

Document downloaded from:

<http://hdl.handle.net/10251/88143>

This paper must be cited as:

Bas Cerdá, MDC.; Ortiz Moragón, J.; Ballesteros Pascual, L.; Martorell Alsina, SS. (2017). Evaluation of a multiple linear regression model and SARIMA model in forecasting 7Be air concentrations. *Chemosphere*. 177:326-333. doi:10.1016/j.chemosphere.2017.03.029.



The final publication is available at

<http://doi.org/10.1016/j.chemosphere.2017.03.029>

Copyright Elsevier

Additional Information

1 EVALUATION OF A MULTIPLE LINEAR REGRESSION MODEL AND SARIMA MODEL IN 2 FORECASTING ⁷BE AIR CONCENTRATIONS

3 1. Introduction

4 ⁷Be is widely used as an atmospheric radiotracer due to its relatively short life ($T_{1/2} = 53.3$ days) and ease
5 of measurement by γ -spectrometry, which provides important information on atmospheric air mass
6 motions. A better understanding of its distribution would facilitate refinement and validation of global
7 atmospheric circulation models (Dueñas et al. 2015). ⁷Be forecasting can thus be adopted as a target value
8 in analyzing fluctuations or deviations that could imply important atmospheric changes.

9 It is generally accepted that the ⁷Be production rate depends on a number of atmospheric factors. Several
10 studies have pointed out that the intensity of galactic cosmic rays in the Earth's orbit is affected by solar
11 activity and the geomagnetic field, which is under constant cosmic ray bombardment from space (O'Brien,
12 1979; Vogt et al., 1990; Hötzl et al., 1991; Ioannidou&Papastefanou, 1994). In particular, an increase in
13 solar activity and the geomagnetic field reduce the galactic cosmic ray flux, which is followed by reduced
14 ⁷Be production.

15 In addition to the above-mentioned sources of variability, ⁷Be concentrations in the lower layers of the
16 atmosphere present temporal variations caused by solar radiation and meteorological parameters that can
17 affect regional weather patterns (temperature, relative humidity, precipitations, wind speed and wind
18 direction) (Feely et al., 1989; Baeza et al., 1996).

19 Many research studies have analyzed the relation between ⁷Be air concentrations and the meteorological
20 and atmospheric variables using a simple correlation analysis (e.g. Dueñas et al., 1999; Ioannidou et al.,
21 2006; Piñero-García & Ferro-García, 2013; Ceballos et al., 2016; Neroda et al.; 2016). Furthermore, some
22 of these studies have applied Multiple Linear Regression (MLR) analysis to develop an explanatory and
23 predictive model for ⁷Be air concentrations using the atmospheric and meteorological variables as
24 predictors (Table 1).

25

26

Location	Period	Significant variables used in MLR	R ²	Source
Málaga, Spain	1992-1995	- Maximum Temperature - Rainfall - Relative Humidity - Hours of sunshine	27%	Dueñas et al. (1999)
Thessaloniki, Greece	1987-2001	- Temperature - Relative Humidity - Sunspot Number	38.5%	Ioannidou et al. (2006)
Granada, Spain	1993-2001	- Temperature - Rainfall - Sunspot Number	71%	Azahra et al. (2004)
Málaga, Spain	1997-2007	- Solar energetic proton - Aerosol optical depth	34%	Dueñas et al. (2015)
Granada, Spain	1996-2010	- Temperature - Relative Humidity - Sunspot Number	52%	Piñero-García & Ferro-García (2013)
Granada, Spain	2005-2009	- Temperature - Relative Humidity - Rainfall	72.16%	Piñero-García et al. (2012)
Plymouth, UK	2009-2010	- Rainfall	94%	Taylor et al. (2016)
Granada, Spain	2011-2014	- Solar Irradiance - Total suspended particles	66.9%	Essaid et al. (2015)
Vladivostok, Russia	2013-2014	- Altitude - Precipitation - Temperature - Aerosol concentration - Trajectories in the pacific (North-East)	55%	Neroda et al. (2016)

27 Table 1. ⁷Be predictive models for different time periods at different locations.

28 Each study uses several predictors to explain ⁷Be air concentration in different time periods at different

29 locations. The explicative power of the model, measured by the R square coefficient, is, in general, less

30 than 50%. The studies that get the highest R², use a historical data range of less than five years, which may

31 not be enough information to forecast the ⁷Be air concentration for the following year. In addition to

32 explanatory power, it is very important to compute accuracy measurements with data that have not been

33 used to develop the model. This procedure is not applied in the above MLR models and is important in

34 measuring the validity and forecasting power of the model, which is one of the aims of the present study.

35 Several authors recommend the use of time series modelling techniques instead of multiple linear regression

36 when monitoring correlated process data (Alwan & Roberts 1988; Harris & Ross 1991; Wardell et al. 1994).

37 Classical regression is often insufficient for explaining all the interesting dynamics of a time series. For

38 instance, the estimated autocorrelation function (ACF) of the residuals of the regression model could reveal

39 additional structure in the data that the regression did not capture. Instead, the introduction of Box-Jenkins

40 models could deal with the limitations of classical regression in time series (Shumway & Stoffer, 2006).

41 A recent study applied a decomposition of the ^{7}Be time series into a trend-cycle, a seasonal and an irregular
42 component in order to separate the inter- and intra-annual patterns of ^{7}Be variability (Bas et al, 2016). The
43 results of this study showed the suitability of applying time series analysis to correlated data in order to
44 separate the different sources of variability of ^{7}Be concentrations and to develop a forecasting model.

45 The aim of this study is to propose two models to explain and forecast ^{7}Be air concentrations: i) a Seasonal
46 Autoregressive Integrated Moving Average (SARIMA) model and ii) a Multiple Linear Regression (MLR)
47 model using meteorological and atmospheric variables. Both the time series and multiple linear regression
48 models are evaluated by comparison with real ^{7}Be air concentrations for the city of Valencia in 2007-2014
49 and with out-of-sample tests for the 12 months of the year 2015, using the Root Mean Square Error (RMSE)
50 and the Adapted Mean Absolute Percentage Error (AMAPE) as forecasting accuracy measures. Finally, the
51 results of the accuracy measurements of both models are compared.

52

53 **2. Material and methods**

54 *2.1. Study area and sampling*

55 Airborne particulate samples were collected weekly on the campus of the Universitat Politècnica de
56 Valencia from January 2007 to December 2015. Valencia is situated on the east coast of Spain (15m above
57 sea level) in the western Mediterranean Basin (39°28'50" N, 0°21'59" W) and has a relatively dry
58 subtropical Mediterranean climate with very mild winters and long hot summers. The sampling point was
59 located approximately 2 km away from the coastline.

60 Aerosol samples were collected using Eberlyne G21DX and Saic AVS28A air samplers placed
61 approximately 1 m above ground level. The aerosol particles were retained on a cellulose filter of 4.2×10^{-2}
62 m effective diameter and 0.8 μm pore size. The filters were changed weekly and the average volume
63 ranged from 300 to 400 m^3 per week. Each filter was put inside a plastic box and kept in a desiccator until
64 it was measured.

65

66 2.2. ⁷Be activity measurements

67 A monthly composite sample containing 4-5 filters was measured by γ -spectrometry to determine specific
68 ⁷Be activities using an HPGe detector (ORTEC Industries, USA) n-type with relative efficiency of 18% for
69 ⁶⁰Co gamma-ray. A certificated standard containing radionuclides with energies ranging from 59 to 1836.1
70 keV was used for preparing the calibrated filters, which were placed inside their plastic boxes on the top of
71 the detector. The counting time was 60000s and the γ -line 477.7 KeV was used to calculate the activity.
72 ORTEC Gamma-Vision software was used for acquisition and analysis. Concentration activities were
73 corrected for the radioactive decay to the mid-collection period. The mean measured uncertainties (K=2)
74 were around 10 %.

75

76 2.3. Statistical analysis

77 **SARIMA MODEL**

78 The SARIMA model building process is designed to take advantage of the association in the sequentially
79 lagged relationships that usually exists in data collected periodically. A time series $\{z_t, t = 1, \dots, N\}$ is
80 generated by a SARIMA(p, d, q)(P, D, Q)_s model if:

81
$$\phi_p(B)\Phi_P(B^s)(1-B)^d(1-B^s)^D z_t = \theta_q(B)\Theta_Q(B^s)a_t$$

82 where N is the number of observations; p, d, q, P, D, Q are integers; B is the lag operator (e.g. $w_t = z_t -$
83 $z_{t-s} = (1 - B^s)z_t$); s is the seasonal period length; d is the number of regular differences ($d \leq 2$); D is
84 the number of seasonal differences, and a_t is the estimated residual at time t , which is a usual Gaussian
85 white noise process (WN).

86 $\phi_p(B) = 1 - \phi_1 B - \phi_2 B^2 - \dots - \phi_p B^p$, is the regular autoregressive operator (AR) of order p ,

87 $\theta_q(B) = 1 - \theta_1 B - \theta_2 B^2 - \dots - \theta_q B^q$, is the regular moving average operator (MA) of order q ,

88 $\Phi_P(B^s) = 1 - \phi_1 B^s - \phi_2 B^{s^2} - \dots - \phi_P B^{sP}$, is the seasonal autoregressive operator (SAR) of order P ,

89 $\Theta_Q(B^s) = 1 - \theta_1 B^s - \theta_2 B^{s^2} - \dots - \theta_Q B^{sQ}$, is the seasonal moving average operator (SMA) of order Q .

90

91 As reported by Box & Jenkins (1976) and Shumway &Stoffer (2006), the SARIMA model consists of three
92 main steps:

93 ***Identification and estimation step***

94 First, the periodogram technique was applied to identify the periodic cycle in the time series (Schuster,
95 1898). The periodogram plot should have clear peaks at points corresponding to the ‘hidden periods’ of the
96 cyclic model.

97 The time series should then be differenced in order to be stationary in mean and variance (identifying d and
98 D parameters). The differencing technique can also be used to remove trends, which are usually detected
99 by inspecting the plot of the ${}^7\text{Be}$ data over the period considered. However, they are also characterized by
100 the autocorrelation function.

101 After differencing the time series, a tentative autoregressive moving average (ARMA) process is carried
102 out based on the estimated autocorrelation function (ACF) and the estimated partial autocorrelation function
103 (PACF). The shape of the ACF and PACF of the real time series is compared with the shape of the
104 theoretical model to identify possible different p , q , P and Q parameters of the SARIMA model (Peña,
105 2010; Shumway & Stoffer, 2006). Having specified tentative models in the identification step, the
106 parameters of the candidate models are estimated by a maximum likelihood function (Shine & Lee, 2000).

107 After trying several combinations for parameters p , q , P and Q , the best model was selected, considering
108 the minimum MAPE, AMAPE and RMSE (defined in the section on the Forecasting Step) for the
109 forecasting data of the sample and out-of-sample as accuracy measures of predictive power.

110 The selection of the most parsimonious model is also based on Akaike’s Information Criterion (AIC), which
111 rewards models for good fit and penalize them for complexity. The model with the minimum AIC is chosen
112 as the parsimonious model. The AIC coefficient is defined as follows:

113
$$AIC = 2 \ln(RMSE) + \frac{2(p + q)}{n}$$

114 where p and q are the number of parameters of AR and MA estimates, RMSE is the Root Mean Square
115 Error (defined in the section on the Forecasting Step) and n is the sample size of the data used to fit the
116 model.

117 ***Validation step***

118 In this step, several statistics were used to check the suitability of the identified models. An essential part
 119 of the procedure is to examine the residuals of the SARIMA model, which, if the model is satisfactory,
 120 should be considered as White Noise (WN). We examine some simple tools for checking the hypothesis
 121 that the residuals are WN and the model is valid. If the fit model passes the following tests, it can be used
 122 to make a forecast.

123 - *t-ratio test* to evaluate the significance of the parameters estimated in each model. The parameters are
 124 considered significant with a 95% of confidence level if p-values<0.05.

125 - *Kolmogorov-Smirnov test* applying Lilliefors correction of the residual series to check that the noise
 126 process is Gaussian. The residual series is Gaussian if p-values>0.05.

127 - *Q* Ljung-Box statistic* to check the condition that the residuals can be considered as a WN. The statistic
 128 proposed is:

$$129 \quad Q^* = n(n + 2) \sum_{k=1}^m (n - k)^{-1} r_k(a)$$

130 where $r_k(a)$ is the sample autocorrelation f order k of the residual, n is the length of residual series and m
 131 is the number of lags considered, $Q^* \approx \chi_{m-n}^2$, $n = p + q + P + Q$. The model is considered valid if
 132 $P(\chi^2(m - n) > Q^*) > 0.05$. In this study, the Q^* Ljung-Box statistic is calculated for a large m in each
 133 model, as suggested by Peña (2010).

134

135 ***Forecasting step***

136 To assess the forecasting performance of different models the data set is divided into two samples for
 137 training and testing. This procedure is known as an out-of-sample technique, which means that the training
 138 data used in model fitting are different to the test sample (out-of-sample) used to evaluate the established
 139 model.

140 Several measurement statistics can be used to examine the forecast accuracy of different models. Root
 141 Mean Square Error (RMSE) and Mean Absolute Percentage Error (MAPE) are the criteria most frequently
 142 used to evaluate the performance of the forecasting models. One of the disadvantages of the MAPE criteria
 143 is the adverse effect of small actual values, in which case MAPE criteria will contribute large terms to the

144 MAPE coefficient, even if the difference between the actual and forecast values is small. It is therefore
 145 better to use an adapted MAPE (AMAPE), as defined in various studies (Tsay, 2005; Wu & Shahidehpour,
 146 2010):

$$147 \quad RMSE = \sqrt{\frac{\sum_{t=1}^n (\hat{z}_t - z_t)^2}{n}}$$

$$148 \quad MAPE = \left(\frac{1}{N} \sum_{t=1}^n \left(\frac{|\hat{z}_t - z_t|}{z_t} \right) \right) 100\%$$

$$149 \quad AMAPE = \frac{1}{n} \sum_{t=1}^n \left(\frac{|\hat{z}_t - z_t|}{\frac{1}{n} \sum_{t=1}^n z_t} \right) * 100\%$$

150 where t represents the time and n is the sample size for forecasts; \hat{z}_t is the forecast at t from any mentioned
 151 model and z_t is the actual value at t . The RMSE statistic depends on the scale of the variables and measures
 152 the absolute errors. The MAPE and AMAPE statistics measure the relative errors. The smaller the RMSE,
 153 MAPE and AMAPE the better the accuracy of the model.

154

155 **MULTIPLE LINEAR REGRESSION**

156 Multiple linear regression analysis is a multivariate statistical technique used to examine the relationship
 157 between a single dependent variable and a set of independent variables. The main objectives of MLR are
 158 explanation and prediction. Explanation examines the regression coefficients, their magnitude, sign and
 159 statistical inference, for each independent variable. Prediction involves the extent to which the independent
 160 variables can predict the dependent variable (Hair et al., 2010). MLR forecasting models are expressed in
 161 the following format:

$$162 \quad Y_t = X_t \beta + \varepsilon_t$$

163 where Y_t is the predicted value at time t , $X_t = (1, x_{1t}, x_{2t}, \dots, x_{kt})$ is a vector of k explanatory variables at
 164 time t , $\beta = (\beta_0, \beta_1, \dots, \beta_k)^T$ is the vector of coefficients, and ε_t is a random error term at time t , $t = 1, \dots, N$.

165 The errors terms should be independent and have a Gaussian distribution.

166

167 The assumptions of the MLR model (independent errors and Gaussian error term distribution) could be
168 analyzed by obtaining the Kolmogorov-Smirnov test and the Q* Ljung-Box statistic, as in the time series.

169 The explanatory power of the MLR is commonly measured by the R square coefficient defined as follows:

170
$$R^2 = \left(\frac{\sigma_{\hat{Y}_t, Y_t}^2}{\sigma_{\hat{Y}_t}^2 \sigma_{Y_t}^2} \right) 100\%$$

171 where $\sigma_{\hat{Y}_t, Y_t}^2$ is the covariance of the forecast and actual values; $\sigma_{\hat{Y}_t}^2$ and $\sigma_{Y_t}^2$ the variance of the forecast and
172 actual values respective.

173 The forecasting power could be measured using the same accuracy measurements as in the time series.

174

175 **3. Results in forecasting ⁷Be air concentrations**

176 The first step in developing any forecasting model is to plot the data. In view of the results obtained in a
177 recent study (Bas et al., 2017), the best ⁷Be concentration forecasting results are based on a time window
178 of at least eight years of data. This result supports the training sample of eight years of historical data (2007-
179 2014) and the out-of-sample test for one year (2015) selected in this study. Figure 1 shows the evolution of
180 ⁷Be air concentrations during the entire period 2007-2015. ⁷Be activity concentrations ranged from 2.28 to
181 8.11 mBq/m³ with an arithmetic mean of 4.62 ± 1.19 mBq/m³ during the period studied.

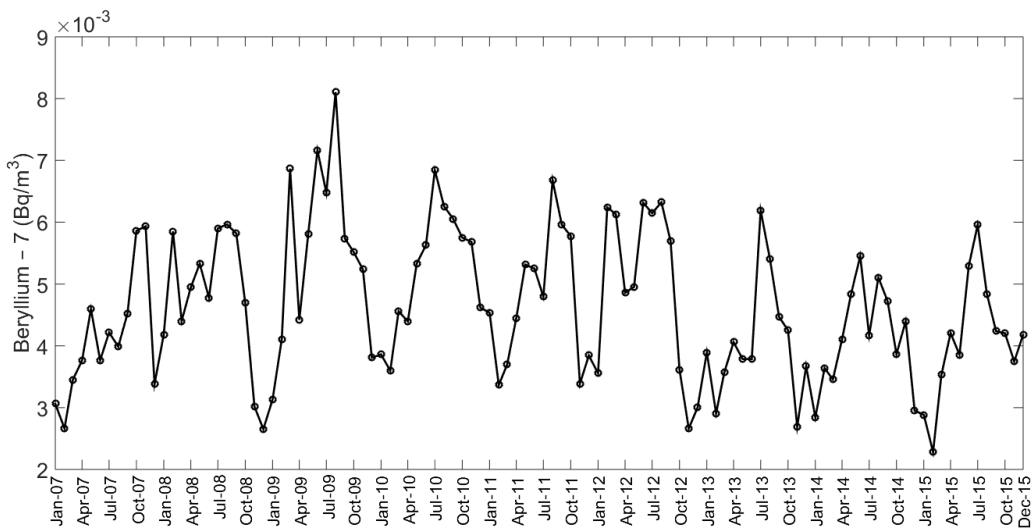
182

183

184

185

186
187
188
189
190
191
192



193

Fig 1. Temporal evolution of ⁷Be air concentration over the period 2007-2015.

194 *3.1. SARIMA model*

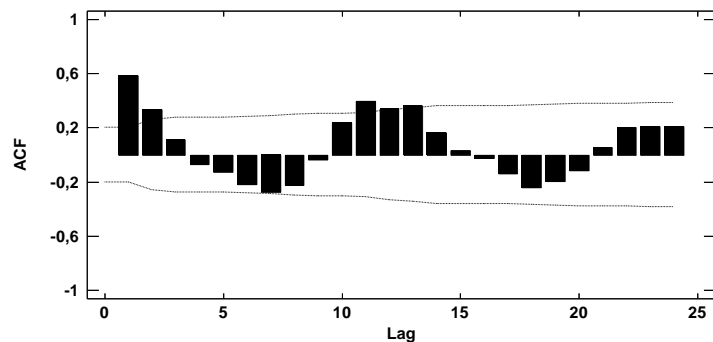
195 The evolution of ⁷Be air concentrations suggests that there exists a seasonal pattern with a sinusoidal trend.

196 The result of the periodogram technique identified a relevant peak corresponding to a period of 12 months
197 (annual periodicity, $s = 12$).

198 For the identification step, a simple ACF (Figure2) that is positive and very slowly decaying in lag 1 and
199 in the seasonal lag 12 suggests a regular and seasonal difference ($d = D = 1$). The

200 SARIMA($p, 1, q$)($P, 1, Q$)₁₂ model is therefore useful for representing ⁷Be air concentrations with a trend.

201 The differenced ⁷Be time series is stationary.



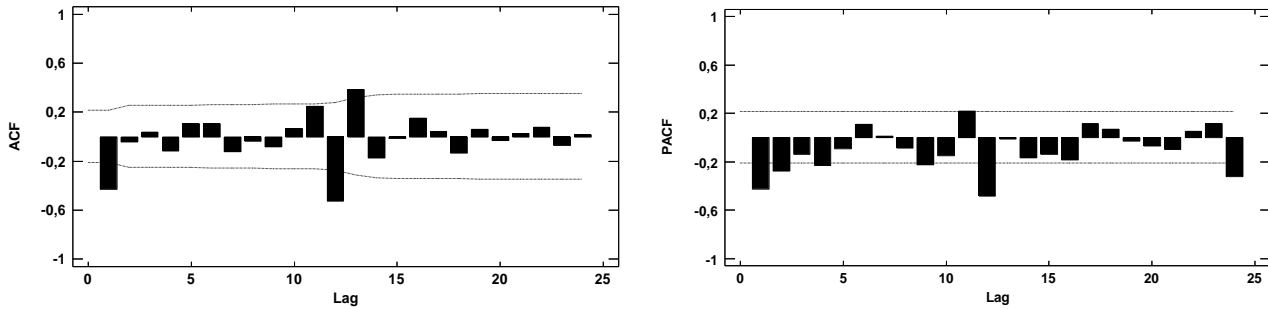
202

203

Fig.2. The sample ACF of the ⁷Be time series

204 After differencing the time series, a tentative autoregressive moving average process is carried out based
 205 on the estimated autocorrelation function (ACF) and the estimated partial autocorrelation function (PACF).
 206 Figure 3 shows that the autocorrelations at lag 1, 12 and 24 are significant in the PACF of the residuals.

207



208 Fig. 3. The sample ACF and PACF of the residuals after applying a regular and seasonal difference.

209 Table 2 reports the results of tentative SARIMA models considering the ACF and PACF of the residuals
 210 (Figure 3) after applying a regular and a seasonal difference ($d = D = 1$). The following accuracy
 211 measurements were used to select the best and most parsimonious model: RMSE, MAPE, AMAPE and
 212 AIC.

213

	Training sample (2007-2014)				Out-of-sample (2015)		
	RMSE	MAPE	AMAPE	AIC	RMSE	MAPE	AMAPE
SARIMA(0, 1, 1)(2, 1, 2)₁₂	0.000722	13.00%	11.89%	-14.36	0.00147	29.19%	27.18%
SARIMA(0, 1, 2)(2, 1, 2)₁₂	0.000723	13.08%	11.98%	-14.37	0.00147	29.59%	27.13%
SARIMA(0, 1, 1)(1, 1, 3)₁₂	0.000678	11.93%	10.74%	-14.49	0.00078	17.75%	17.20%
SARIMA(0, 1, 2)(1, 1, 3)₁₂	0.000672	11.76%	10.70%	-14.49	0.00078	18.05%	17.40%

214 Table 2. Models selection criterion

215 Table 2 show that the AIC criterion is similar in the different models proposed. However, considering that
 216 the RMSE, MAPE and AMAPE coefficients should be minimum, the SARIMA(0,1,1)(1,1,3) and
 217 SARIMA(0,1,2)(1,1,3) models are the best options, considering the analysis in both samples (training and
 218 out-of-samples). Of the two, we propose the SARIMA(0,1,1)(1,1,3) model as it is simpler than the other
 219 and the RMSE, MAPE and AMAPE coefficients in the training sample and in out-of-sample are similar in
 220 both models.

221 Having specified the best model in the identification step, the parameters are estimated by a maximum
 222 likelihood function and the estimated model can be written as follows:

$$223 \quad (1 + 0.814B^{12})(1 - B)(1 - B^{12})z_t = (1 - 0.665B)(1 - 0.555B^{12} - 0.932B^{24} + 0.687B^{36})a_t$$

224 where $a_t \approx WN(0, \sigma = 8.07E - 04)$. WN=White Noise.

225 The parameters estimated in the model are significant (p-values<0.05) (Table 3). The residuals obtained
 226 from fitting a SARIMA(0,1,1)x(1,1,3)₁₂ model to ⁷Be concentration data for a time window of eight years
 227 (2007-2014) are normally distributed (K-S test, $p - value > 0.05$) with mean zero and standard deviation
 228 $\sigma = 8.07E - 04$. Moreover, significant autocorrelation is not found in the residuals (Q^* test, $p - value >$
 229 0.05), therefore the residuals can be considered as WN and the SARIMA(0,1,1)x(1,1,3)₁₂ can be
 230 considered a suitable forecasting model.

	t-ratio test (p-value)	K-S Lilliefors (p-value)	Q^* Ljung-Box (p-value)
θ_1	8.0431 (<0.000001)	D = 0.079983	$\chi^2 = 47.809$
Φ_1	-9.4919 (<0.000001)	(p-v: 0.2119)	(p-v: 0.2837)
θ_1	8.08721 (<0.000001)		m=48
θ_2	19.5035 (<0.000001)		
θ_3	-13.7005 (<0.000001)		

231 Table 3. Validation of the proposed SARIMA model.

232 The predictive model obtained, after developing the above expression, is:

$$233 \quad \hat{z}_t = z_{t-1} + 0.186z_{t-12} - 0.186z_{t-13} + 0.814z_{t-14} - 0.814z_{t-25} - 0.665a_{t-1} - 0.555a_{t-12}$$

$$234 \quad \quad \quad + 0.369a_{t-13} - 0.932a_{t-24} + 0.619a_{t-25} + 0.687a_{t-36} - 0.456a_{t-37} + a_t$$

235 Figure 4 shows the comparison between measured and forecast values using a
 236 SARIMA(0,1,1)x(1,1,3)₁₂ in the 2007-2014 training sample and in the out-of-sample data in 2015. The
 237 time series proposed explains 70.88% of the variability of the actual data.

238
 239
 240
 241
 242
 243
 244

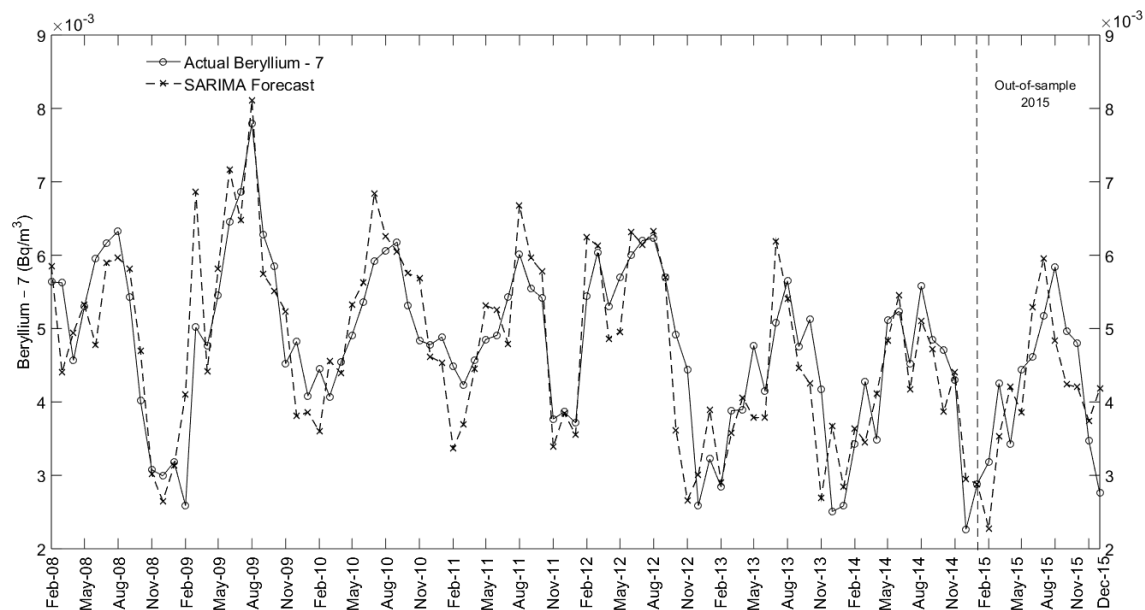


Fig. 4. Comparison between measured and forecast (SARIMA) power

245

246

247

248 *3.2. MULTIPLE LINEAR REGRESSION analysis model*

249

250

251

252

253

254

255

A multiple linear regression model is developed to explain and forecast ^7Be air concentrations. The atmospheric parameters studied in the present work are: sunspot number (SSN), temperature (T) (in tenths of $^{\circ}\text{C}$), precipitation (PP) (in tenths of a millimetre), relative humidity (RH) (in %) and wind speed (WS) (in km/h). The meteorological factors were collected by the Universitat Politècnica de Valencia's weather station, which was also the sampling point for ^7Be activity. The sunspot number parameter (SSN) was collected daily during the period 2007-2015 by the World Data Center SILSO, Royal Observatory of Belgium, in Brussels (SILSO, 2015).

256

257

258

259

We selected these variables after taking into account the atmospheric parameters that mainly affect Valencia weather, with a relatively dry subtropical Mediterranean climate, very mild winters and long hot summers, and considering the variables adopted in a previous study (Bas et al, 2016) and the variables most frequently considered to study ^7Be activity in the literature.

260

261

262

A logarithmic transformation of the ^7Be variable is applied to better identify a Gaussian distribution in the data. In this study we considered the mean monthly values of temperature, relative humidity, wind speed, and sunspot number. The precipitation factor was considered as the number of rainy days per month due to

263 the particular rainfall regime in Valencia, with few days of torrential rainfall and many dry days. Solar
 264 activity was considered as measured by the sunspot number parameter.

265 The R^2 obtained for the regression given below is significant at the 95% confidence level, however this
 266 model explains only 48.76% of the ^7Be variability. The predictive model obtained is:

267
$$LN(^7\text{Be}) = -5.3631 + 0.0025 * T - 0.0602 * WS - 0.0112 * PP - 0.0018 * SSN$$

<i>Parameter</i>	<i>Estimation</i>	<i>St. Error</i>	<i>t-statistic</i>	<i>p-value</i>
β_0	-5.3631	0.1449	-37.0094	<0.00001
T	0.0025	0.0004	6.3614	<0.00001
WS	-0.0602	0.0176	-3.4199	0.0009
PP	-0.0112	0.0049	-2.2957	0.0240
SSN	-0.0018	0.0004	-4.1275	0.0001

268 Table 4. Estimated parameters and its significance in the MLR model.

269 The significant variables that affect ^7Be air concentration are: temperature, wind speed, precipitation and
 270 sunspot number (Table 4). However, the relative humidity variable is positively correlated with temperature
 271 ($r = 0.67, p - value < 0.00001$), so that both variables explain the same behaviour of ^7Be activity and
 272 the multiple regression technique selected the variable most highly correlated with ^7Be . Note that all the
 273 variables have an inverse influence on ^7Be activity, except temperature, which has a positive effect.

274 The Kolmogorov-Smirnov test was applied to check the normality of the residuals, obtaining $D = 0.056379$
 275 with a p-value of 0.604. The residuals can therefore be considered Gaussian. Finally, the Ljung-Box test
 276 was also applied to check the randomness of the residuals. In this case, the p-value obtained for any lag (m)
 277 considered is less than 0.05, which means that the residuals are not random and this result reveals additional
 278 structure in the data that the regression could not capture.

279 Figure 5 shows the comparison between measured and forecast values using an MLR in the training sample
 280 2007-2014 and in the forecasting data in 2015.

281

282

283

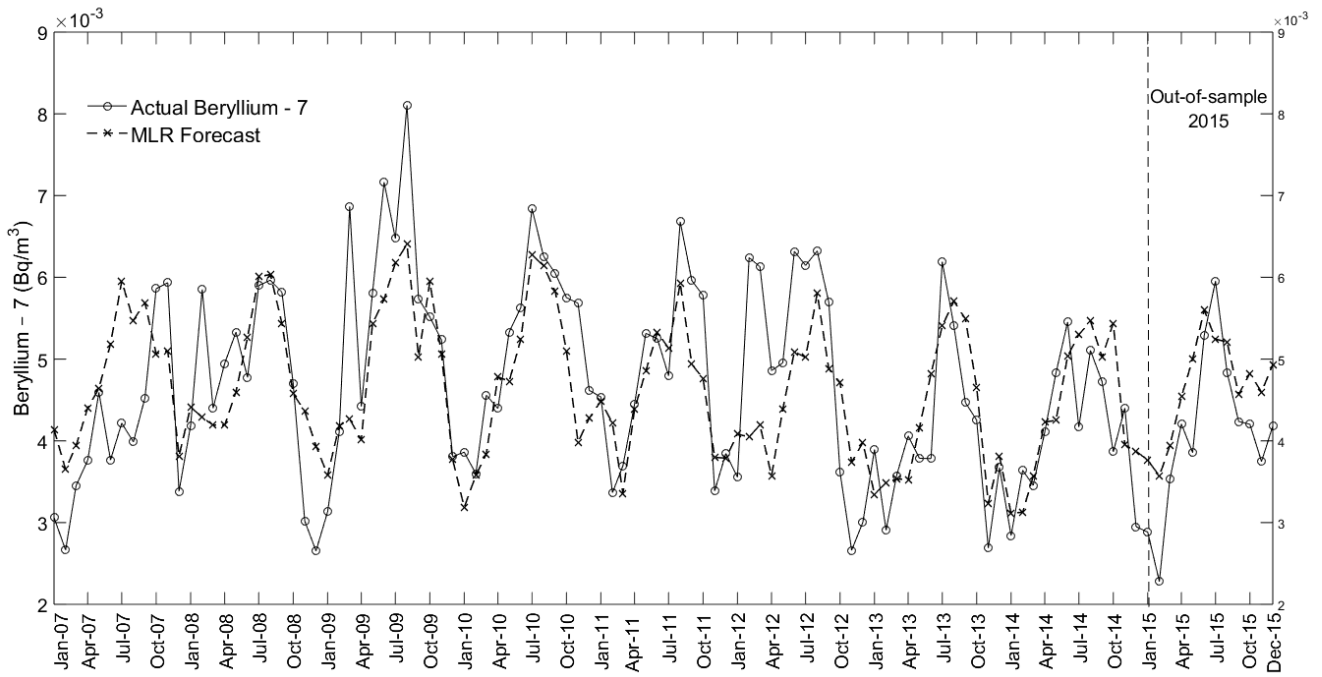


Fig. 5. Comparison between measured and forecast (MLR) power

284

285

286 *3.3. Comparison of the forecasting performance of the SARIMA and MLR models*

287 Table 5 shows the explanatory and forecasting power of the SARIMA and MLR models. For the former
 288 power we used the R^2 coefficient. The accuracy measures used to analyze the validity of the model are the
 289 RMSE and AMAPE coefficients, considering the following sample sizes for the out-of-sample forecasts:
 290 $n = 1, n = 3, n = 6, n = 9,$ and $n = 12$ months. As can be seen in Table 5, the RMSE value for $n = 1$ is
 291 very different to that of $n > 1$, suggesting that predictions for 1-month periods are uncertain. The selection
 292 model criteria are therefore based on forecasts for at least three months.

Model	Explanatory power	Forecasting power Out-of-sample Year=2015				
		RMSE and AMAPE				
		R^2	$n = 1$	$n = 3$	$n = 6$	$n = 9$
MLR	48.76%	0.00088	0,00093 29,68%	0,00083 19,81%	0,00073 15,59%	0,00074 16,27%
SARIMA(0, 1, 1)(1, 1, 3) ₁₂	70.88%	1E-07	0.00067 18.70%	0.00068 16.61%	0.00073 16.66%	0.00078 17.20%

293 Table 5. Comparison of explicative and forecasting power between SARIMA and MLR

294 In the MLR model the atmospheric variables explain 48.76% of the variability of the ^7Be air concentration,
 295 whereas the SARIMA model explains 70.88%. The predictive model cannot explain more variability in ^7Be
 296 activity due to the joint effect of the parameters considered, which masks the intra and inter annual

297 components of the time series. This result agrees with observations made in previous studies (Piñero-García
298 & Ferro-García, 2013, Dueñas et al., 2015, Bas et al., 2016).

299 Considering the forecasting power in the out-of-sample year, the RMSE and AMAPE accuracy measures
300 are very similar for $n = 9$ and $n = 12$ in both models, although slightly lower in the MLR model. However,
301 these coefficients are much lower for $n = 3$ and $n = 6$ in the SARIMA than in the MLR model.
302 Furthermore, the RMSE and AMAPE coefficients are more constant in the SARIMA than in the MLR
303 model. This is an important property that a predictive model should have in order to control the errors
304 associated with different predictions.

305

306 **4. Conclusions**

307 ^7Be forecasting models can be adopted as a target value in analyzing fluctuations or deviations that could
308 imply important atmospheric changes. In this study an explicative and forecasting model of ^7Be air
309 concentrations is proposed, using two different statistical techniques: the SARIMA time series and the MLR
310 model. In both models, the historical data used to develop the model was for the period 2007-2014. The
311 data for the 12 months of the year 2015 was used to measure the validity of the models.

312 Considering the forecasting power measured by the RMSE, MAPE and AMAPE accuracy coefficients, and
313 the simplicity of the model measured by the AIC coefficient, a $\text{SARIMA}(0,1,1)\times(1,1,3)_{12}$ time series is
314 proposed. The analysis of the residuals in the validation step reveals that the model is suitable for
315 forecasting.

316 The MLR model was developed considering the meteorological variables that mainly affect the climatology
317 of Valencia. The significant variables obtained to predict ^7Be activity are: sunspot number, temperature,
318 precipitation and wind speed, which explain only 48.76% of ^7Be variability. The predictive model cannot
319 explain a higher degree of variability of ^7Be activity due to the joint effect of the variables considered,
320 which may mask the intra and inter annual components of the time series. In addition, the analysis of the
321 residuals in the validation step reveals additional structure in the data that the regression did not capture.
322 MLR also has the disadvantage of requiring forecast meteorological parameters to predict ^7Be air
323 concentrations.

324 The comparison between SARIMA and MLR reveals the greater explanatory power of the SARIMA model
325 (70.88%), while its accuracy measurements are consistently lower for both short terms (3-6 months) and
326 long terms (9-12 months) in the out-of-sample period. The MLR model performs well in the long term, but
327 its errors are less consistent in short terms. The proposed SARIMA model can therefore be considered a
328 good forecaster of ^7Be air concentrations. However, the MLR model provides information on significant
329 meteorological variables that affect these concentrations, which could be useful in identifying
330 meteorological or atmospheric changes that could cause deviations in ^7Be concentrations.

331

332 5. Acknowledgements

333 This study has been partially supported by the REM program of the Nuclear Safety Council of Spain
334 (SRA/2071/2015/227.06). We are also grateful to the UPV's weather station for providing the atmospheric
335 information used in this study.

336

337 Bibliography

- 338 Alwan, L.C., Roberts, H.V.1988. Time series modelling for statistical process control. *Journal of Business and Economic Statistics*,
339 6, 87-95.
- 340 Azahra, M., López-Peñalver, J.J., Camacho García, C., González-Gómez, C., El Bardouni T., Boukhal, H. 2004a. Atmospheric
341 concentrations of ^7Be and ^{210}Pb in Granada, Spain. *Journal of Radioanalytical and Nuclear Chemistry*, 261,401-405.
- 342 Baeza, A., Del Río, L.M., Jiménez, A., Miró, C., Paniagua, J.M., Rufo, M., 1996. Analysis of the temporal evolution of atmospheric
343 ^7Be as a vector of the behavior of other radionuclides in the atmosphere. *Journal of Radioanalytical and Nuclear Chemistry*, 207,
344 331-344.
- 345 Bas, M.C., Ortiz, J., Ballesteros, L., Martorell, S. 2016. Analysis of the influence of solar activity and atmospheric factors on ^7Be
346 air concentration by seasonal-trend decomposition. *Atmospheric Environment*, 145, 147-157.
- 347 Bas, M.C., Ortiz, J., Ballesteros, L., Martorell, S. 2017. Forecasting ^7Be concentrations in surface air using time series analysis.
348 *Atmospheric Environment*, 155, 154-161.
- 349 Box, G.E.P., Jenkins, G.M.1976. *Time series analysis: forecasting and control*. San Francisco: Holden Day.
- 350 Ceballos, M.R., Borràs, A., Gomila, E., Estela, J.M., Cerdà, V., Ferrer, L., 2016. Monitoring of ^7Be and gross beta in particulate
351 matter of surface air from Mallorca Island, Spain. *Chemosphere*, 152, 481-489.

352 Dueñas, C., Fernández, M.C., Cabello, M., Gordo, E., Liger, E., Cañete, S., Pérez, M., 2015. Study of the cosmogenic factors
353 influence on temporal variation of ^7Be air concentration during the 23rd solar cycle in Málaga. *Journal of Radioanalytical and*
354 *Nuclear Chemistry*, 303, 2151-2158.

355 Dueñas, C., Fernández, M.C., Liger, E., Carretero, J. 1999. Gross alpha, gross beta activities and ^7Be concentrations in surface air:
356 analysis of their variations and prediction model. *Atmospheric Environment*, 33, 3705-3715.

357 Essaid, C., Piñero-García, F., Ferro-García, M.A., Azahra, M., El Bardouni, T. 2015. Monitoring of ^7Be in Surface air of Granada
358 and their variations with Solar Irradiance and meteorological parameters. In *Proceedings of 4th SEFM-SEPR Congress*, Valencia.

359 Feely, H.W., Larsen, R.J., Sanderson, C.G., 1989. Factors that cause seasonal variations in Beryllium-7 concentrations in surface
360 air. *Journal of Environmental Radioactivity*, 9, 223-249.

361 Hair, J.F., Black, W.C., Babin, B.J., Anderson, R.E. 2010. *Multivariate Data Analysis (7th Edition)*. New Jersey: Pearson Prentice
362 Hall.

363 Harris, T.J., Ross, W.H. 1991. Statistical process control procedures for correlated observations. *Canadian Journal of Chemical*
364 *Engineering*, 69, 48–57.

365 Hötzl, H., Rosner, G., Winkler, R., 1991. Correlation of ^7Be concentrations in surface air and precipitation with the solar cycle.
366 *Naturwissenschaften*, 78, 215-217.

367 Ioannidou, A., Papastefanou, C., 1994. Atmospheric Beryllium-7 concentrations and sun spots. *Nuclear Geophysics*, 8, 539-543.

368 Ioannidou, A., Papastefanou, C., 2006. Precipitation scavenging of ^7Be and ^{137}Cs radionuclides in air. *Journal of Environmental*
369 *Radioactivity*, 85, 121-136.

370 Neroda, A.S., Goncharova, A.A., Goryachev, V.A., Mishukov, V.F., Shlyk, N.V. 2016. Long-range atmospheric transport
371 Beryllium-7 to region the Sea of Japan. *Journal of Environmental Radioactivity*, 160, 102-111.

372 O'Brien, K., 1979. Secular variations in the production of cosmogenic isotopes in the earth's atmosphere. *Journal of Geophysical*
373 *Research*, 84, 423-431.

374 Peña, D. (2010). *Análisis de series temporales*. Madrid: Alianza Editorial.

375 Piñero-García, F., Ferro-García, M.A. 2013. Evolution and solar modulation of ^7Be during the solar cycle 23. *Journal of*
376 *Radioanalytical and Nuclear Chemistry*, 296, 1193–1204.

377 Piñero-García, F., Ferro-García, M.A., Azahra, M. 2012. ^7Be behaviour in the atmosphere of the city of Granada January 2005 to
378 December 2009. *Atmospheric Environment*, 47, 84-91.

379 Schuster, A. 1898. On the investigation of hidden periodicities with application to a supposed 26 day period of meteorological
380 phenomena, *Terrestrial Magnetism*, 3, 13-41.

381 Shine, D.W., Lee, J.H., 2000. Consistency of the maximum likelihood estimators for nonstationary ARIMA regressions with time
382 trends. *Journal of Statistical Planning and Inference*, 87, 55-68.

383 Shumway, R.H., Stoffer, D.S., 2006. *Time series analysis and its applications: With R Examples*. Springer Texts in Statistics. New
384 York: Springer-Verlag.

385 SILSO, 2015. World Data Center e Sunspot Number and Long-term Solar Observations, Royal Observatory of Belgium, On-line
386 Sunspot Number Catalogue. <http://www.sidc.be/SILSO/>.

387 Taylor, A., Keith-Roach, M.J., Iurian, A.R., Mabit, L., Blake, W.H. 2016. Temporal variability of beryllium-7 fallout in southwest
388 UK. *Journal of Environmental Radioactivity*, 160, 80-86.

389 Tsay, R. S. 2005. *Analysis of Financial Time Series*, 2nd ed. New York: Wiley.

390 Vogt, S., Herzog, G.F., Reedy, R.C., 1990. Cosmogenic nuclides in extraterrestrial materials. *Reviews of Geophysics*, 28, 253-275.

391 Wardell, D.G., Moskowitz, H., Plante, R.D. 1994. Run length distributions of special-cause control charts for correlated processes.
392 *Technometrics*, 36, 3-17.

393 Wu, L., Shahidepour, M. 2010. A Hybrid Model for Day-Ahead Price Forecasting. *IEEE Transactions on power systems*, 25,
394 1519-1530.