivas-Martínez, S., Penas, Á., del Río, S., Díaz González, T.E., Rivas-Sáenz, S.,
812 2017. Bioclimatology of the Iberian Peninsula and the Balearic Islands, in:
813 Loidi, J. (Ed.), *The Vegetation of the Iberian Peninsula*. Springer, Cham,
814 Utrecht, Netherlands, pp. 29–80. [https://doi.org/10.1007/978-3-319-54784-](https://doi.org/10.1007/978-3-319-54784-8_2)
815 [8_2](https://doi.org/10.1007/978-3-319-54784-8_2)
- 816 Rojo, J., Orlandi, F., Pérez-Badia, R., Aguilera, F., Ben Dhiab, A., Bouziane, H.,
817 Díaz de la Guardia, C., Galán, C., Gutiérrez-Bustillo, A.M., Moreno-Grau, S.,
818 Msallem, M., Trigo, M.M.M., Fornaciari, M., 2016. Modeling olive pollen
819 intensity in the Mediterranean region through analysis of emission sources.
820 *Sci. Total Environ.* 551–552, 73–82.
821 <https://doi.org/10.1016/j.scitotenv.2016.01.193>
- 822 Rojo, J., Oteros, J., Pérez-Badia, R., Cervigón, P., Ferencova, Z., Gutiérrez-
823 Bustillo, A.M., Bergmann, K.C., Oliver, G., Thibaudon, M., Albertini, R.,
824 Rodríguez-De la Cruz, D., Sánchez-Reyes, E., Sánchez-Sánchez, J., Pessi,
825 A.M., Reiniharju, J., Saarto, A., Calderón, M.C., Guerrero, C., Berra, D.,
826 Bonini, M., Chiodini, E., Fernández-González, D., García-Sánchez, J., Trigo,
827 M.M., Myszkowska, D., Fernández-Rodríguez, S., Tormo-Molina, R.,
828 Damialis, A., Kolek, F., Traidl-Hoffmann, C., Severova, E., Caeiro, E.,
829 Ribeiro, H., Magyar, D., Makra, L., Udvardy, O., Alcázar, P., Galán, C.,
830 Borycka, K., Kasprzyk, I., Newbiggin, E., Adams-Groom, B., Apangu, G.P.,
831 Frisk, C.A., Skjøth, C.A., Radišić, P., Šikoparija, B., Celenk, S., Schmidt-

- 832 Weber, C.B., Buters, J., 2019. Near-ground effect of height on pollen
833 exposure. *Environ. Res.* 160–169.
834 <https://doi.org/10.1016/j.envres.2019.04.027>
- 835 Rojo, Jesús, Picornell, A., Oteros, J., 2019. AeRobiology: the computational tool
836 for biological data in the air. *Methods Ecol. Evol.* 10, 1371–1376.
837 <https://doi.org/10.1111/2041-210x.13203>
- 838 Rojo, J., Picornell, A., Oteros, J., Werchan, M., Werchan, B., Bergmann, K.C.,
839 Smith, M., Weichenmeier, I., Schmidt-Weber, C.B., Buters, J., 2021.
840 Consequences of climate change on airborne pollen in Bavaria, Central
841 Europe. *Reg. Environ. Chang.* 21, 9. [https://doi.org/10.1007/s10113-020-](https://doi.org/10.1007/s10113-020-01729-z)
842 [01729-z](https://doi.org/10.1007/s10113-020-01729-z)
- 843 Rubin, D.B., 1976. Inference and missing data. *Biometrika* 63, 581–592.
844 <https://doi.org/10.1093/biomet/63.3.581>
- 845 Ruiz-Valenzuela, L., Aguilera, F., 2018. Trends in airborne pollen and pollen-
846 season-related features of anemophilous species in Jaen (south Spain): A
847 23-year perspective. *Atmos. Environ.* 180, 234–243.
848 <https://doi.org/10.1016/j.atmosenv.2018.03.012>
- 849 Schouten, R.M., Lugtig, P., Vink, G., 2018. Generating missing values for
850 simulation purposes: a multivariate amputation procedure. *J. Stat. Comput.*
851 *Simul.* 88, 2909–2930. <https://doi.org/10.1080/00949655.2018.1491577>
- 852 Skjøth, C.A., Damialis, A., Belmonte, J., De Linares, C., Fernández-Rodríguez,
853 S., Grinn-Gofroń, A., Jędrzycka, M., Kasprzyk, I., Magyar, D., Myszkowska,
854 D., Oliver, G., Páldy, A., Pashley, C.H., Rasmussen, K., Satchwell, J.,
855 Thibaudon, M., Tormo-Molina, R., Vokou, D., Ziemianin, M., Werner, M.,

- 856 2016. *Alternaria* spores in the air across Europe: abundance, seasonality
857 and relationships with climate, meteorology and local environment.
858 *Aerobiologia*. 32, 3–22. <https://doi.org/10.1007/s10453-016-9426-6>
- 859 Stekhoven, D.J., Buhlmann, P., 2012. MissForest--non-parametric missing value
860 imputation for mixed-type data. *Bioinformatics* 28, 112–118.
861 <https://doi.org/10.1093/bioinformatics/btr597>
- 862 Su, Y.-S., Gelman, A., Hill, J., Yajima, M., 2011. Multiple Imputation with
863 Diagnostics (mi) in R: Opening Windows into the Black Box. *J. Stat. Softw.*
864 45, 1–31. <https://doi.org/10.7916/D8VQ3CD3>
- 865 Todorov, V., 2020. rrcovNA: Scalable Robust Estimators with High Breakdown
866 Point for Incomplete Data.
- 867 Valencia, J.A., Astray, G., Fernández-González, M., Aira, M.J., Rodríguez-Rajo,
868 F.J., 2019. Assessment of neural networks and time series analysis to
869 forecast airborne *Parietaria* pollen presence in the Atlantic coastal regions.
870 *Int. J. Biometeorol.* 63, 735–745. <https://doi.org/10.1007/s00484-019-01688->
871 [z](https://doi.org/10.1007/s00484-019-01688-z)
- 872 van Buuren, S., Groothuis-Oudshoorn, K., 2011. MICE: Multivariate imputation
873 by chained equations in R. *J. Stat. Softw.* 45, 1–68.
874 <https://doi.org/10.18637/jss.v045.i03>
- 875 Velasco-Jiménez, M.J., Alcázar, P., Domínguez-Vilches, E., Galán, C., 2013.
876 Comparative study of airborne pollen counts located in different areas of the
877 city of Córdoba (south-western Spain). *Aerobiologia*. 29, 113–120.
878 <https://doi.org/10.1007/s10453-012-9267-x>

This is the accepted version of the manuscript. The published version is available at:
<https://doi.org/10.1016/j.envres.2021.111391>

879 Verstraeten, W.W., Kouznetsov, R., Hoebeke, L., Bruffaerts, N., Sofiev, M.,
880 Delcloo, A.W., 2021. Modelling grass pollen levels in Belgium. *Sci. Total*
881 *Environ.* 753, 141903. <https://doi.org/10.1016/j.scitotenv.2020.141903>

882