

Testing machine learning algorithms for the prediction of depositional fluxes of the radionuclides ^7Be , ^{210}Pb and ^{40}K

P. De La Torre Luque^a, C. Dueñas^b, E. Gordo^{c,*}, S. Cañete^c

^a The Oskar Klein Centre for Cosmo Particle Physics, Alba Nova, 10691, Stockholm, Sweden

^b Department of Applied Physics I, Faculty of Sciences, University of Malaga, 29071, Malaga, Spain

^c Central Research Services (SCAI), University of Malaga, 29071, Malaga, Spain

ARTICLE INFO

Keywords:

Natural radionuclides

Depositional fluxes

Machine learning techniques

ABSTRACT

The monthly depositional fluxes of ^7Be , ^{210}Pb and ^{40}K were measured at Malaga, (Southern Spain) from 2005 to 2018. In this work, the depositional fluxes of these radionuclides are investigated and their relations with several atmospheric variables have been studied by applying two popular machine learning methods: Random Forest and Neural Network algorithms. We extensively test different configurations of these algorithms and demonstrate their predictive ability for reproducing depositional fluxes. The models derived with Neural Networks achieve slightly better results, in average, although similar, having into account the uncertainties. The mean Pearson-R coefficients, evaluated with a k-fold cross-validation method, are around 0.85 for the three radionuclides using Neural Network models, while they go down to 0.83, 0.79 and 0.8 for ^7Be , ^{210}Pb and ^{40}K , respectively, for the Random Forest models. Additionally, applying the Recursive Feature Elimination technique we determine the variables more correlated with the depositional fluxes of these radionuclides, which elucidates the main dependences of their temporal variability.

1. Introduction

The use of natural radionuclides as markers for studying the atmospheric circulation provides valuable information about the complex mechanisms involved. It is common to employ different natural radionuclides as tracers and chronometers in aquatic and atmospheric systems (Wogman et al., 1968; Martell, 1970; Schuler et al., 1991) and it was demonstrated to be very useful in studies dedicated to understanding the mechanisms and rates of removal of aerosols (Baskaran et al., 1993). In this work, we aim at the study of a predictive model for the depositions of radionuclides ^7Be , ^{210}Pb , and ^{40}K , whose different origins allow us to infer important features of the atmospheric circulation, erosion processes, transportation and deposition of soils and sediments from episodic to long-term timescales.

^7Be is a cosmogenic radionuclide originated by spallation reactions of cosmic rays with light atmospheric nuclei, such as nitrogen and oxygen (Lal et al., 1958) that has a decay half-life of $T_{1/2} = 53$ days. Thus, this nuclide is mostly produced in the stratosphere and reaches the troposphere in periods of air exchange between these two layers. This is why the production of ^7Be is dependent on altitude, latitude and solar

cycle but has negligible dependence on longitude (Baskaran et al., 1993; Dueñas et al., 2017).

In contrast, ^{210}Pb , with a decay half-life of $T_{1/2} = 22.3$ yr, is produced from the radioactive decay of ^{222}Rn , the only gaseous decay product of ^{238}U series. Therefore, ^{210}Pb is found in larger concentrations near the ground and with important dependence on the distribution of land and seas (Moore et al., 1973; Wilkening et al., 1975; Preiss et al., 1996; García-Orellana et al., 2006).

The atmospheric ^{40}K ($T_{1/2} = 1.3 \cdot 10^9$ yr) is related to a crustal origin, from most kinds of soil, which is usually found in association with other re-suspended materials, as PM_{10} (particulate matter with diameter 10 μm) from the African continent (Karlsson et al., 2008; Dueñas et al., 2011).

Several works in the past have been dedicated to study the relations between the concentrations or depositional fluxes of these radionuclides with different environmental variables for different latitudes and longitudes. In this work, we employ a large dataset (168 monthly measurements, from January 2005 to December 2018) of environmental variables and the fluxes of ^7Be , ^{210}Pb , and ^{40}K radioactivity in the Mediterranean coastal region of Málaga (Southern Spain).

* Corresponding author.

E-mail addresses: pedro.delatorreluque@fysik.su.se (P. De La Torre Luque), mcduenas@uma.es (C. Dueñas), elisagp@uma.es (E. Gordo), scanete@uma.es (S. Cañete).

Similar studies were carried out in the same zone in the past and reported some important results, such as correlations with particulate material (PM₁₀ levels) or with other environmental variables included in this work (Dueñas et al., 2004, 2009, 2011, 2017).

Here, we are exploring new methods of studying the complex relations between the depositional flux of these radionuclides and atmospheric variables, using machine learning algorithms. Machine learning (ML) techniques (Carbonell et al., 1983) provide a promising tool in the prediction of any magnitude which depends on many variables and exhibits complex relations with them. Particularly, we are focused here on the implementation of these methods for the prediction of depositional fluxes of the mentioned radionuclides. These models allow us to identify subtle long-term relationships between the temporal variability of the depositional fluxes and other environmental cycles, like the Solar cycle or atmospheric cycles. Additionally, reproducing these fluxes allow us to discern the real agents affecting the depositions of these radionuclides and could provide another tracer of anomalous (artificial) radiation episodes. In addition, we argue that these kinds of models could be extended to different zones, always that measurements are available, to study relations with other variables not yet considered.

2. Materials and methods

2.1. Study area

Málaga (4° 28' 8" W; 36° 43' 40" N), is the major coastal city in the Andalusian region situated in the south-east of Spain (see Fig. 1), on the Mediterranean coast and, therefore, has a climate influenced by continental and maritime air masses. The predominant winds are easterly (SE) and westerly (NW). The climate is temperate, with contrasting wet (approximately October–April) and dry (approximately May–September) periods (Dueñas et al., 2012).

Due to its geographical proximity to the African continent, our study area is frequently affected by intrusions of air masses with high concentrations of atmospheric particulate matter (Escudero et al., 2005). The sampling point is located on the flat roof of the Central Research Services (SCAI) building at the University of Málaga, at a height of 10 m above the ground and approximately at 5 km from the coastline, near the airport and surrounded by roads with traffic exhaust.

2.2. Data extraction

Bulk deposition samples were collected from January 2005 to December 2018. Samples were collected monthly using a collector that is a slightly tilted stainless steel tray 1 m² in area and a polyethylene vessel of 60 L capacity for a rainwater sample reservoir. A volume of 6 L

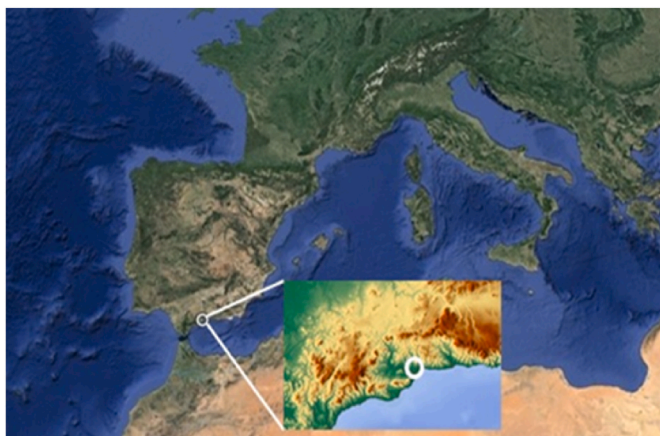


Fig. 1. Physical map showing the location of the study area. The zoomed window shows the exact position of the study area, in Málaga.

of the bulk deposition (the sum of wet deposition flux and the gravitational sedimentation fraction of the dry deposition) was reduced via evaporation to approximately 1 L and transferred to a Marinelli geometry container for gamma counting. The method and processing procedures were described previously (Dueñas et al., 2011).

The atmospheric fluxes were calculated using the expression (1).

$$F = A / St \text{ (Bq m}^{-2} \text{ month}^{-1}\text{)}, \quad (1)$$

where A is the activity in the sample obtained from the gamma spectra, S is the surface area of the collector and t is the duration of sampling time. Additionally, aerosol samples were collected weekly in cellulose filters of 0.8 μm pore size and 47 mm diameter with an air sampler (Radeco, mod AVS-28 A) at a flow rate of 40 L min⁻¹. A monthly composite sample containing 4 or 5 filters (depending on the number of weeks each month) was formed for the gamma analysis.

Radiometric measurements were performed by low-level gamma spectrometry with a coaxial-type germanium detector (Canberra Industries Inc., USA), with a relative efficiency of 20% and it was calibrated using certified reference gamma ray cocktail that contains 8 radionuclides with a range of energy from 59.54 to 1836.05 KeV. Each sample was measured for 172,000 s. Gamma spectra analyses were performed with the Genie2K spectrometry software version 2.0 (Canberra Industries Inc., USA). The characteristic gamma peaks selected for the determination of the different radionuclides were: 477.6 keV for ⁷Be, 1460.81 keV for ⁴⁰K and 46.5 keV for ²¹⁰Pb. To validate the methods, our lab routinely participates in interlaboratory comparisons to measure gamma-emitting radionuclides, in different types of samples, organized by the International Atomic Energy Agency (IAEA), the Joint Research Centre (JRC), and the Spanish Nuclear Safety Council (CSN). Further details of the low-background gamma-ray detection system have been previously described by (Dueñas et al., 1999, 2004).

The meteorological data (temperature, relative humidity, distance travelled monthly by the wind and precipitation) used in this study were obtained from the nearest station network of the Spanish Meteorological Agency (AEMET) (500 m away from the sampling site). Days affected by African dust outbreaks have been obtained from the CALIMA project (www.calima.es). The monthly sunspots number was obtained from NOAA's Space Weather Prediction Center (SWPC).

Additionally, data of daily concentrations of particulate matter fraction PM₁₀ were obtained from Carranque (36° 43' 40" N; 4° 28' 4" W), a monitoring station belonging to the regional Atmospheric Pollution Monitoring network managed by the Environmental Health Service of the Andalusian Government.

The full data-set is available through <https://zenodo.org/record/4685954>.

2.3. Description of the algorithms applied and cross-framework

ML techniques have demonstrated their predictive power in a variety of fields, from medicine (e.g. (Lapedes et al., 1988)), to astrophysics (e.g. (Schaefer et al., 2018; Graff et al., 2014)), used for both classification (as in (Williams et al., 2006)) and numerical forecasting (see, for example (Sarkar et al., 2009; Stencil and Stastny, 2011)). Generally, ML methods are used to find the relation between a set of input variables and an output variable one is interested in. These variables are usually called features and labels, respectively. In the present study, the labels are the monthly depositional fluxes collected from 2005 to 2018 and the features are the atmospheric variables gathered in the same period. Earlier studies have demonstrated that it is possible to find linear relations between atmospheric variables and the depositional fluxes of these radionuclides, although the uncertainties related to this determination become too large to have accurate predictions. Using these methods, we aim at obtaining more precise predictions on the depositional fluxes that could be used, e.g., to reliably detect the emission of artificial radiation or other non-expected radiation sources.

The relation between features and labels is progressively adjusted by iterating over the amount of data samples given to the algorithm, therefore the larger the number of samples used to feed (or train) the algorithm the better the predictions become. The data sample used to adjust the algorithm is called training dataset and this adjustment process is known as the training phase, which basically consists of tuning some training parameters in order to predict the correct labels given. The algorithm adjusts itself in each iteration by comparing its predicted label with the correct label. Then, in order to evaluate the performance of the model one must provide it with new input data (i.e., these features must be different from the training data to ensure unbiased or over-fitted evaluations of the algorithm effectiveness). In this way, we can “grade” or “score” the model performance by comparing the predicted outputs with the real labels in what is called the test phase. The new set of data used in this phase is called test data.

Two different supervised algorithms have been implemented in this study, Neural Networks and Random Forest techniques, and their ability to predict depositional fluxes has been extensively tested for different configurations and for the depositional fluxes of the ^7Be , ^{210}Pb and ^{40}K radionuclides. Very few works have been published using ML techniques to predict depositional fluxes and none of them systematically analyzing their performance. An example of these studies can be ref. (Chham et al., 2018), but deeper research on the efficiency of these techniques is necessary.

The most popular ML algorithm is the Artificial Neural Network (ANN) model. Neural networks can learn complex patterns using layers of neurons which mathematically transform the data. The layers between the input and output are referred to as “hidden layers”. A Neural Network can learn relationships between the features that other algorithms cannot easily discover, including also complex non-linear relations.

Moreover, we used an alternative and less demanding (in terms of resources) technique, the Random Forest algorithm, which, in turn, is not able to consider non-linear features in the relations between the features. This algorithm relies on an ensemble of decision trees which are combined to get averaged predictions. Each tree uses a sub-sample of the full data set, randomly selected, and progressively divides it into different nodes (or leaves) depending on certain quantitative (or qualitative, in case the tree is applied for a classification problem) criteria decided by the algorithm.

We have divided our collected data set into a training set, containing the 80–85% of the full data set, and a test set that allows us to quantify the performance of our predictions. The list of features (meteorological or atmospheric variables) employed is based on monthly averages (or monthly accumulated) and it consists of: Air temperature (in $^{\circ}\text{C}$), relative humidity level (%), number of days affected by African dust outbreaks (intrusions), distance travelled monthly by the wind (in km), pressure (hPa), sunspot number, amount of rainfall (dm), PM_{10} level ($\mu\text{g m}^{-3}$), seasonal factor (from 1, for winter, to 4, for spring), monthly factor (from 1, for January, to 12, for December), total rainfall duration (min), humid days, dry days and time between rains (in days). For both algorithms, the labels (depositional fluxes) are normalized, since this allows a better performance of the algorithm.

A Neural Network in which the input features first result into 8 units (1st hidden layer) and then into 4 units (second hidden layer) have been found to be the most adequate, as it is depicted in the appendix A. The implementation of the Neural Network has been achieved by using the Python Keras (Chollet, 2015) library. The connections between the input features and the first hidden layer, as well as between the first and second hidden layers use the Rectified Linear Unit (ReLU) as activation function and the connections from the second hidden layer and the output units are calculated with a linear activation function.

The model performance was optimized including a step of batch normalization and dropout (finding the best results adjusting it to the 10% of the sample) after each of the hidden layers. In addition, the adaptive moment estimation optimizer, or Adam optimizer, was found

to get the best performance for every one of the radionuclides. On top of this, the best results were found when taking the natural logarithm of the values for the features, as expected, and setting the mean absolute error metrics as the loss function.

Different configurations of the neural networks models and the hyperparameters involved (i.e., the values needed to control the learning process in ML algorithms) were refined by applying a simple random search method (Bergstra and Bengio, 2012). This technique consists of comparing the performance of the algorithm for different combinations of hyperparameters (in an equally spaced grid of values) and choosing the set of hyperparameters that provides the best performance. In Table 1, we show the main hyperparameters tuned for the Adam optimizer for each radionuclide. The rest of hyperparameters needed by the optimizer were set to their default values given by the keras method.

For the Random Forest algorithm, it was found that using the features values normalized (i.e. subtracting the mean and dividing by the standard deviation of the measurement), instead of their natural logarithm, gave better results. Then, the main hyper-parameters were adjusted for each of the nuclides, setting the mean absolute error (MAE) as criterion for splitting the nodes and a minimum number of samples required to split an internal node (min samples split) to 3. The number of decision trees (also known as number of estimators) used in the model was set to be 680 for ^7Be and ^{210}Pb and 280 for ^{40}K .

The results from both algorithms and for the three radionuclides are shown and compared in the next section, in which we fully demonstrate their ability for reproducing the data and systematically explore the statistical errors around these predictions as well as the main features involved.

3. Results

As a first step before running our models, we randomly shuffle the features and labels and then, they are divided into a training and a test set. Once the model is trained, we rate its performance by comparing the predictions with the test labels, corresponding to a 15–20% of the full data sample, using the mean percentage error and the Pearson-R index value. While the former is an indicator of the quantitative differences between test labels and predictions, the latter is a good indicator of the trend similarities between the two sets.

In order to compare these results with a reference model, we applied the same kind of evaluation to the model found in the linear regression analysis presented in ref (Dueñas et al., 2017). for the ^7Be radionuclide. This analysis yields a linear relation (see equation 2) between the depositional flux of ^7Be and the amount of rainfall, in L (the variable which shows the largest correlation with the depositional flux of every radionuclide).

$$\text{FluxBe} = 6.33 + 2.6 \times \text{rainfall} (\text{Bq m}^{-2} \text{ month}^{-1}), \text{sl} \quad (2)$$

This comparison is carried out by using a portion of 25 randomly selected measurements (similar to the number of samples in the test sets used for the ML algorithms applied) of ^7Be and amount of rainfall (corresponding to the same order to have a robust idea on the value of these metrics, we repeated this for 100 times (analogous to what is done in section IV A), with different randomly selected samples of 25 measurements, and computed the average value. These metrics result in a mean R index of $\sim 0.45 \pm 0.4$ and a maximum R index of 0.95, while the

Table 1

Main hyperparameters (i.e., the values needed to control the learning process in ML algorithms) used in the Adam optimizer, adjusted for each of the radionuclides studied.

	^7Be	^{210}Pb	40K
Learning raterowhead	$2.1 \cdot 10^{-3}$	$2.1 \cdot 10^{-3}$	$2.2 \cdot 10^{-3}$
Decay raterowhead	$5 \cdot 10^{-6}$	$5 \cdot 10^{-5}$	$2.4 \cdot 10^{-6}$

mean percentage errors were of $103 \pm 150\%$. Having these reference metric values is necessary to compare to the quantitative results of the Random Forest and Neural Network algorithms studied here. In Fig. 8 (appendix B), we display the comparison between the predictions from the reference model and the depositional flux measurements for one of these samples.

In Fig. 2, we show some of the best results acquired from the Neural Network and Random Forest algorithms for all the studied radionuclides, which demonstrates that these algorithms can allow us to significantly improve our predictions on depositional fluxes with respect to traditional methods. Here, we highlight that these are predictions obtained from their corresponding atmospheric variables and remark the importance of evaluating these predictions with data not used for the training phase, since this highly biases our evaluation. As we can see by the Pearson-R value, these predictions are able to suitably reproduce the labels trend with respect to the atmospheric variables. In addition, we find mean absolute errors of the order 50% usually, which are well below the error levels found using linear regressions (as shown above) and are similar to the experimental uncertainties in the determination of these fluxes, which can be O (10%), as shown in refs. (Herranz et al., 2008; Heydorn, 2004). In this case, it has been observed that high-flux values are difficult to be matched, which may be related to periods of anomalous radiation doses. Nevertheless, this requires a dedicated study of those points and their temporal behavior, which is beyond the scope of this paper. Further sources of uncertainty in these comparisons mainly come from the statistical uncertainties related to the measurement of the atmospheric variables and variables not included in the model.

Surprisingly, the models make good predictions also for the ^{40}K nuclide, even with a considerably smaller number of samples available for it. On top of this, we found that the absolute percentage errors follow a similar distribution for each radionuclide and both algorithms. They are well described with a Gamma probability distribution, which exhibits a slightly negative mode and a slightly positive median. This is likely since the distribution of depositional fluxes is also very well reproduced with a Gamma function. A representative example of these distributions for the Neural Network and Random Forest algorithms is shown in Fig. 3 for the ^7Be radionuclide after gathering several

repetitions for different test sets used. The fact that these errors follow such distribution can be used to statistically diagnose anomalous episodes of radiation doses. We noticed that the Random Forest models produce slightly larger median values and mode values more deviated from 0, but no significant differences between the same algorithms for different nuclides was detected.

Nevertheless, the evaluation of the models is highly dependent on the data set used. From one side, the larger the test set, the more reliable is the model performance evaluation, but at the cost of reducing the number of samples used in the training set. On the other side, if the test set is too short, the model performance evaluation will be very uncertain. In this case, we observed that using around 20% of the full data set allowed us to make consistent evaluations. Even though they are still short enough to make our evaluation very dependent on the data test used.

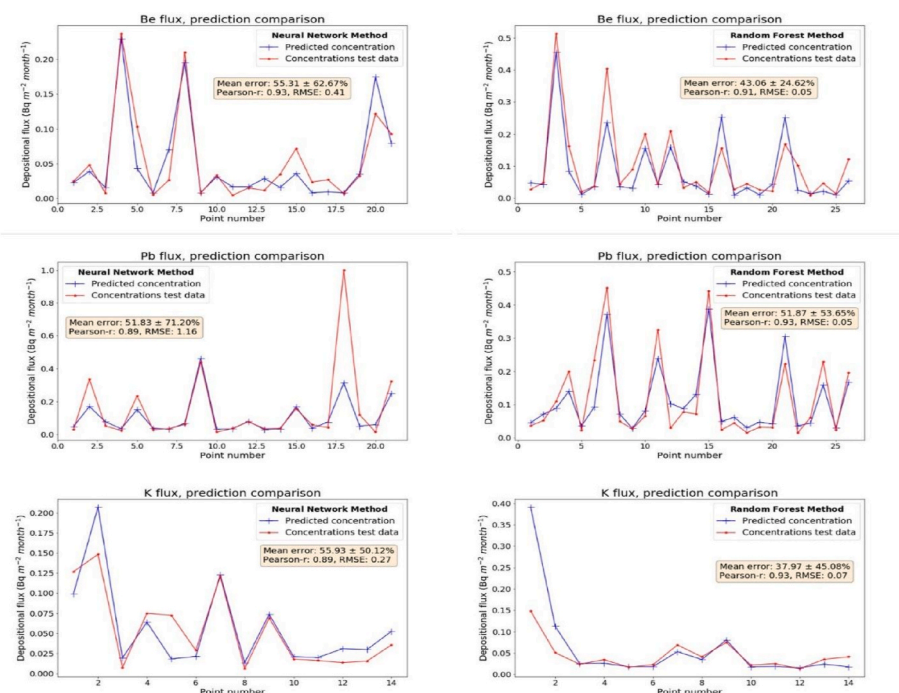
This issue is well known by the ML community and there are many possible strategies to deal with it to have an unbiased evaluation of our model (Raschka, 2018) and its predictions uncertainties, as it is explored in the next section.

3.1. Statistical evaluation

To prevent biasing our model evaluation by the small amount of test data used and take into account the full uncertainty involved; we evaluate the algorithms by means of a k-fold procedure. In this process the data set is divided into k subsets. Each time, one of the k subsets is used as a test set and the other k-1 subsets form the training set. Then, we statistically combine the results to get solid conclusions.

At this point, another difference between the Neural Network and the Random Forest algorithms should be taken into account to correctly manage the full uncertainties involved: while the training process exactly results in the same model for the Random Forest algorithm, this is subject to further fluctuations in the Neural Network algorithm. This is due to the optimization procedure necessary for finding the minimum error or loss when evaluating the examples in the training dataset. The main problems usually faced are getting stacked in local minimal or local optima (i.e., regions where the loss is relatively low, but it is not the lowest), saddle or flat points (regions where adjustments of the training

Fig. 2. Example of the results of the predictions found from the Neural Network (left panels) and Random Forest (right panels) models for random test samples. These predictions are limited to the test sample, which is chosen to be around 20% of the full data set. We also include the values of the metrics used to evaluate the predictive ability of these methods, which are the Pearson-R correlation coefficient and the mean absolute error and its standard deviation. The root mean square error (RMSE), in units of $\text{Bq m}^{-2} \text{ month}^{-1}$, is also included for completeness.



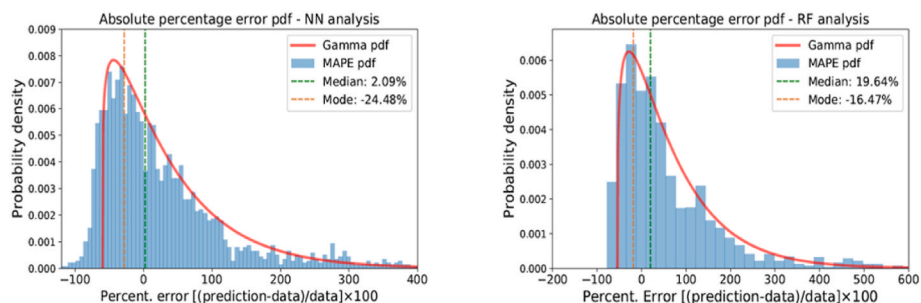


Fig. 3. Probability distribution for the percentage errors found for various evaluations with (around 20) different test sets. The left plot shows the results of these evaluations for the Neural Network algorithm and the right plot those for the Random Forest algorithm.

weights do not lead to an appreciable change in the loss) and other issues more related with the loss function, gradients and the dimensionality involved. More precise information about these problems can be found, e.g., in (Bengio et al., 2012).

Therefore, each time the Neural Network is trained, specially when the number of samples is not large enough, it is subject to small variations in the model predictions. For this reason, a good evaluation of the uncertainties involved in the predictions of the Neural Network model requires adding these fluctuations.

We repeated the training and test phase for 5 times with the same test and training datasets. Then, we perform the evaluations with 20 different randomly selected test sub-datasets following the k-fold procedure. This means that we carry out a total of 100 training and evaluation steps to determine the Pearson-R value and the mean percentage error of our predictions with respect to the experimental data, as well as the uncertainties related to these determinations for the Neural Network model. In turn, as the Random Forest algorithm does not suffer from those training fluctuations, we performed 60 evaluations of the model, employing a different test and training subsets, accordingly, in each evaluation.

These results are shown in Fig. 4, where we represent the mean Pearson-R index values and the 1σ uncertainty related to its

determination for both, the Neural Networks and Random Forest algorithms and for the three nuclides with respect to the number of iterations employed in the training phase. In general, we observe that the mean Pearson-R index values are larger for the ⁷Be and ²¹⁰Pb radionuclides than for ⁴⁰K, although ⁴⁰K shows larger errors due to the smaller number of samples available. In addition, the uncertainties related to the determination of the R index value from the Random Forest algorithm is slightly larger than that from the NN algorithm. The mean Pearson-R index values obtained are between 0.75 and 0.88 for ⁷Be and ²¹⁰Pb, but around 0.7–0.8 for ⁴⁰K, although the errors are still high for every radionuclide. In particular, the determination of ⁷Be seems to be the most accurate in general, showing a 1σ uncertainty in the determination of the R index value around ±0.065 for the NN algorithm and ±0.08 for the RF algorithm. A maximum mean R index value of around 0.87 and 0.88 are found for ⁷Be and ²¹⁰Pb, respectively, at 1400 and 1300 iterations. The maximum mean R index value obtained for ⁴⁰K is slightly above 0.8, found with the RF algorithm.

As expected, the performance of these methods in reproducing depositional fluxes improves when having more samples, obtaining larger Pearson-R index values and lower uncertainties related. Nevertheless, we observed that the NN algorithm seems to accuse the smaller number of samples with respect to the RF technique.

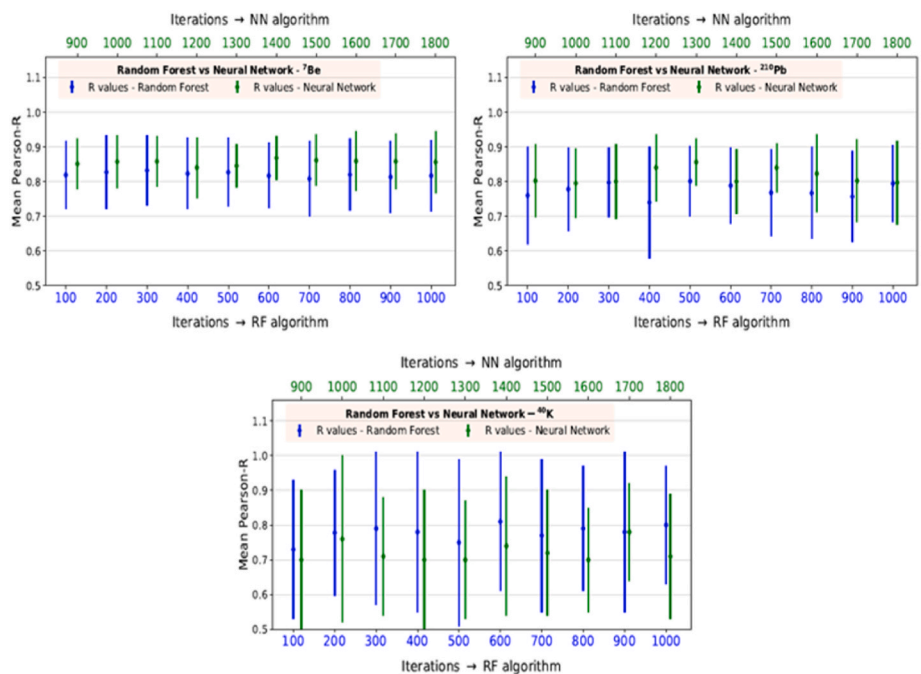


Fig. 4. Results from the k-fold evaluation of the Pearson-R correlation coefficient for the Neural Network and Random Forest algorithms for the depositional fluxes of ⁷Be (upper panel), ²¹⁰Pb (middle panel) and ⁴⁰K (lower panel). The results obtained from the NN algorithm are shown in green while the results from the RF algorithm are shown in blue. (For interpretation of the references to colour in this figure legend, the reader is referred to the Web version of this article.)

3.2. Selecting the main variables

To fully exploit the capability of ML techniques in improving our predictions in the depositional fluxes, we determined which are the most important features using the recursive feature elimination algorithm (RFE), which allows us to reduce the complexity and needed CPU time of the Neural Network and Random Forest algorithms and prevents from over-fitting our results. In addition, we compared the results obtained with these features with those obtained when using all the features. Specifically, we used the RFECV method from the sklearn.features election python package. The RFE algorithm is a feature selection method that allows a model to progressively eliminate the weakest features and find the best scoring combination of features.

In Fig. 5 we show the optimal important features found by the RFE algorithm, along with their relative importance. As expected, the rainfall duration and rainfall volume are selected by the three radionuclides. Then, we observe that other atmospheric variables are present, such as the number of humid or dry days, the average monthly pressure, or the mean air temperature. On the other hand, the PM₁₀ level and sunspot number are selected as important for the ⁴⁰K nuclide.

The fact that the sunspot number arises as one of the most important variables describing the depositional fluxes of ²¹⁰Pb and ⁴⁰K is unexpected and there is no current theoretical explanation for it. In principle, this variable is expected to be relevant to produce ⁷Be since it is related with the solar activity (this is, the Sun's magnetic field), which plays an important role on the flux of cosmic rays reaching the atmosphere (Yoshimori et al., 2003). This fact is probably due to the mild correlations between sunspot number and other atmospheric variables, but more data samples are needed to get a solid conclusion, since the sunspot number follows cycles of 11 and 22 years, following the solar magnetic cycles (Cliver, 2015). This could be explained by the fact that there are other correlations found between the solar cycle and other atmospheric variables, as the atmospheric temperature (Qu et al., 2012) and correlations with the cosmic-ray intensity at Earth, which is known to be related to climate and directly involved in processes of cloud formation (Veretenenko et al., 2018; Svensmark et al., 2013; Marsh and Svensmark, 2000). This motivates further campaigns of data measurements to reveal the nature of this correlation.

Once these features have been selected, we proceed to compare the NN and RF algorithms explored in this work using all the features and using just the important features, as displayed in Fig. 6. From this figure, we can see that the NN models for ⁴⁰K have significantly improved, restricting our features to be just the important ones. This means that some of the eliminated features were over-fitting the model. This can be related to the fact that this radionuclide actually comes from African zones and reaches coastal zones of Southern Spain after it is transported by winds in the correct direction. Therefore, some of the atmospheric variables measured in the zone of Malaga could not be suitable to describe its amount and depositions in Malaga. Even so, the amount of rainfall should still be crucial to make the African dust to definitely fall in the study region. Furthermore, the presence of the sunspot number as an important feature has not been pointed out in the past (Dueñas et al., 2015), which may mean that there are other atmospheric variables with a considerable role in the amount and depositional flux of ⁴⁰K found in the Mediterranean coastal zone of Southern Spain.

On the other hand, we see that for ⁷Be and ²¹⁰Pb the results remain very similar to the case with all the features, which is quite remarkable given the number of variables needed. In addition, the uncertainties related to the determination of the Pearson-R correlation coefficient have been considerably reduced in the NN models for ⁴⁰K, while they seem to be almost identical for all other cases.

In general, these results show that the use of these ML methods allow our predictions to be more complex and better adapt to the variability related to the depositional fluxes of different radionuclides.

4. Discussion

Modern computer algorithms allow us to refine our measurements and model predictions via new statistical tools or artificial intelligence. In this work, we have made use of two common machine learning algorithms, Neural Networks and Random Forests, in order to predict and analyze the depositional fluxes of ⁷Be, ²¹⁰Pb and ⁴⁰K. This work has shown, first, that these methods can be successfully applied to study the depositional fluxes of different radionuclides from atmospheric variables such as the amount of rainfall, pressure, or air temperatures. Second, we have evaluated the performance of these models using a k-

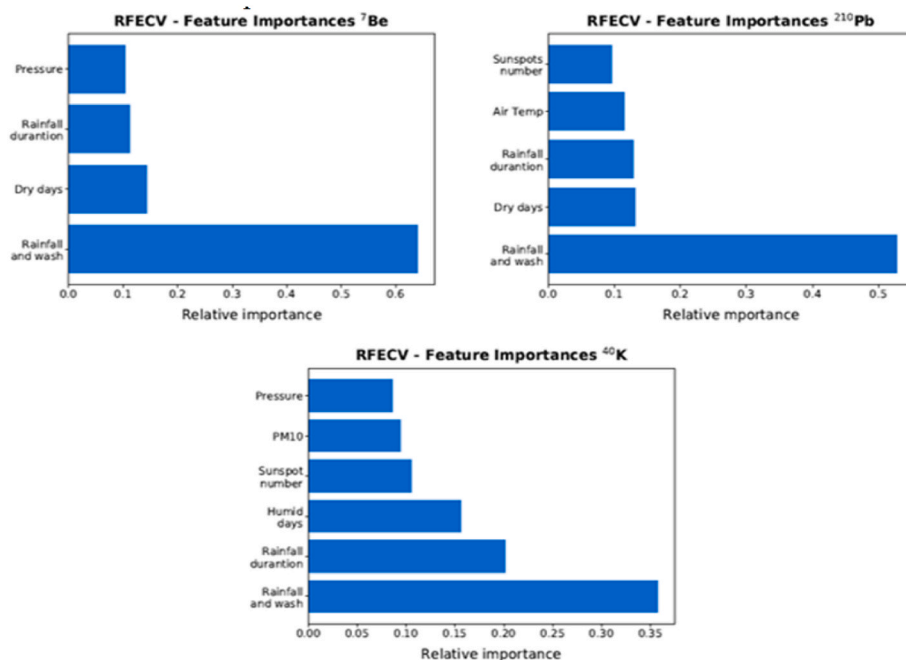


Fig. 5. Histograms with the important features found with the implemented recursive feature elimination algorithm for the depositional fluxes of ⁷Be (upper panel), ²¹⁰Pb (middle panel) and ⁴⁰K (lower panel) with their relative importance.

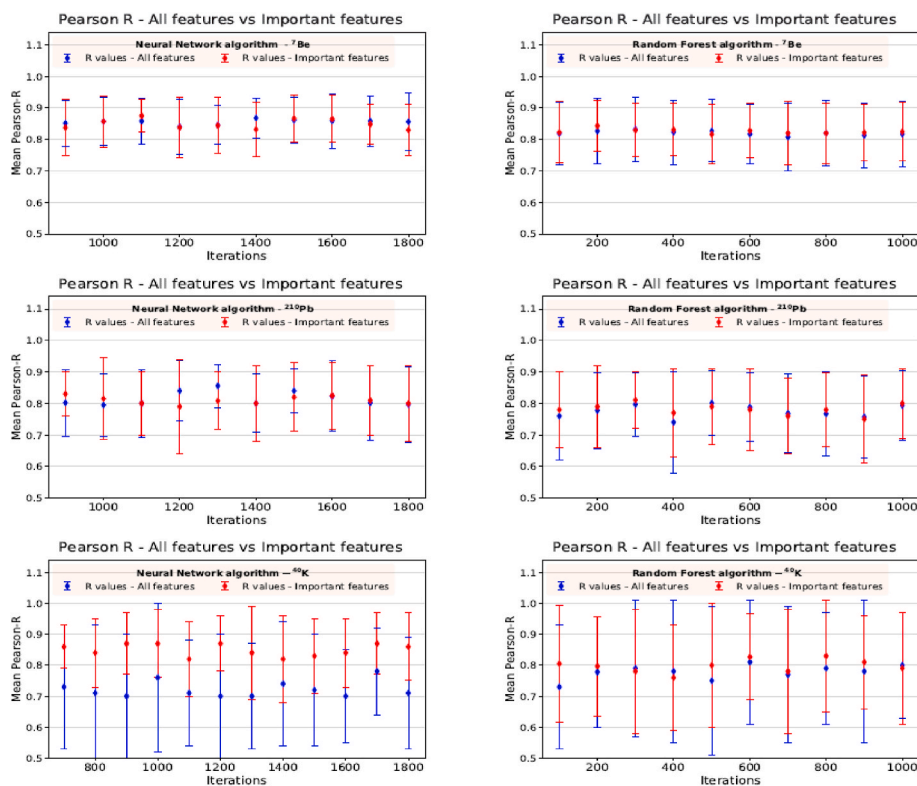


Fig. 6. Same as in Fig. 4 but comparing now the results obtained using the main variables obtained from the RFE algorithm and those obtained from the models trained with all the available variables in the data set.

fold method and the Pearson-R coefficient and mean absolute error as metrics finding that these techniques can significantly improve old predictions made from multivariate linear regression analyses.

As expected, the performance of these methods in reproducing depositional fluxes improves when having more samples, obtaining larger Pearson-R index values and lower uncertainties related. This, in fact, confirms the prospects on future models, with a larger number of samples measured. This is mainly related to the long times involved in the natural cycles of atmospheric variables, as, for example, the sunspot number, which is known to follow 11 or 22-years periods (solar magnetic cycles). Nonetheless, we have demonstrated that the algorithms employed here are able to reproduce the experimental depositional fluxes using monthly-averaged variables and that these predictions can help identify periods of anomalous radiation doses. Interestingly, we found that both the depositional fluxes of ^{210}Pb and ^{40}K , seem to be correlated with the Sunspot number.

The Neural Network models seem to reach higher mean Pearson-R index values, calculated using a k-fold cross-validation treatment, almost reaching 0.9, although the uncertainties are still quite high. Furthermore, the use of a Recursive Feature Elimination algorithm has been used to find the variables that perform the best predictions and allow us to reduce to 4, 5 and 6 the number of variables used for predicting the depositional fluxes of ^{7}Be , ^{210}Pb and ^{40}K , respectively. The training of the Neural Network and Random Forest models with these variables resulted in a negligible difference in the Pearson-R index values and the uncertainties related to its determination except for the ^{40}K nuclide in the Neural Network model, which showed a significant improvement. Even with this reduced number of variables used for training our methods, we were able to obtain mean values for the Pearson-R index value above 0.80 for all the three nuclides and both algorithms. A maximum mean R index value around 0.87 is found for ^{7}Be , ^{210}Pb and ^{40}K , respectively, at 1400, 1300 and 1200 iterations for the Neural Network method. For the Random Forest method, the maximum mean R index value of sim0.81 is found around 500 and 600

iterations for ^{210}Pb and ^{40}K and of almost 0.85 for the ^{7}Be radionuclide.

5. Conclusions

In conclusion, we demonstrate that Random Forest and Neural Networks methods are able to improve our current knowledge and predictions on the depositional fluxes of radionuclides in the Mediterranean coastal zone of Malaga and these models can be extended to other zones too, in order to build a more complex ensemble that could refine the existent knowledge on deposition of different radionuclides. Thus, this work constitutes the first step into the study of a large-scale (in terms of geographical areas) model able to make predictions on depositional fluxes for different geographical zones thanks to the adaptability of these algorithms. The implementation of a recurrent neural network applied to the prediction of depositional fluxes can improve these models and will be also investigated in a next work.

Author contributions

Conceptualization, P. De La Torre Luque, C. Dueñas, E. Gordo, and S. Cañete.; methodology, P. De La Torre Luque, C. Dueñas, E. Gordo, and S. Cañete; formal analysis, P. De La Torre Luque.; investigation P. De La Torre Luque, C. Dueñas, E. Gordo, and S. Cañete.; data curation, P. De La Torre Luque, E. Gordo; writing—original draft preparation, P. De La Torre Luque, E. Gordo; writing—review and editing P. De La Torre Luque, C. Dueñas, E. Gordo, and S. Cañete; project administration, C. Dueñas; funding, C. Dueñas. All authors have read and agreed to the published version of the manuscript.

Funding

This research was funded by Consejo de Seguridad Nuclear (Spain).

Declaration of competing interest

The authors declare no conflict of interest. The funders had no role in the design of the study; in the collection, analyses, or interpretation of data; in the writing of the manuscript; or in the decision to publish the

Appendix A

In this appendix, we show a sketch of the general structure of the Neural Network model employed and an example of a branch of a decision tree from the Random Forest algorithm investigated in this work.

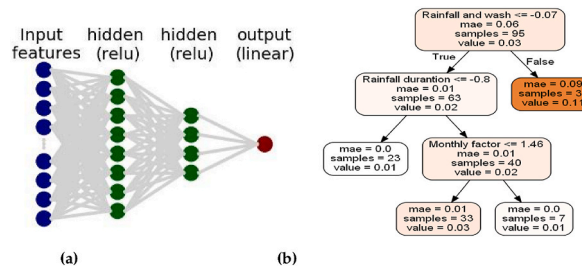


Fig. A1. a): Sketch of the Neural Network model used, where there are two hidden layers that use the ReLU activation function and an output unit that linearly combines the nodes of the last hidden layer. b): Example of a decision tree used as part of a Random Forest model.

Appendix B

This appendix shows a comparison between the predictions from the reference model and the depositional flux measurements for one of these samples. It is crucial to have a reference model evaluated in the same way as for the ML algorithms studied in the paper, since this kind of evaluation is rather peculiar from ML algorithms. As we see, traditional models, based in linear regressions, are unable to reproduce the depositional fluxes behavior, because of the complex relationships between variables.

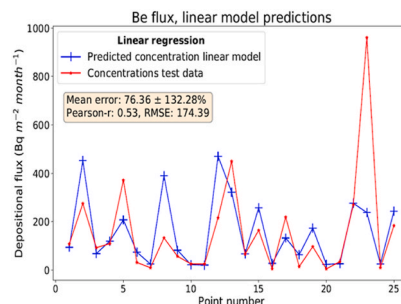


Fig. A2. Predictions found from the reference linear model on one of the 25-length data samples, using the same evaluation as for the Random Forest and Neural Network algorithms studied in this work. Units of RMSE are of Bq m⁻² month⁻¹.

References

- Baskaran, M., Coleman, C.H., Santschi, P.H., 1993. Atmospheric depositional fluxes of ⁷Be and ²¹⁰Pb at Galveston and college station, Texas. *J. Geophys. Res.: Atmosphere* 98 (D11), 20555–20571.
- Bengio, Y., 2012. Practical recommendations for gradient-based training of deep architectures. In: Montavon, Grégoire, Orr, Genevieve B., Klaus-Robert Müller (Eds.), *Neural Networks: Tricks Of The Trade - Second Edition*, Volume 7700 of *Lecture Notes In Computer Science*, pp. 437–478.
- Bergstra, J., Bengio, Y., 2012. Random search for hyper-parameter optimization. *J. Mach. Learn. Res.* 13, 281–305.
- Carbonell, J.G., S Michalski, R., M Mitchell, T., 1983. Machine learning: a historical and methodological analysis. *AI Mag.* 4 (3), 69, 69.
- Chhañ, E., Piñero-García, F., Brattich, E., El Bardouni, T., Ferro-García, M.A., 2018. ⁷Be spatial and temporal pattern in southwest of Europe (Spain): evaluation of a predictive model. *Chemosphere* 205, 194–202.
- Chollet, F., 2015. *Keras*. <https://github.com/fchollet/keras>.
- Clinger, E.W., 2015. *The Extended Cycle of Solar Activity and the Sun's 22-year Magnetic Cycle*, vol. 53. *Space Sciences Series of ISSI*.
- Dueñas, C., Fernández, M.C., Liger, E., Carretero, J., Gross, alpha, 1999. Gross beta activities and ⁷Be concentrations in surface air: analysis of their variations and prediction model. *Atmos. Environ.* 33 (22), 3705–3715.
- Dueñas, C., Fernández, M.C., Carretero, J., Liger, E., Cañete, S., 2004. Long-term variation of the concentrations of long-lived Rn descendants and cosmogenic ⁷Be and determination of the MRT of aerosols. *Atmos. Environ.* 38 (9), 1291–1301.
- Dueñas, C., Fernández, M.C., Cañete, S., Pérez, M., 2009. ⁷Be to ²¹⁰Pb concentration ratio in ground level air in Málaga (36.7°N, 4.5°W). *Atmos. Res.* 92 (1), 49–57.
- Dueñas, C., Fernández, M.C., Gordo, E., Cañete, S., Pérez, M., 2011. Gross alpha, gross beta activities and gamma emitting radionuclides composition of rainwater samples and deposition to ground. *Atmos. Environ.* 45 (4), 1015–1024.
- Dueñas, C., Fernández, M.C., Gordo, E., Cañete, S., Pérez, M., 2012. Chemical and radioactive composition of bulk deposition in Málaga (Spain). *Atmos. Environ.* 62, 1–8.
- Dueñas, C., Fernández, M.C., Cabello, M., Gordo, E., Liger, E., Cañete, S., Pérez, M., 2015. Study of the cosmogenic factors influence on temporal variation of ⁷Be air concentration during the 23rd solar cycle in Málaga (South Spain). *J. Radioanal. Nucl. Chem.* 303, 2151–2158.
- Dueñas, C., Gordo, E., Liger, E., Cabello, M., Cañete, S., Pérez, M., 2017. P. de la Torre Luque. ⁷Be, ²¹⁰Pb and ⁴⁰K depositions over 11 years in Málaga. *J. Environ. Radioact.* 178–179, 325–334, 11.
- Escudero, M., Castillo, S., Querol, X., Ávila, A., Alarcón, M., Viana, M., Alastuey, A., Cuevas, E., Rodríguez, S., 2005. Wet and dry African dust episodes over eastern Spain. *J. Geophys. Res.* 110, 18, 8, 09.
- García-Orellana, J.A., Sánchez-Cabeza, P., Masqué, A., Ávila, E., Costa, M.D., Loyé-Pilot, Bruach-Menchén, J.M., 2006. Atmospheric fluxes of ²¹⁰Pb to the western Mediterranean Sea and the Saharan dust influence. *J. Geophys. Res.: Atmosphere* 111 (D15).
- Graff, P., Feroz, F., Hobson, M.P., SkyNet, A. Lasenby, 2014. An efficient and robust neural network training tool for machine learning in astronomy. *Mon. Not. Roy. Astron. Soc.* 441 (2), 1741–1759, 05.

- Herranz, M., Idoeta, R., Legarda, F., 2008. Evaluation of uncertainty and detection limits in radioactivity measurements. *Nucl. Instrum. Methods Phys. Res. Sect. A Accel. Spectrom. Detect. Assoc. Equip.* 595, 526–534, 10.
- Heydorn, K., 2004. Evaluation of the uncertainty of environmental measurements of radioactivity. *J. Radioanal. Nucl. Chem.* 262, 249–253, 10.
- Karlsson, L., Hernández, F., Rodríguez, S., López-Pérez, M., Hernández-Armas, J., Alonso-Pérez, S., Cuevas, E., 2008. Using ^{137}Cs and ^{40}K to identify natural Saharan dust contributions to PM_{10} concentrations and air quality impairment in the Canary Islands. *Atmos. Environ.* 42 (30), 7034–7042.
- Lal, D., Malhotra, P.K., Peters, B., 1958. On the production of radioisotopes in the atmosphere by cosmic radiation and their application to meteorology. *J. Atmos. Terr. Phys.* 12 (4), 306–328.
- Lapedes, A., Barnes, C., Burks, C., Farber, R., Sirotkin, K., 1988. Application of Neural Networks and Other Machine Learning Algorithms to DNA Sequence Analysis. Los Alamos National Lab., NM (USA). *Technical report*.
- Marsh, Nigel, Svensmark, Henrik, 2000. Cosmic rays, clouds, and climate. *Space Sci. Rev.* 94, 215–230, 11.
- Martell, E.A., 1970. Transport Patterns and Residence Times for Atmospheric Trace Constituents vs. Altitude, ume 93. ACS Publications.
- Moore, H.E., Poet, S.E., Martell, E.A., 1973. ^{222}Rn , ^{210}Pb , ^{210}Bi , and ^{210}Po profiles and aerosol residence times versus altitude. *J. Geophys. Res.* 78 (30), 7065–7075, 1896–1977.
- Preiss, N., Mélière, M.A., Pourchet, M., 1996. A compilation of data on lead 210 concentration in surface air and fluxes at the air-surface and water-sediment interfaces. *J. Geophys. Res.* 101, 862. D22:28,847–28.
- Qu, W., Zhao, J., Huang, F., Deng, S., 2012. Correlation between the 22-year solar magnetic cycle and the 22-year quasi cycle in the earth's atmospheric temperature. *Astron. J.* 144 (1), 6.
- Raschka, S., 2018. Model Evaluation, Model Selection, and Algorithm Selection in Machine Learning, p. 11. *ArXiv preprint*. <https://arxiv.org/abs/1811.12808v3>.
- Sarkar, K., Ben Ghalia, M., Wu, Z., Bose, S.C., 2009. A neural network model for the numerical prediction of the diameter of electro-spun polyethylene oxide nanofibers. *J. Mater. Process. Technol.* 209 (7), 3156–3165.
- Schaefer, C., Geiger, M., Kuntzer, T., Kneib, J., 2018. Deep convolutional neural networks as strong gravitational lens detectors. *A&A* 611, A2.
- Schuler, C., Wieland, E., Santschi, P.H., Sturm, M., Lueck, A., Bollhalder, S., Beer, J., Bonani, G., Hofmann, H.J., Suter, M., Wolfli, W., 1991. A multitracer study of radionuclides in lake Zurich, Switzerland: 1. comparison of atmospheric and sedimentary fluxes of ^7Be , ^{10}Be , ^{210}Pb , ^{210}Po , and ^{137}Cs . *J. Geophys. Res.: Oceans* 96 (C9), 17051–17065.
- Stencl, M., Stastny, J., 2011. Artificial neural networks numerical forecasting of economic time series. 15+. *IntechOpen* 4.
- Svensmark, H., Enghoff, M.B., Olaf Pepke Pedersen, J., 2013. Response of cloud condensation nuclei (≥ 50 nm) to changes in ion-nucleation. *Phys. Lett.* 377 (37), 2343–2347.
- Veretenenko, S., Ogurtsov, M., Lindholm, M., Jalkanen, R., 2018. Galactic cosmic rays and low clouds: possible reasons for correlation reversal. In: Szadkowski, Zbigniew (Ed.), *Cosmic Rays*. *IntechOpen, Rijeka* (chapter 5).
- Wilkening, M.H., Clements, W.E., Stanley, D., 1975. Radon 222 flux measurements in widely separated regions. *Natural radiation environment II* 717–730.
- Williams, N., Zander, S., Armitage, G., 2006. A preliminary performance comparison of five machine learning algorithms for practical Ip traffic flow classification. *SIGCOMM Comput. Commun. Rev.* 36 (5), 5–16.
- Wogman, N.A., Thomas, C.W., Cooper, J.A., Engelmann, R.J., Perkins, R.W., 1968. Cosmic ray-produced radionuclides as tracers of atmospheric precipitation processes. *Science* 159 (3811), 189–192.
- Yoshimori, M., Hirayama, H., Mori, S., Sasaki, K., Sakurai, H., 2003. Be-7 nuclei produced by galactic cosmic rays and solar energetic particles in the earth's atmosphere. *Adv. Space Res.* 32 (12), 2691–2696.