ORIGINAL PAPER

# Forecasting *Betula* and Poaceae airborne pollen concentrations on a 3-hourly resolution in Augsburg, Germany: toward automatically generated, real-time predictions

Anna Muzalyova · Jens O. Brunner · Claudia Traidl-Hoffmann · Athanasios Damialis

**Abstract** Airborne allergenic pollen impact the health of a great part of the global population. Under climate change conditions, the abundance of airborne pollen has been rising dramatically and so is the effect on sensitized individuals. The first line of allergy management is allergen avoidance, which, to date, is by rule achieved via forecasting of daily pollen concentrations. The aim of this study was to elaborate on 3-hourly predictive models, one of the very few to the best of our knowledge, attempting to forecast pollen concentration based on near-real-time automatic pollen measurements. The study was conducted in Augsburg, Germany, during four years (2016–2019) focusing on *Betula* and Poaceae pollen, the most abundant and allergenic in temperate climates. ARIMA and dynamic regression models were employed, as well as machine learning techniques, viz. artificial neural networks and neural network autoregression models. Air temperature, relative humidity, precipitation, air pressure, sunshine duration, diffuse radiation, and wind speed were additionally considered for the development of the models. It was found that air temperature and precipitation were the most significant variables for the prediction of airborne pollen concentrations. At such fine temporal resolution, our forecasting models performed well showing their ability to explain most of the variability of pollen concentrations for both taxa. However, predictive power of *Betula* forecasting model was higher achieving $R^2$ up to 0.62, whereas Poaceae up to 0.55. Neural autoregression was superior in forecasting

A. Muzalyova
University Hospital of Augsburg, Augsburg, Germany

J. O. Brunner
Chair of Health Care Operations/Health Information Management, Faculty of Business and Economics, University of Augsburg, Augsburg, Germany

C. Traidl-Hoffmann · A. Damialis (✉)
Department of Environmental Medicine, Faculty of Medicine, University of Augsburg, Neusaesser Str. 47, Augsburg 86156, Germany
e-mail: thanos.damialis@tum.de

C. Traidl-Hoffmann · A. Damialis
Chair of Environmental Medicine, Technical University of Munich, Augsburg, Germany

C. Traidl-Hoffmann · A. Damialis
Institute of Environmental Medicine, Helmholtz Centre Munich, Research Center for Environmental Health, Augsburg, Germany

C. Traidl-Hoffmann
Christine Kühne - Center for Research and Education (CK-CARE), Davos, Switzerland

*Betula* pollen concentrations, whereas, for Poaceae, seasonal ARIMA performed best. The good performance of seasonal ARIMA in describing variability of pollen concentrations of both examined taxa suggests an important role of plants' phenology in observed pollen abundance. The present study provides novel insight on per-hour forecasts to be used in real-time mobile apps by pollen allergic patients. Despite the huge need for real-time, short-term predictions for everyday clinical practice, extreme weather events, like in the year 2019 in our case, still comprise an obstacle toward highly performing forecasts at such fine timescales, highlighting that there is still a way to go to this direction.

**Keywords** Aerobiology · Diurnal pollen distribution · Dynamic regression · Environmental health · Neural networks · Time series analysis

## Abbreviations

| | |
|---|---|
| $\omega$ | Periodic term |
| $B$ | Difference operator |
| $d$ | Non-seasonal difference |
| $p$ | Order of the non-seasonal autoregressive model |
| $q$ | Order of the non-seasonal moving average model |
| $P$ | Order of the seasonal autoregressive model |
| $Q$ | Order of the seasonal moving average model |
| $\varphi$ | Parameter of non-seasonal autoregressive model |
| $\theta$ | Parameter of non-seasonal moving average model |
| $\Phi$ | Parameter of seasonal autoregressive model |
| $\Theta$ | Parameter of seasonal moving average model |
| $\eta$ | Error term following ARIMA process |
| $f$ | Activation function |
| $x$ | Input of a neural network |
| $y$ | Output of a neural network |
| $w$ | Weight of a neural network structure |
| $b$ | Bias |
| MAE | Mean absolute error |
| RMSE | Root mean square error |
| $R^2$ | Coefficient of determination |

# 1 Introduction

Airborne pollen dispersion is part of plant phenology, following yearly seasonal cycles with the aim of successful reproduction. While elementary for the ecosystem, pollen grains are known to be a trigger for allergic reactions in sensitized individuals (Sofiev and Bergmann, 2013). The current prevalence of allergic diseases worldwide remains high, ranging from 15 to 25% (Passali et al. 2018), with industrialized countries affected more by this negative trend (Pawankar 2014). The ongoing increase in air temperature and the overall effect of climate change have been increasing steadily the abundances of airborne pollen across the globe and, at the same time, have been shifting earlier the pollen seasons for several allergenic taxa (Ziska et al. 2019). The World Allergy Organization has warned that, because of climate change, plants will be stressed to flower and pollinate earlier within the year and in higher amounts, thus increasing the natural pollen exposure of sensitized individuals and, consequently, increasing the severity of associated symptoms (Pawankar 2014).

Being mostly not a life-threatening condition, pollen allergic symptoms can significantly reduce health-related quality of life and workplace productivity of people concerned because of profound physical and psychological complications (Blaiss et al. 2018; Devillier et al. 2016; Haanpää et al. 2018). Allergic individuals have several possibilities to control allergic symptoms, with allergen avoidance being one of the most effective measures (Glacy et al. 2013). However, since severity of occurring symptoms significantly depends on the current concentration of aeroallergens in the ambient environment (Bastl et al. 2013), to be effective, allergen avoidance strategies make sense only if performed when concentration of airborne allergenic pollen is high. Consequently, pollen information provided for example via pollen applications to the target population of allergic individuals might become an important aid in avoiding exposure to allergenic pollen, and in planning medication and outdoor activities (Kmenta et al. 2014). As airborne pollen has been identified as a biological weather parameter, a network of nearly 400 Hirst-type pollen traps is currently monitoring the airborne pollen in Europe (Berger et al. 2013). However, for pollen information to be useful for allergy management, it has to be delivered on time,

shortly after the measurement took place, to reflect the actual pollen abundance. Therefore, in order to provide up-to-date information on pollen concentration, a more rapid, and preferably instantaneous technique in pollen monitoring than a conventional pollen trap of Hirst-type is needed. Automated pollen monitoring in real time might be a solution covering this urgent need.

Such novel approaches have been implemented very rarely, as by Chappuis et al. (2020), who used data deriving from an automatic pollen monitoring system. The importance of integrating hourly resolution pollen measurements to forecasting models and, even more, using real-time data from novel, automatic monitoring devices has been suggested and discussed by Sofiev (2019), highlighting that such an approach could boost the predictive power of future models. To be fair, it is also pointed out (and we agree, as our current results also show) that there is still a way to reach operational predictions for everyday practice. Toward the same direction, Geller-Bernstein and Portnoy (2019) reviewed that automatic, real-time pollen monitoring information would be valuable for short-term operational forecasts for allergic individuals, which, otherwise, is currently provided via daily predictive models with no detailed information on the intradiurnal variation for everyday activities and planning.

For this reason, the Bavaria State in south Germany has developed a network based on the automatic pollen monitoring devices of BAA500 type (Bio Aerosol Analyzer 500) (Oteros et al. 2019), as described in more technical detail in (Oteros et al. 2015). The BAA 500 is an automatic system for air particle collection (among others, pollen and fungal spores), analysis, and automatic data transmission to a data bank, with pollen information available three hours after observation. Automatic pollen monitoring is a promising tool in pollen season monitoring, as it provides pollen information nearly up-to-date with a high sampling rate of up to 8 pollen measurements per day. The BAA 500 operated in Munich, Germany, was reported to be a functional pollen monitoring device with 93.3% of pollen automatically classified by that device to be correctly identified (Oteros et al. 2015). Automatic pollen monitoring is a new technique, which is yet not widely used. At the moment, only few countries stand out developing innovative monitoring sites. Among those are Japan (Kawashima et al. 2017),

the USA (Buters et al. 2018), Switzerland (Crouzy et al. 2016), and Germany, the latter of which has been operating automatic pollen monitors for the last half decade.

The circadian pathophysiology of pollen allergy is well documented already (Nakao et al. 2015), with symptoms worsening over night or in the early morning. Because of the lack of real-time, high-resolution (hourly) pollen measurements, this phenomenon remains poorly researched. Most commonly, aerobiologists work on daily data, predicting the pollen concentration for the next day or several days ahead. The novel automatic pollen monitoring devices, with the near-real-time pollen data, allow to go beyond the current state-of-the-art and to develop reliable short-term pollen predictions. Pollen forecasting at this scale can be the cornerstone of operational diurnal allergy risk alerts for allergic individuals.

To achieve such operational forecasts, apart from the real-time, high-resolution pollen data, sophisticated mathematical and statistical tools need to be employed. Scientific works examining the diurnal pollen variation in the air only seldom apply deterministic predictive models, narrowing their efforts down to descriptive methods and correlation analysis (Chappuis et al. 2020; Fernández-Rodríguez et al. 2014; Ščevková et al. 2015). The most common predictive techniques used so far are linear or nonlinear regressions, with significant steps having been made the last few years (Nowosad et al. 2018; Piotrowska, 2012; Ritenberga et al. 2016), and time-series analysis, based on Box–Jenkins methods (García-Mozo et al. 2014; Ocana-Peinado et al. 2008; Valencia et al. 2019). Also, variables like meteorological factors are frequently considered, as they have been proven as significant predictors of airborne pollen concentrations. Meteorological factors, such as solar radiation (Iglesias-Otero et al. 2015; Nowosad et al. 2018), sunshine duration (Myszkowska & Majewska, 2014; Rodríguez-Rajo et al. 2006), and air temperature (Howard & Levetin, 2014; Nowosad et al. 2018; Ščevková et al. 2015), are positively correlated with airborne pollen concentrations, whereas variables like relative humidity (Ščevková et al. 2015; Makra et al. 2011), and precipitation (Piotrowska, 2012; Rodríguez-Rajo et al. 2006) show a negative association with airborne pollen abundances. Some articles examined the relationship with

wind vectors and found them to be of significant influence (Astray et al. 2010).

Nowadays, novel and more sophisticated forecasting techniques are starting to be employed, as in the case of machine learning, which is increasingly gaining scientific interest. Several aerobiological studies have implemented machine learning algorithms, at various scales of analysis, such as artificial neural networks (Iglesias-Otero et al. 2015; Puc, 2012; Valencia et al. 2019), random forests (Nowosad et al. 2018; Zewdie et al. 2019), and support vector machines (Zewdie et al. 2019).

Each pollen forecasting technique exhibits pros and cons, and their selection is based on the research question per case and on data availability and quality. Therefore, regression analysis allows for inclusion of co-factors, but neglects the serial autocorrelation of all variables. On the contrary, Box–Jenkins models consider the autocorrelation of the dependent variable, but neglect the effect of other potential co-factors. Dynamic regression, albeit a statistical approach using the advantages of both above-mentioned forecasting techniques (Pankratz, 2012), has been seldom adopted in airborne pollen forecasting (Ocana-Peinado et al. 2008). Overall, forecasting of pollen concentrations is challenging due to the data complexity, intense seasonality with numerous 'out of season' zero values, high skewness and level of irregularity and extreme outliers. The above are mixed in a double-periodic pattern, within-day and within-year, with different factors influencing each periodicity and pollen distribution. The relationships are often nonlinear and the affecting co-factors usually collinear and sometimes confounding. This challenge could be answered by machine learning algorithms, like artificial neural networks, as they have a high ability to assess complex relationships (Twomey & Smith, 1995). To ensure the sound interpretation of the acquired results produced by the artificial neural network, it then makes sense to cross-validate the model output with that of 'conventional' forecasting techniques, as time series analysis and dynamic regression.

The aim of the present study was to assess and forecast the diurnal variability of airborne pollen concentrations and the development of short-term predictive models using near-real-time 3-hourly *Betula* and Poaceae pollen data. Both pollen taxa were selected because of their high atmospheric abundance in Bavaria (Oteros et al. 2019), and of their high

prevalence in sensitization rates among the study area population (Muzalyova et al. 2019). To our best of knowledge, there is very limited research focusing on forecasting of diurnal pollen concentrations based on data provided by automatic pollen measurement systems. Therefore, this is the first paper using a 3-h sampling frequency of airborne pollen detected by an automatic pollen monitoring to develop and compare different predictive models. Knowledge of variation of pollen quantity on hourly scale is very important for people suffering from pollen allergies, as it can help them to avoid exposure to allergenic pollen. Incorporating real-time, automatic pollen measurements in airborne pollen forecasts is expected to dramatically improve the efficiency of allergy management.

## 2 Materials and methods

### 2.1 Data

Pollen data for *Betula and* Poaceae were acquired by use of an automatic pollen monitoring device BAA500, located in Augsburg, Germany. The automatic pollen monitor was situated at the Bavarian State Office for the Environment (B*ayerisches Landesamt für Umwelt—LFU Bayern*) (coordinates 48°32′ 60.29″N, 10°90′30.77″E), located in a suburban environment in Augsburg, Germany. The pollen data were collected in 3-h intervals for the years 2016–2019. Accordingly, each day (24-h period) encompasses 8 data points beginning with the first pollen measurement performed at midnight (0.p.m), and the last performed at 9 p.m. Pollen concentrations are expressed in grains per $m^3$ with a time step $n$ corresponding to a 3-h interval. Missing data (8.4%

**Table 1** Descriptive statistics of the pollen measurements of examined taxa (grains per $m^3$)

|  | *Betula* | Poaceae |
|---|---|---|
| Mean (SD) | 71.9 (138.1) | 12.6 (29.7) |
| Median | 30.0 | 4.0 |
| Min | 0 | 0 |
| Max | 1582.0 | 750.0 |
| Skewness | 5.9 | 10.0 |
| Kurtosis | 49.1 | 181.9 |

*Betula* and 7.8% Poaceae) were imputed based on regression analysis using 5 data points of the corresponding time period before and after the data gap. Scattered missing points were imputed by averaging closest measurement before and after the data gap. The normal distribution of the data was tested using the Kolmogorov–Smirnov and Shapiro–Wilk tests, where it was concluded that the hourly data did not follow a normal distribution being extremely right-skewed (Table 1).

Meteorological data were retrieved from the German Weather Institute (*Deutscher Wetterdienst—DWD*, https://opendata.dwd.de/climate_environment/CDC/), recorded at the airport of Augsburg (coordinates 48°21′57.564″ N, 10°53′ 34.944″ E), located approximately 11 km north of LFU. The following meteorological parameters were available for data analysis: air temperature [°C], relative humidity [%], air pressure [hPA], precipitation [mm], sunshine duration [min], solar radiation [J/cm$^2$], and wind speed [m/s]. A Spearman correlation test was used to analyze associations between the examined meteorological variables. The statistical analysis included the 3-hourly data set from March to September (main pollen season of *Betula* and Poaceae) and was performed with the SPSS 25.0 statistical package.
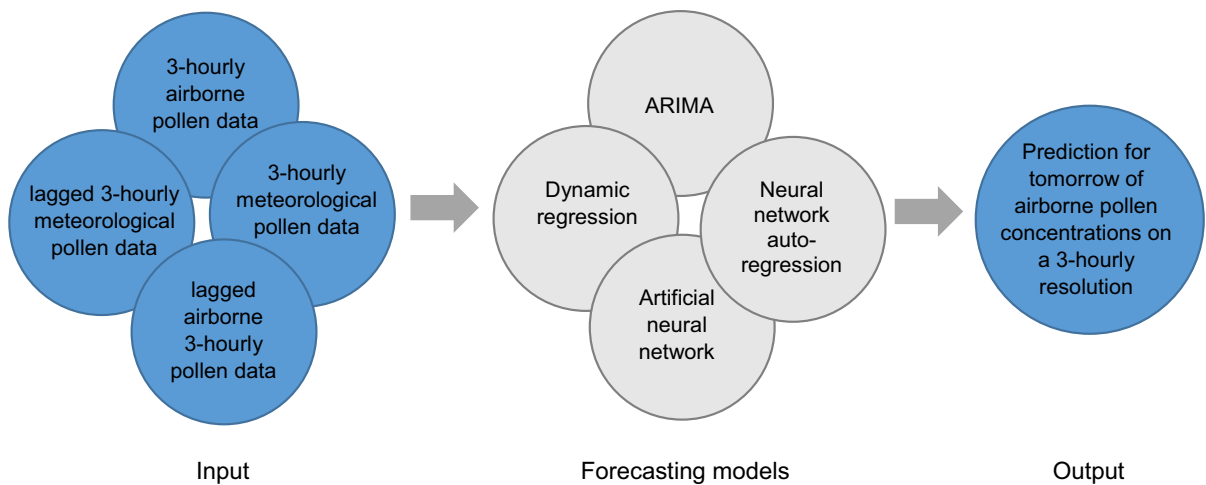
The analysis of the diurnal distribution of pollen concentrations and development of predictive models was performed based on pollen data of the main pollen season for each pollen taxa and each year. Accordingly, the following phenological features were determined for each available study year: Pollen Season Start (PSS), Pollen Season Peak (PSP), Pollen Season End (PSE), Pollen Season Duration (PSD), and the annual Pollen Integral (PI) in line with (Galan et al. 2017). The PSS was defined in line with European Aeroallergen Network pollen season definition. Due to this, the PSS was the first day achieving 5% of the cumulative daily pollen concentrations over the whole year. The PSE was determined as a day reaching 95% of the accumulated daily pollen concentrations throughout the whole pollen season (Bastl et al. 2018). The PI was specified as the sum of the daily average pollen concentrations per cubic meter over the whole year. The PSP was defined as the day with the highest daily pollen concentration. The overview of the data used and the development of the forecasting models is given in Fig. 1.

## 2.2 Autoregressive integrated moving average (ARIMA)

ARMA or ARIMA (also known as Box–Jenkins model) represent a combination of autoregressive and moving average models (Box et al. 2016). For modeling of seasonal time series, ARIMA $(p, d, q)(P, D, Q)_\omega$ is known as multiplicative ARIMA model (Cowpertwait & Metcalfe, 2009). Due to this, six parameters, namely *p, d, q, P, D,* and *Q,* have to be determined to be included in the forecasting model. This step was performed based on the analysis of the Partial Auto-Correlation Function (PACF) and Auto-Correlation Function (ACF). The Akaike Information Criteria (AIC) and the Bayesian Information Criterion (BIC) were the adjustment criteria used for selection of the best model for each examined pollen species. Additional confidence in the best fitting model was gained by deliberately overfitting the model by including further parameters and observing increase in the AIC and BIC. After the best fitting model was found, the correlogram of the residuals was verified as white noise.

## 2.3 Dynamic regression (DR)

A dynamic regression is an extension of a regression model allowing errors from the regression to contain autocorrelations (Pankratz, 2012). A dynamic regression uses advantages of the Box–Jenkins method modeling the autoregression between successive observations of the time series and allows for the inclusion of the external influencing variables like a conventional regression. Additionally, dynamic regression can be applied to seasonal data (Harvey & Scott, 1994) and also allows for lagged effect of the predictors (Pankratz, 2012). In the present study, the order of the autoregressive and moving average components for the dynamic regression modeling was determined based on the evaluation of the PACF and ACF. The external predictors were selected based on backward elimination using Julian day, and 16 lags (two days) of each available meteorological variable. Similar to ARIMA, AIC and BIC were used as adjustment criteria for the best fitting model along with the significance of the selected parameters.
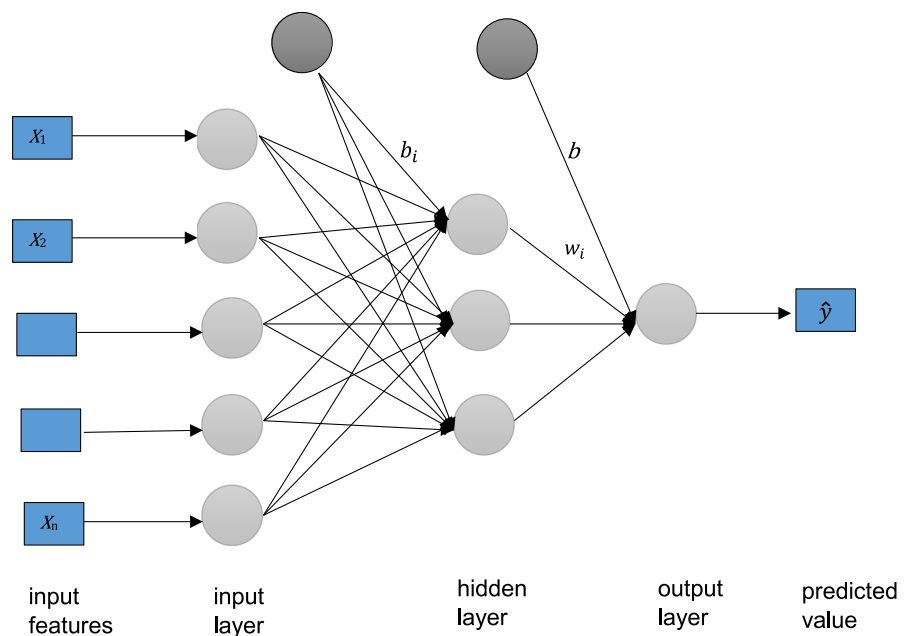
**Fig. 1** Process flowchart of the forecasting model development based on 3-hourly data

### 2.4 Artificial neural network (ANN)

Artificial neural networks are forecasting methods based on a simple mathematical model inspired by information flow in the human brain. A neural network consists of a system of artificial neurons organized in layers. A common neural network incorporates an input, an output layer, as well as one or several intermediate layers containing so-called hidden neurons. A network can incorporate one to many hidden layers, and one to many neurons in each. The number

of input neurons is defined by the number of used input features, and the number of output neurons is defined by the number of required output. The idea of a neural network is to model the response variable, representing the output, based on nonlinear combination of several input variables. A neuron receives information from other neuron or from an external influencing variable and computes a function $f$ based on the weighted sum of the inputs (Goodfellow et al. 2016). The output of a neural network structure having three neurons in the hidden layer shown in Fig. 2 where $x_j$

**Fig. 2** Neural network structure with three hidden neurons

represents the input, $w_{ij}$ is the weight from neuron $j$ to neuron $i$, and $b$ denotes bias. The function $f$ represents an activation function which determines the output activity of the neuron. Through the activation function, the neuron and, thus, the model maps from a linear input to a nonlinear output. Neural network development requires a big implementation of models with different number of neurons in the hidden layer. Designing an optimal schema involves finding the structure with the smallest size network (parsimonious network), which produces optimal errors for trained as well as untrained cases (Astray et al. 2016). During the training phase of the model development, bias values and weights are modified to minimize the error between outputs produced by the model and target values using Mean Squared Error Loss function for linear problems, as given in the preset study.

In the present study, the Julian day of the measurement and available meteorological variables with up to 16 lags of each (up to two-day delay-effect) was used as input variables for the neural networks developed. As the pollen data are usually strongly autocorrelated, the pollen concentrations detected in the previous time periods reflect this time series and were included as influencing variables in order to improve the prediction capacity of the neural network. Since measurements of available input parameters were made on different scales, the parameter were normalized to lie between 0 and 1 before being imputed to the neural network.

## 2.5 Neural network autoregression (NNAR)

Neural network autoregression has a similar theoretical foundation as the ANN explained above. However, this type of an artificial neural network was specifically developed for autoregressive time series and represents a hybrid architecture comprising an ARIMA model and a neural network (Hyndman & Athanasopoulos, 2018). Those combined methods are argued to give better forecasts by taking advantage of each model's capability (Taskaya-Temizel & Casey, 2005). Due to its neural network part of architecture, it is capable of estimating nonlinear relationships, and due to its underlying ARIMA part, the algorithm explicitly uses lagged values of the time series as inputs.

A neural network autoregression is denoted as NNAR $(p, P, k)$ with $p$ indicating the number of lagged inputs, $P$ indicating the number of seasonal lagged inputs, and $k$ representing nodes in the hidden layer. For example, an NNAR $(2, 1, 3)_8$ uses inputs $y_{t-1}$, $y_{t-2}$, and $y_{t-8}$, has three neuron in the hidden layer and is complementary to ARIMA $(2,0,0)(1,0,0)_8$ but without the restriction on the parameters that ensure stationarity.

In the present pollen study, the order of $p$ and $P$ was determined based on the PACF analysis with Julian day and meteorological variables used as external influencing variables similar to the deployment of the ANN. The number of neurons in hidden layer was established similar to the ANN by a trial and error process based on the prediction accuracy of several tested models.

## 2.6 Model validation

It is a common practice in the data modeling to test the predictive power of the established forecasting models on unknown data, not deployed for the fitting process (Goodfellow et al. 2016). For this purpose, the available pollen dataset was split into a training and test datasets as following: the dataset representing the main pollen season in the last year, 2019, was used for the test of the developed predictive models, and the remaining three years of pollen data were applied for the model fitting and training. The predictive accuracy and validity of each established forecasting model was determined based on the comparison of predicted and observed pollen concentration values. Two accuracy metrics, namely mean absolute error ($MAE$) and root mean squared error ($RMSE$), were used as criteria for evaluation of the performance of the established forecasting models:

$$\text{MAE} = \frac{\sum_{i=1}^{N} |\hat{y}_i - y_i|}{N}$$

$$\text{RMSE} = \sqrt{\frac{\sum_{i=1}^{N} (\hat{y}_i - y_i)^2}{N}}$$

Generally, the RMSE stronger punishes deviation between predicted value $\hat{y}_n$ and observed variable $y_n$ due to squaring the difference. It is therefore better suited for modeling on data with strong peak and outliers (Twomey & Smith, 1995).

As both introduced accuracy metrics are based only on error term $e_t$, they are therefore scale-dependent and allow to make comparison between time series that involve different units. In order to compare the performance of predictive models based on pollen data of *Betula* and Poaceae, the coefficient of determination ($R^2$) was used. $R^2$ describes the proportion of variance explained by the model to the total variance in the data and can be defined using the following formula:

$$R^2 = 1 - \frac{\sum_{i=1}^{n}(y_i - \hat{y}_i)^2}{\sum_{i=1}^{n}(y_i - \overline{y})^2}$$

All forecasting techniques were implemented in RStudio, version 1.0.143 using *tseries, fpp2, lmtest, neuralnet* libraries.

## 3 Results

The characteristics of the examined pollen seasons are outlined in Table 2. The main pollen season of *Betula* average started by the end of March (from 15/03 to 8/04), and lasted on average 38 days (SD = 10.2). The main pollen season of Poaceae average started by the end of April (from 20/04 to 12/05) and had a

**Table 2** Descriptive statistics of the examined pollen seasons

|  | 2016 | 2017 | 2018 | 2019 | Average |
|---|---|---|---|---|---|
| *Betula* | | | | | |
| API | 2,491 | 3,828 | 2,121 | 1,860 | 2,732 (5,720) |
| PSS | 05/04 | 30/03 | 04/03 | 01/04 | 27/03 (15/03–8/04) |
| PSP | 12/04 | 01/04 | 16/04 | 09/04 | 09/04 (04/04–14/04) |
| PSE | 09/05 | 13/05 | 24/04 | 25/04 | 27/05 (19/05–04/06) |
| PSD | 34 | 44 | 51 | 24 | 38 (10) |
| Poaceae | | | | | |
| API | 2,018 | 1,490 | 856 | 378 | 1,197 (5,199) |
| PSS | 07/05 | 11/05 | 21/04 | 27/04 | 01/05 (20/04–12/05) |
| PSP | 07/06 | 09/06 | 01/06 | 08/06 | 06/06 (03/06–09/06) |
| PSE | 02/08 | 28/07 | 18/07 | 13/08 | 31/07 (22/07–09/08) |
| PSD | 87 | 78 | 107 | 108 | 95 (13) |

In parenthesis, the range of calendar dates per season feature is provided, along with the standard deviation for API

API; annual pollen integral, PSS: pollen season start date, PSP: pollen season peak date, PSE: pollen season end date, PSD: pollen season duration (in number of calendar days)

comparably longer duration of 95 days (SD = 12.9). Considering *Betula*, the PD of the pollen season usually occurred shortly after the PSS on the 13th day of the main pollen season (from 04/04 to 14/04), whereas the PD of Poaceae was situated closer to the middle of the main pollen season and occurred on the 40th day (from 03/06 to 09/06) of the main pollen season. Furthermore, *Betula* usually had one well-defined peak, whereas Poaceae was characterized by several peaks of variable amplitude within the main pollen season. Generally, the pollen release of *Betula* was more intensive in absolute terms, peak values and also average pollen concentration per time period, compared to that of Poaceae.

Regarding inter-annual variability, the pollen season of the year 2018, interestingly, stands out among analyzed pollen seasons due to the earlier PSS and PSE for both investigated allergenic species (Table 2). In particular, the main pollen season of *Betula* started already by the beginning of April and lasted for more than fifty days. The PSS of Poaceae occurred 10 days earlier of the average date and ended by the middle of July. Furthermore, the intensity of the Poaceae pollen seasons was continuously decreasing across examined years, with 2019 exhibiting the lowest pollen abundance of all years (Fig. 3).

The diurnal distribution of pollen concentrations of both taxa is depicted in Fig. 4. The pollen load of *Betula* was relatively constant during the day with the highest levels occurring at 3 p.m. Kruskal–Wallis-Test revealed a significant difference between time periods ($H$ (7) = 28.590, $p < 0.01$). However, due to high standard deviation and none well-defined diurnal patterns a post hoc test (Dunn–Bonferroni) revealed only difference between 12 a.m. and 3 p.m. to be significant ($z = -3.137$, $p = 0.048$), whereas all other differences of pollen concentration between considered time periods were non-significant. On the contrary, the pollen concentration of Poaceae was noticeably peaking twice a day at 9 a.m. and 3 p.m., with relatively low abundance during the night hours. The Kruskal–Wallis Test revealed a significant difference between groups [$H$ (7) = 317.982, $p < 0.01$], and a pairwise comparison showed significant differences between pollen concentrations measured between the night hours and early morning (9 p.m.–6 a.m.) and those observed beginning with morning until evening (9 a.m.–6 p.m.). As pollen concentration of Poaceae is higher during the warmer parts of the day, it
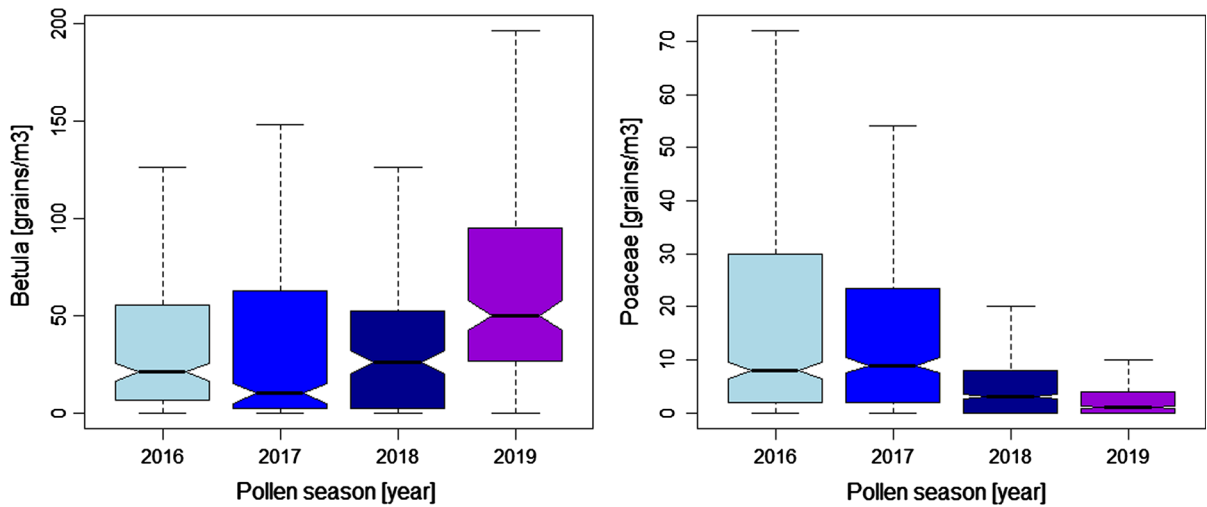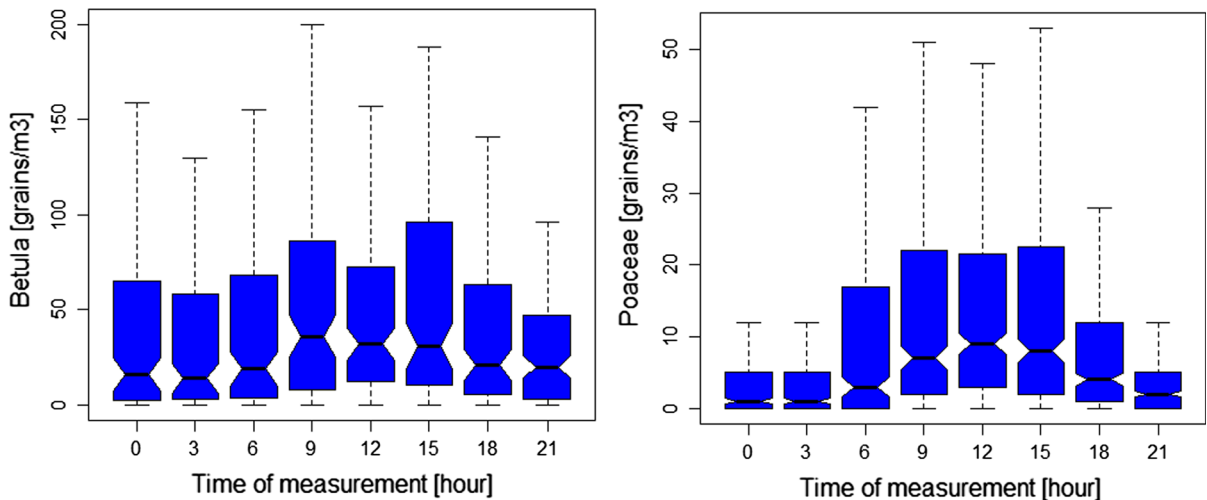
**Fig. 3** Boxplots of averaged annual concentrations of *Betula* and Poaceae pollen, where the horizontal line denotes the median of all concentrations throughout the year, while the box and vertical lines signify the quartiles 25–75%
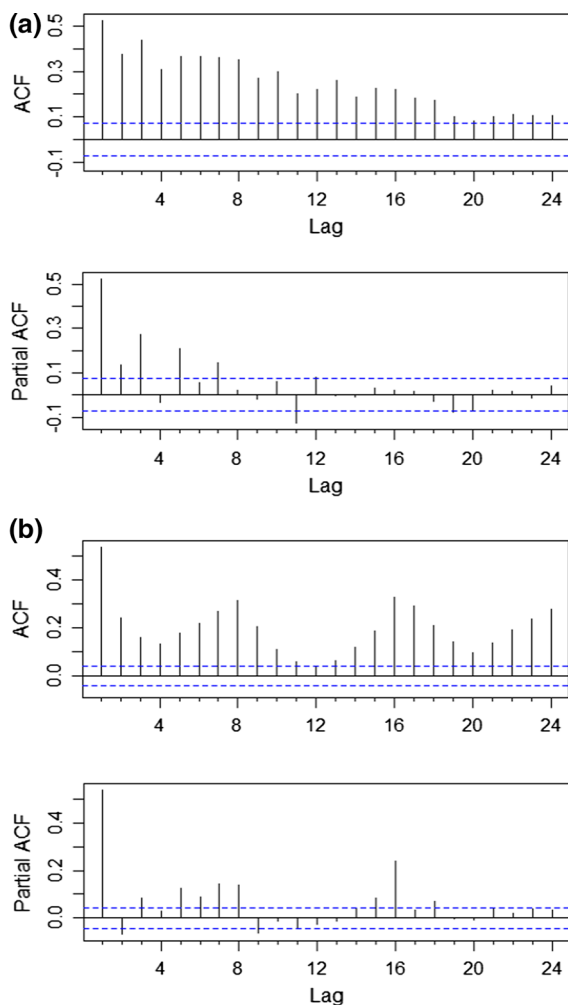


**Fig. 4** 3-hourly average distribution of *Betula* and Poaceae pollen concentrations over 24 h

suggests a stronger relationship between airborne pollen concentrations of this allergenic species and heat-related meteorological variables like air temperature, sunshine duration, and solar radiation.

Correlogram analysis (Fig. 5) of *Betula* pollen concentrations showed a steady decrease with no well-defined peaks at daily cycle (at 8th lags), concluding that a seasonal term has to be included in the model. Inspection of the PACF-correlogram suggested that choosing one non-seasonal, and none seasonal autoregressive and moving average parameters were sufficient. However, due to rising significant correlation

between 5 and 8th lags in the PACF, up to seven non-seasonal autoregressive and moving average parameters were tested. Correlogram of the Poaceae pollen data depicts a tendency similar to *Betula*'s, including a well-defined seasonality of the data at 8th lags; however, the decrease across the lags occurs considerably slower in comparison with *Betula,* presumably due to the shorter main pollen season duration. Inspection of the partial autocorrelation function showed a high significant correlation at lag one. According to this analysis one non-seasonal, as well as, up to two seasonal autoregressive and moving

**Fig. 5** ACF (autocorrelation function) and PACF (autocorrelation function) for *Betula* (**a**) and Poaceae (**b**) pollen data

average parameters might be sufficient for the ARIMA model. However, in order to estimate the effect of overfitting up to seven, both, autoregressive and moving average parameters, and one seasonal autoregressive and moving average parameters were tested. The best fitting model for each pollen species was chosen based on the lowest AIC and BIC statistics.

It is not possible to mention all relevant results of tested ARIMA model structures in this paper; therefore, the structures related singly to the best-fitted models are presented in Table 4. Thus, the best ARIMA model of *Betula* pollen concentration corresponded to ARIMA $(7,1,3)(1,1,1)_{[8]}$ and contained seven non-seasonal autoregressive and three moving average parameter, and one of each seasonal
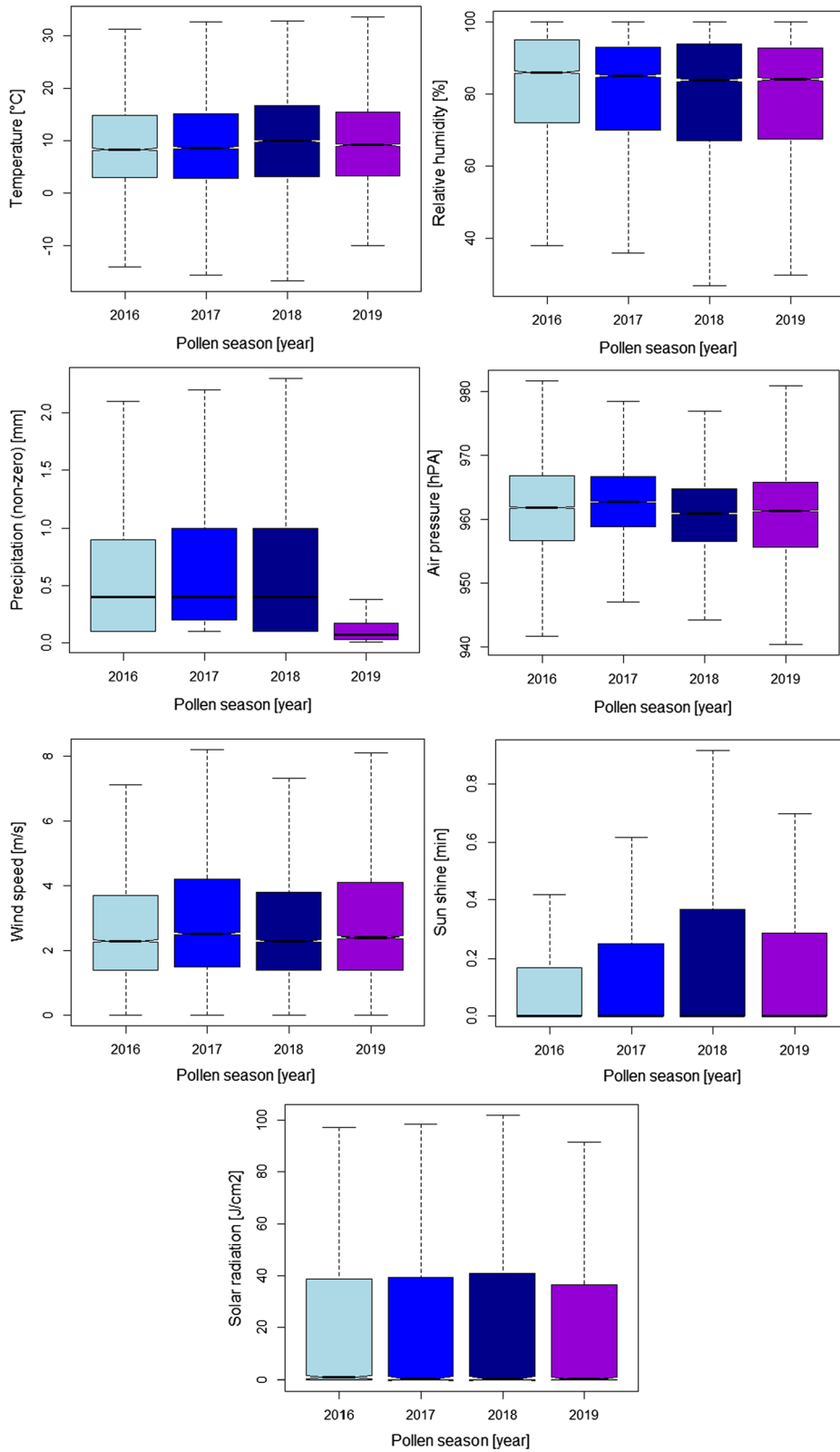
**Fig. 6** Boxplots of averaged annual values of different ▶ meteorological factors, where the horizontal line denotes the median of all measurements throughout the year, while the box and vertical lines signify the quartiles 25–75%

autoregressive and moving average parameter. Regarding Poaceae, the best fitting model was given by ARIMA $(1,1,2)(1,0,1)_{[8]}$ and consisted of one non-seasonal autoregressive parameters, two non-seasonal moving average parameter, and one of each seasonal autoregressive and moving average parameter.

Descriptive analysis of the meteorological variables (Fig. 6) reviled a significant difference between the years used in the training data set and the year 2019 representing the test data set. Of particular note, the year 2019 was significantly drier in comparison with all other considered pollen seasons.

The Spearman's correlation analysis was used preliminary to DR development in order to discover relationships between pollen concentrations and available meteorological parameters. The results of correlation analysis are given in Table 3. Generally, correlations were significant in a large number of cases. As expected, Poaceae pollen concentrations were strongly related to air temperature, sunshine, and solar radiation in comparison with *Betula* pollen counts. The air pressure was found to be significantly correlated only to *Betula*, whereas precipitation and humidity were negatively related to the pollen concentrations of both pollen taxa. No significant relationship between wind speed and pollen concentrations was detected.

The order of autoregressive and moving average parameters in DR was determined based on ACF and PCF analysis, however, with regard to previous ARIMA modeling. All available weather data were imputed in the dynamic regression using up to 16 lagged values representing two days as well as Julian day of the measurement and tested for significance. A step-by-step procedure was followed, and a backward stepwise removal of all non-significant influencing variables was processed, beginning with the highest p-value. The best fitting model was determined using AIC and BIC values along with significance of the certain influencing variables. The final result depicting the best fitting model for both pollen taxa can be taken out of Table 4.

Similar to correlation analysis, air temperature had a positive significant effect on the airborne pollen concentrations for both examined pollen species, with regression coefficient being higher for *Betula*.

However, multiple lagged time periods of temperature measurement were found significant for Poaceae, suggesting airborne pollen concentration of this species to be more sensitive to air temperature.

**Table 3** Correlation coefficients of pollen concentrations and meteorological variables (only training dataset)

| T (°C) | RH (%) | P (mm) | AP(hPa) | WS (m/s) | S (min) | R(J/cm$^2$) |
|---|---|---|---|---|---|---|
| *Betula* | | | | | | |
| 0.248** | − 0.238*** | − 0.230*** | 0.163** | 0.018 | 0.207** | 0.208** |
| Poaceae | | | | | | |
| 0.471** | − 0.369** | − 0.151** | 0.002 | 0.079 | 0.370** | 0.387** |

T: temperature, RH: relative humidity, P: precipitation, AP: air pressure, WS: wind speed, S: sun shine, R: diffuse radiation

Significance levels: 0.001 ***, 0.01 **, 0.05 *

**Table 4** Coefficients obtained for ARIMA and DR models

| | ARIMA | | DR | |
|---|---|---|---|---|
| | *Betula* | Poaceae | *Betula* | Poaceae |
| $\varphi_1$ | − 0.59*** | 0.24*** | 0.59*** | 0.18*** |
| $\varphi_2$ | − 0.38*** | | − 0.41*** | |
| $\varphi_3$ | 0.53*** | | 0.51*** | |
| $\varphi_4$ | 0.10 | | 0.07 | |
| $\varphi_5$ | 0.26*** | | 0.29*** | |
| $\varphi_6$ | 0.09* | | 0.16*** | |
| $\varphi_7$ | 0.25*** | | 0.25*** | |
| $\Phi_1$ | 0.17*** | 0.92*** | 0.17*** | 0.88*** |
| $\theta_1$ | 0.03 | − 0.72*** | − 0.07*** | − 0.69*** |
| $\theta_2$ | − 0.29*** | − 0.26*** | − 0.23 | − 0.28*** |
| $\theta_3$ | − 0.74*** | | − 0.84 | |
| $\Theta_1$ | − 1.00*** | − 0.80*** | − 1.00** | − 0.75*** |
| $T$ | | | 5.47** | 2.31*** |
| $T_1$ | | | | − 0.97** |
| $T_5$ | | | | 0.91*** |
| $T_8$ | | | | − 0.74** |
| $P$ | | | − 71.03*** | |
| $P_1$ | | | − 39.20** | − 2.77**** |
| $P_2$ | | | 55.74*** | − 2.15*** |
| $AP$ | | | 12.85** | |
| $AP_1$ | | | − 12.31** | |
| *Ljung–Box test* | | | | |
| $Q^*(df)$ | 5.39(4) | 48.21(11) | 15.71(3) | 59.19(15) |
| *p-value* | 0.24 | 0.00 | 0.00 | 0.00 |
| *Goodness-of-fit* | | | | |
| $R^2$ | 0.42 | 0.38 | 0.46 | 0.42 |

Significance levels: 0.001 ***, 0.01 **, 0.05 *

Precipitation had a substantially greater impact on pollen abundance for both examined species reflected in higher calculated parameters with this effect lasting up to 6 h. Interestingly, the effect of rain occurred immediately on *Betula* airborne pollen concentration and with a delay of three hours on Poaceae. The air pressure was a significant predictor but only for airborne *Betula* pollen. Furthermore, only examined meteorological variables representing at most 8th lag were determined as significant predictors of the airborne pollen concentration, suggesting that only most current values have an influence on the pollen levels in the air.

Generally, extension of ARIMA model by meteorological variables has improved the performance of the predictive models in terms of higher coefficient of determination ($R^2$); however, this effect was small despite highly significant relationships.

Statistical results obtained from the ARIMA and DR analysis were used as a starting point for the setup of the both neural networks. Particularly, the best fitting order of the autoregressive parameters served as a starting framework for definition of the NNAR structure. Accordingly, for *Betula* NNAR structure was defined as NNAR $(7,1,k)_{[8]}$, and NNAR $(1,1,k)_{[8]}$ for Poaceae. Lagged pollen counts corresponding to the $p$ and $q$ order of the ARIMA model were also employed as input features in the ANN. All available meteorological variables including its lagged values were deployed as influencing variables in NNAR and ANN, as well as Julian day of the measurement. The number of neurons in the hidden layer $k$ for each neural network was determined iteratively by testing different neuron schemas. After the trial-and-error process, structures providing better results in terms of the model accuracy were obtained for each of examined allergenic species, and each neural network used for data modeling. The final neuron structures, as well as, goodness-of-fit criteria can be taken out from Table 5. Interestingly, the best NAAR structure for predicting *Betula* airborne pollen counts was given by one autoregressive non-seasonal component in comparison with ARIMA and DR having the order of 7. The most important meteorological variables for NNAR and ANN were Julian day, air temperature, precipitation, and solar radiation, whereas the NNAR prediction of Poaceae pollen levels was dominated singly by precipitation. It is also remarkable that neural networks predicting *Betula* pollen counts achieved

substantially higher $R^2$ coefficients in training process in comparison with Poaceae.

The predictive models fitted to the training data set were applied on the test data set in order to determine their predictive accuracy. Overall, the ARIMA and DR could achieve higher coefficients of determination in the test run in comparison with the training of the models. Furthermore, ARIMA and DR performed almost equally well; thus, the deployment of additional meteorological parameters has not changed the predictive accuracy significantly.

On the contrary, the high coefficient of determination achieved when fitting neural networks using training data were only partly reproduced in the independent test. The goodness-of-fit of the independent model test can be seen in Table 6. The NNAR produced better predictions for *Betula*, whereas simple seasonal ARIMA outperformed all other predictive methods in forecasting airborne Poaceae pollen concentrations. Furthermore, DR exhibited low predictive power for Poaceae pollen levels. Despite substantially higher values of *RMSE* and *MAE*, forecasting models predicting *Betula* pollen concentrations performed better, achieving $R^2$ in the range between 0.13 and 0.62. On the contrary, predictive models of Poaceae achieved coefficients of determination between 0.03 and 0.55. The high *RMSE* and *MAE* for *Betula* pollen concentrations were predetermined by higher intensity of airborne pollen levels in comparison with Poaceae.

Figure 7 shows the comparison of predicted values based on four applied modeling techniques and observed *Betula* pollen concentrations. The test prediction was made using roughly 25% of the available data and, in total, consisted of 200 data points. The figure depicts predictions provided by each of the tested forecasting models. The black line shows the observed pollen counts for considered data points, and the other lines depict prediction made by ARIMA, DR, NNAR, and ANN. As can be seen, for the hold-out year 2019, *Betula* had no single well-defined peak in this test data set, and the highest pollen level was achieved closer to the middle of the pollen season. The salient finding was that in terms of pollen season occurrence, ARIMA, DR, and NNAR performed remarkably well; however, they consistently underestimated the pollen abundances. Practically, all models' overall performances were lower than expected, because of not ever managing to predict the highest

**Table 5** Neural networks schemas of fitted models

| Model | $R^2$ |
|---|---|
| *Betula* | |
| NNAR (1,1,8) | 0.74 |
| ANN (13,6,1) | 0.91 |
| Poaceae | |
| NNAR(1,1,4) | 0.56 |
| ANN (12,8,1) | 0.63 |

**Fig. 7** Results of independent test for *Betula* pollen data (3-▶ hourly sequence) using 4 forecasting techniques

peak of the season. Noticeably, the test year, 2019, was the driest one of all examined years (Fig. 6) and, potentially for this reason with the highest annual pollen integral of *Betula* pollen compared to all years in the study period (Fig. 3). The ANN exhibited the highest irregularities, by underestimating the first peak in the overall cluster and overestimating the second peak, whereas it performed quite well in the lower concentrations. Furthermore, it is noticeable that ANN tended either to strongly overestimate, or miss several peaks inside the season, whereas the NNAR generally captured this pollen behavior but underestimated it. Additionally, the DR also showed a tendency to slightly overestimate the airborne pollen concentrations.

The test of the established predictive models for Poaceae pollen concentrations was performed using the main pollen season 2019 representing roughly one fourth of the data and contained in total 872 data points. Figure 8 depicts a representative section of the observed pollen counts beginning with the start of the considered pollen season and predicted values using four forecasting techniques. As shown in the graphical representation, the observed peaking behavior of pollen counts was overestimated by all applied forecasting techniques except for ARIMA.

Additionally, DR showed a tendency to strongly overestimate the variability of the airborne pollen concentrations, whereas both neural networks predict values clearly above the actually observed pollen concentrations, however, capturing the pollen behavior in terms of its amplitude. This result can be traced back to the lowest intensity of the Poaceae pollen season among all examined pollen seasons. ARIMA describes well the pollen behavior of low pollen concentrations of low pollen levels, and the ANN outperformed in forecasting the peaking behavior beginning with the time period 161.
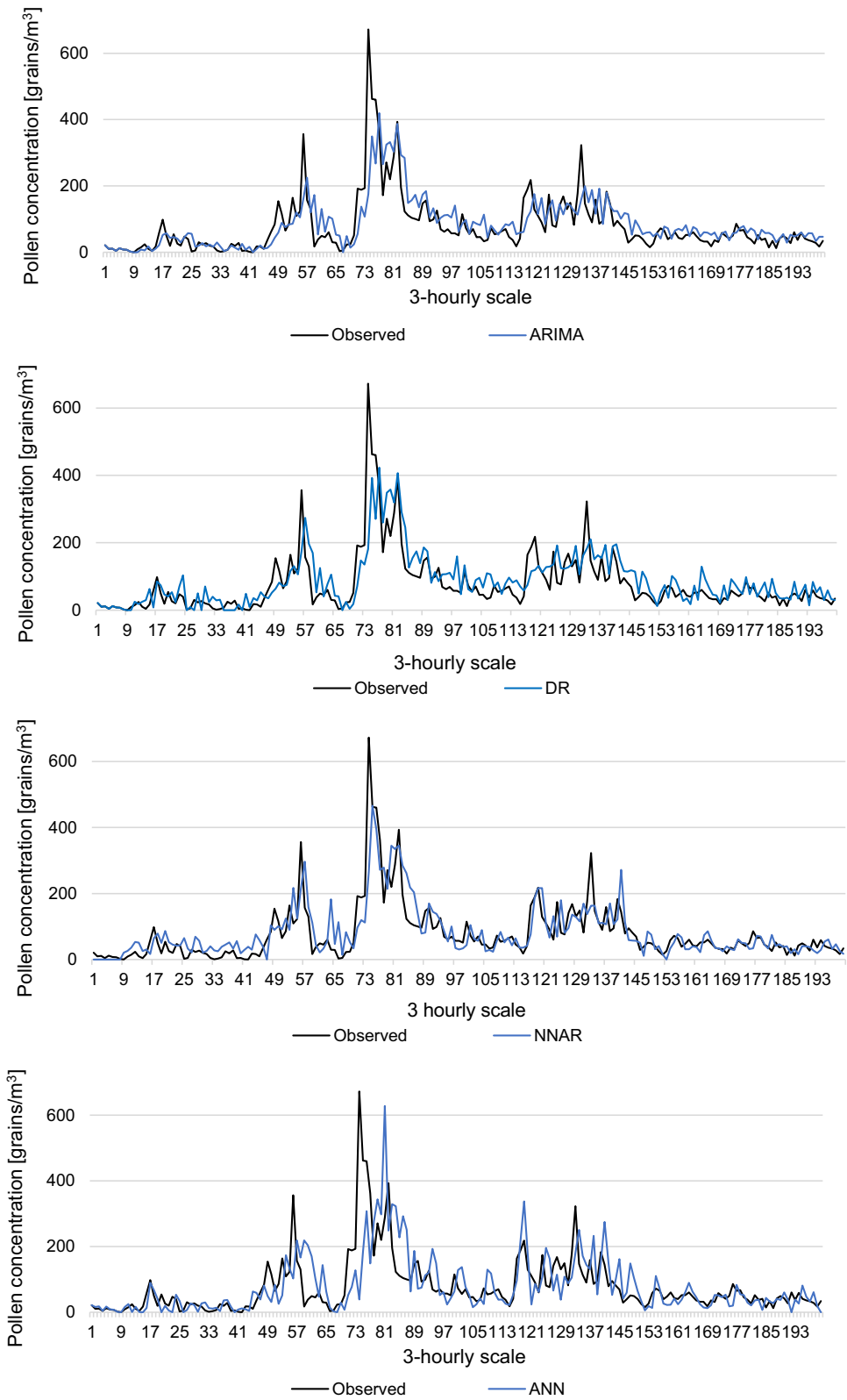
## 4 Discussion

In the present study, we elaborated novel, automated, near-real-time pollen data, on a 3-hourly time resolution, attempting to predict pollen concentrations on a diurnal horizon. In contrast to the current tendency to forecast the start, peak and end of the main pollen season, here we attempted to define the pollen concentration *after* the start of and *within* the main pollen season, so as to potentially provide real-time, operational allergy risk alerts. To achieve this, we used a variety of statistical techniques, among which time series analysis and machine learning. Most forecasting attempts have been using much simpler tools or less fine time resolution. From an operational and clinical point of view, allergic patients and their practitioners actually need the diurnal distribution of air pollen

**Table 6** Prediction capacity of four models with a one-day forecast horizon on a 3-hourly resolution

| *Betula* | | | | Poaceae | | | |
|---|---|---|---|---|---|---|---|
| ARIMA | DR | NNAR | ANN | ARIMA | DR | NNAR | ANN |
| $R^2$ | | | | | | | |
| 0.56 | 0.56 | 0.62 | 0.13 | 0.55 | 0.03 | 0.29 | 0.45 |
| RMSE | | | | | | | |
| 59.15 | 59.41 | 55.76 | 84.94 | 3.81 | 6.86 | 4.85 | 4.26 |
| MAE | | | | | | | |
| 34.22 | 37.03 | 35.50 | 49.22 | 2.23 | 5.33 | 3.70 | 3.21 |

$R^2$: the coefficient of determination, RMSE: Root Mean Square Error, MAE: Mean Absolute Error

abundances every day so as to plan their daily activities, including exposing (or not) themselves to expected airborne pollen concentrations and receiving the appropriate medication.

Considering that sensitization rates to airborne pollen account up to 25% worldwide (Passali et al. 2018), and pollen allergies comprise according to the World Allergy Organization one of the emerging diseases of the century (Pawankar 2014), the above-mentioned information will undoubtedly be the cornerstone for the pollen allergen avoidance on a regular basis, if disseminated operationally. Specifically in the study area's country, Germany, almost 15% of adult population are suffering from at least one allergic disorder (Bergmann et al. 2016), and allergic individuals account for about 12.6% among German children (Schmitz et al. 2014). This additionally highlights the necessity for the elaboration of such prophylaxis and management toolkits.

This study employed advanced statistical methods, namely ARIMA, dynamic regression, and machine learning, such as neural autoregression and artificial neural network, to predict pollen concentrations of *Betula* and Poaceae in Augsburg, Germany. The mathematical modeling techniques were used by integrating meteorological factors and past pollen observations. Such approaches are quite common in this research area, but not on such a fine, 3-hourly, timescale.
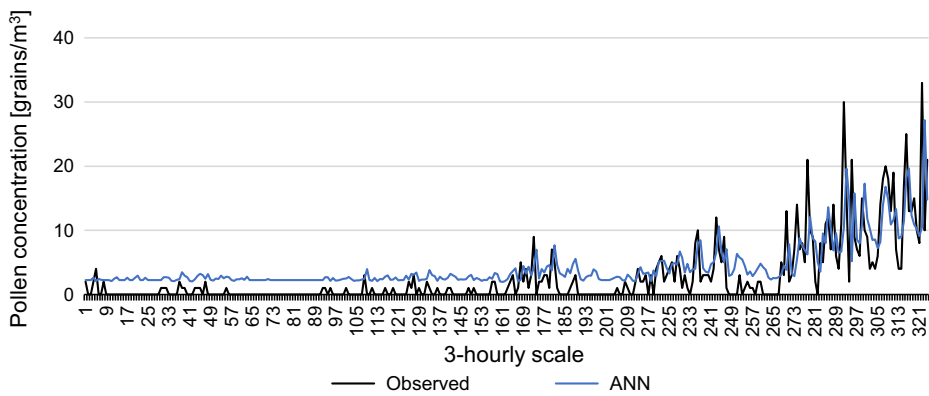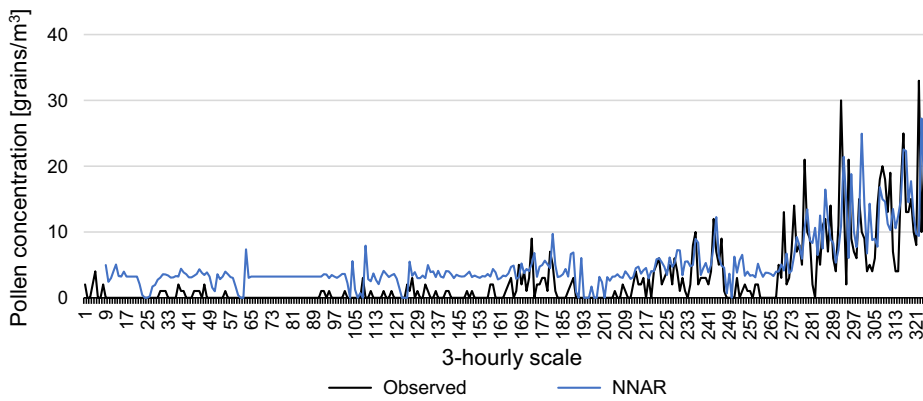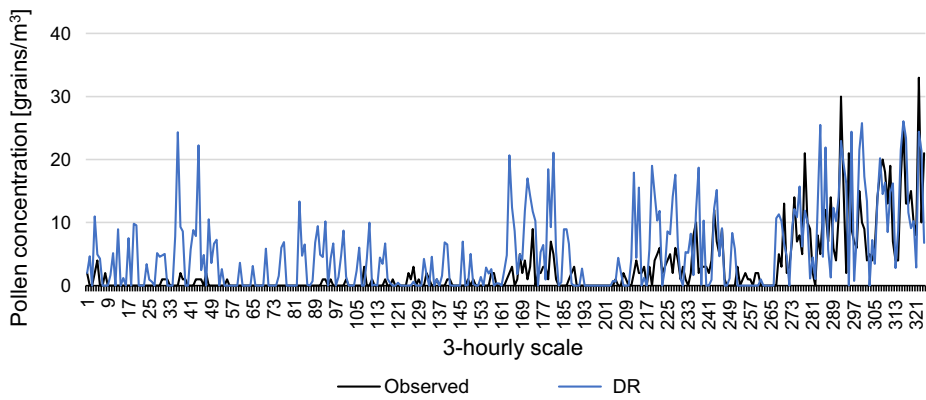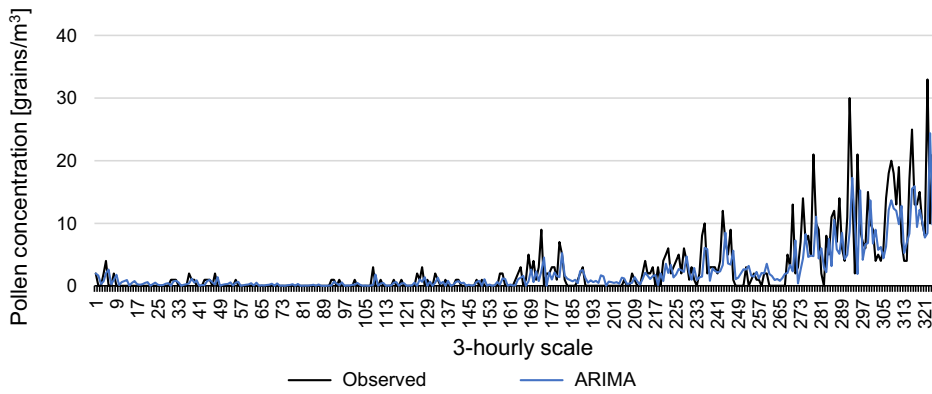
Among the statistical forecasting techniques employed in the present study, dynamic regression considered autocorrelations both with the dependent pollen data and the values of meteorological variables and performed better in pollen prediction based on the training dataset. This finding agrees with Sanchez et al. (2005), who showed the combination of meteorological factors and previous pollen data to yield better results in pollen forecasting, than using alone pollen data or meteorological variables (Sánchez Mesa et al. 2005). In the present study, two meteorological variables, namely air temperature and precipitation, were determined as significant predictors in DR modeling for both examined pollen species, with precipitation having a stronger effect on the airborne pollen concentration. In the current aerobiological research, the most of pollen forecasting studies apply some meteorological variables as input parameters. Among them, air temperature is one of the most studied meteorological factors which is discovered to

have a significant effect on airborne pollen concentrations across different pollen taxa (García-Mozo et al. 2014; Ziello et al. 2012). For example, Howard and Levetin (2014) used air temperature and precipitation to predict pollen concentrations too (Howard & Levetin, 2014), discovering that temperature was one of the most significant repressor. Iglesias-Otero et al. (2015) employed precipitation, sunshine duration, and humidity, with rainfall being the most sensitive variable in the predictive model (Iglesias-Otero et al. 2015). These findings agree with our result, showing the precipitation to have an even stronger impact on airborne pollen levels. That suggests that rainfall simply washes out the pollen grains from the air with this effect lasting for several hours. It is worth pointing out here that different meteorological and climatic indices would provide variable predictive capacity to our developed models, and this is highly relying on the timescale examined. It is well known, as documented, i.e., by Damialis et al. (2005) that even by conventional Hirst-type monitoring techniques, wind vectors and precipitation are the leading determining factors, with the effects lasting for at least four hours (Damialis et al. 2005). In the present study, although an extension of the simple ARIMA model by meteorological variables provided better results in the training, the predictive power of ARIMA and DR models were similar for *Betula* pollen counts, whereas DR failed in predicting airborne pollen concentration of Poaceae, achieving coefficient of determination of only $R^2 = 0.03$. Furthermore, considering Poaceae pollen concentrations, the simplest among applied forecasting techniques, namely ARIMA, clearly outperformed all other models.

Regarding the two machine learning techniques used for the development of predictive models, neural autoregression substantially outperformed the artificial neural network, for *Betula* pollen data, and delivered the best predictive power in terms of $R^2$ of all applied forecasting techniques. On the contrary, Poaceae was better predicted by ANN. In general, the forecasting model developed for *Betula* pollen performed better in terms of obtained coefficient of determination $R^2$ in independent test, despite higher

variation in the data. Possibly, it can be explained by a shorter main pollen season for *Betula* and a clearer pattern of pollen behavior consisting of only one well-defined peak. Furthermore, it is worth mentioning, that intensity of the Poaceae pollen season was decreasing across examined years with 2019 having the lowest pollen abundance. Also, 2019 was the driest year, especially in the pollen season of Poaceae, both in terms of precipitation and relative humidity, which contributed to it being also the longest pollen season. Consequently, both applied machine learning techniques were constantly overestimating observed airborne pollen counts, especially in the weakly abundant beginning of the pollen season. An additional inclusion of a parameter reflecting expected intensity of the pollen season might have an essential effect on the accuracy of the predictive models. Consequently, more historical pollen data are obviously needed to more thoroughly investigate the intensity of the Poaceae pollen seasons.

Overall, a good predictive performance of simple seasonal ARIMA model in comparison with advanced forecasting techniques suggests that the phenology of the plants, reflected by the lagged pollen concentrations, is the most relevant predictor for the observed airborne pollen concentration. This insight is also supported by diurnal pollen concentration patterns discovered in the present study, especially for Poaceae, showing significant differences in airborne pollen concentrations depending on the time period of the measurement. This finding highlights the importance of the further, scrupulous investigation of the diurnal variation of the airborne pollen concentration and its influencing factors. Investigation of airborne pollen concentrations on hourly scales represents a promising research direction, since it accommodates one of the most urgent/important objectives, namely that of the delivery of pollen information to people suffering from pollen-induced allergies. Given that under ongoing climate change conditions, increasing and more intense extreme weather events influence the abundance and seasonality and circadian periodicity of airborne pollen, developing accurate short-term forecasts is a real challenge. In our results, this is highlighted by the fact that the significantly drier year 2019 led to reduced predictive capacity of most models and signified past pollen records as the most reliable, in this dataset, predictor of future, on an hourly scale, pollen concentrations. It is anticipated

that unexpected and extreme weather incidents may be already causing unpredictable pollen seasons and diurnal distribution, which is worth to be investigated more thoroughly.

When developing a forecasting model to notify the pollen allergic individuals about expected airborne pollen levels for supporting their pollen allergy management, one has to keep in mind the needs of the target population. Allergic individuals might be hardly interested in pollen forecasting expressed in absolute values. On the contrary, they might be interested in notifications of critical pollen values, or expected symptom severity induced by the airborne pollen levels. Firstly, this consideration can also affect the definition of the main pollen, taking into account pollen thresholds inducing allergic reaction of different severity in sensitized individuals (Karatzas et al. 2019). Secondly, there are several studies focusing on prediction of certain levels of airborne pollen concentration (Brighetti et al. 2014; Castellano-Méndez et al. 2005) or even expected season severity (Sánchez Mesa et al. 2005). Pollen level inducing an allergic reaction of a certain severity in allergic individuals might be variable for different locations due to different climatic conditions (Weger et al. 2013). The forecasting of critical values might be very useful for allergic individuals. However, so far, it lacks scientific efforts in this direction in Germany, and, thus, it lacks knowledge of pollen thresholds triggering allergic symptoms in sensitive individuals.

Overall, after comparing the performance of the four models used here with the actual observed pollen concentrations, it is concluded that all models had a satisfactory capacity to predict the timing of pollen occurrence on a 3-hourly scale, but most of them were not as well-performing when they had to forecast highly peaked pollen concentrations. Given the avowedly short length of the overall examined time series, as well as the weather peculiarities in the hold-out year 2019, we consider all models exhibiting adequate forecasting performance. In similar statistical approaches on a 2-hourly timescale, only few researchers, like Chappuis et al. (2020), reported some significant correlations between hourly pollen concentrations and hourly weather variables, most frequently weaker than the ones presented in our work. Noticeably, Chappuis et al. (2020) also used data deriving from an automatic pollen monitoring system, even though from a different manufacturer. Some

other studies that attempted to predict pollen concentrations with weather variables on a diurnal level were those by Simoleit et al. (2016), also in Germany (Berlin), who found by rule weaker relationships with most meteorological factors. Likewise, Ríos et al. (2016) in Mexico, and Alba et al. (2000) in Spain also detected much weaker predictions, almost by rule.

An important limiting factor of the present study is the volume of the available pollen data. As BAA500 has been operating in Augsburg for only half a decade, only four complete pollen seasons were available for the data analysis. Environmental data are known to be very complex to model due to underlying inter-relations (Zewdie et al. 2019); hence, the time-series available for the present research might be too short to determine the seasonal phenology of examined species or to identify and characterize anomalous pollen seasons. In order to realize and to calibrate forecasting models, long historical series of pollen and meteorological data are necessary. Furthermore, it is worth pointing out that the predictive models presented in this study are based on data provided by an innovative fully automated pollen monitor, which, being a novel device, is still undergoing improvements. Although the pollen monitoring has been reported to show a high accuracy of pollen determination (Oteros et al. 2015), it has been documented already that a further improvement of the recognition algorithm is possible and that, consequently, there is still a lot of room for increasing the accuracy of pollen identification in near future (Schiele et al. 2019). Therefore, we conclude that the key for reliable, short-term pollen predictions, does not necessarily lie on the complexity and how sophisticated the applied statistical techniques are, but on the completeness of the toolkit used toward this purpose, as suggested below:

- good quality of data (reliability)
- long datasets (consistency)
- considerations of the whole multi-factorial design
  - pollen autocorrelations
  - interaction effects with weather and climatic parameters
  - trends and multi-periodicities (within season and within the day).

## References

Alba, F., Díaz De La Guardia, C., & Comtois, P. (2000). The effect of meteorological parameters on diurnal patterns of airborne olive pollen concentration. *Grana, 39,* 200–208. https://doi.org/10.1080/00173130051084340

Astray, G., Fernández-González, M., Rodríguez-Rajo, F. J., López, D., & Mejuto, J. C. (2016). Airborne castanea pollen forecasting model for ecological and allergological implementation. *The Science of the Total Environment.* https://doi.org/10.1016/j.scitotenv.2016.01.035

Astray, G., Rodríguez-Rajo, F. J., Ferreiro-Lage, J. A., Fernández-González, M., Jato, V., & Mejuto, J. C. (2010). The use of artificial neural networks to forecast biological atmospheric allergens or pathogens only as Alternaria spores. *Journal of Environmental Monitoring: JEM.* https://doi.org/10.1039/c0em00248h

Bastl, K., Kmenta, M., & Berger, U. E. (2018). Defining pollen seasons: Background and recommendations. *Current Allergy and Asthma Reports.* https://doi.org/10.1007/s11882-018-0829-z

Bastl, K., Kmenta, M., Jäger, S., Bergmann, K.-C., & Berger, U. (2013). Calculation and application of the symptom load index: Computing the season severity from the allergy sufferer's point of view. *Allergo Journal.* https://doi.org/10.1007/s15007-013-0389-4

Berger, U., Karatzas, K., Jaeger, S., Voukantsis, D., Sofiev, M., Brandt, O., et al. (2013). Personalized pollen-related symptom-forecast information services for allergic rhinitis patients in Europe. *Allergy.* https://doi.org/10.1111/all.12181

Bergmann, K.-C., Heinrich, J., & Niemann, H. (2016). Aktueller Stand zur Verbreitung von Allergien in Deutschland. *Allergo Journal.* https://doi.org/10.1007/s15007-016-1015-z

Blaiss, M. S., Hammerby, E., Robinson, S., Kennedy-Martin, T., & Buchs, S. (2018). The burden of allergic rhinitis and allergic rhinoconjunctivitis on adolescents: A literature review. *Annals of Allergy, Asthma & Immunology: Official Publication of the American College of Allergy, Asthma, & Immunology.* https://doi.org/10.1016/j.anai.2018.03.028

Box, G. E. P., Jenkins, G. M., Reinsel, G. C., & Ljung, G. M. (2016). *Time series analysis: Forecasting and control (Wiley Series in Probability and Statistics).* Wiley.

Brighetti, M. A., Costa, C., Menesatti, P., Antonucci, F., Tripodi, S., & Travaglini, A. (2014). Multivariate statistical forecasting modeling to predict Poaceae pollen critical concentrations by meteoclimatic data. *Aerobiologia*. https://doi.org/10.1007/s10453-013-9305-3

Buters, J. T. M., Antunes, C., Galveias, A., Bergmann, K. C., Thibaudon, M., Galán, C., Schmidt-Weber, C., & Oteros, J. (2018). Pollen and spore monitoring in the world. *Clinical and Translational Allergy*. https://doi.org/10.1186/s13601-018-0197-8

Castellano-Méndez, M., Aira, M. J., Iglesias, I., Jato, V., & González-Manteiga, W. (2005). Artificial neural networks as a useful tool to predict the risk level of Betula pollen in the air. *International Journal of Biometeorology*. https://doi.org/10.1007/s00484-004-0247-x

Chappuis, C., Tummon, F., Clot, B., Konzelmann, T., Calpini, B., & Crouzy, B. (2020). Automatic pollen monitoring: First insights from hourly data. *Aerobiologia*. https://doi.org/10.1007/s10453-019-09619-6

Cowpertwait, P. S. P., & Metcalfe, A. V. (2009). *Introductory time series with R (Use R)*. Springer.

Crouzy, B., Stella, M., Konzelmann, T., Calpini, B., & Clot, B. (2016). All-optical automatic pollen identification: Towards an operational system. *Atmospheric Environment*. https://doi.org/10.1016/j.atmosenv.2016.05.062

Damialis, A., Gioulekas, D., Lazopoulou, C., Balafoutis, C., & Vokou, D. (2005). Transport of airborne pollen into the city of Thessaloniki: The effects of wind direction, speed and persistence. *International Journal of Biometeorology*. https://doi.org/10.1007/s00484-004-0229-z

Devillier, P., Bousquet, J., Salvator, H., Naline, E., Grassin-Delyle, S., & de Beaumont, O. (2016). In allergic rhinitis, work, classroom and activity impairments are weakly related to other outcome measures. *Clinical and Experimental Allergy: Journal of the British Society for Allergy and Clinical Immunology*. https://doi.org/10.1111/cea.12801

Fernández-Rodríguez, S., Tormo-Molina, R., Maya-Manzano, J. M., Silva-Palacios, I., & Gonzalo-Garijo, Á. (2014). Comparative study of the effect of distance on the daily and hourly pollen counts in a city in the south-western Iberian Peninsula. *Aerobiologia*. https://doi.org/10.1007/s10453-013-9316-0

Galán, C., Ariatti, A., Bonini, M., Clot, B., Crouzy, B., Dahl, A., & Sofiev, M. (2017). Recommended terminology for aerobiological studies. *Aerobiologia, 33*(3), 293–295.

García-Mozo, H., Yaezel, L., Oteros, J., & Galán, C. (2014). Statistical approach to the analysis of olive long-term pollen season trends in southern Spain. *The Science of the Total Environment*. https://doi.org/10.1016/j.scitotenv.2013.11.142

Geller-Bernstein, C., & Portnoy, J. M. (2019). The clinical utility of pollen counts. *Clinical Reviews in Allergy and Immunology, 57,* 340–349. https://doi.org/10.1007/s12016-018-8698-8

Glacy, J., Putnam, K., Godfrey, S., Falzon, L., Mauger, B., Samson, D., & Aronson, N. (2013). *Treatments for seasonal allergic rhinitis*. Rockville (MD): Agency for Healthcare Research and Quality (US); 2013 Jul. Report No.: 13-EHC098-EF. PMID: 23946962.

Goodfellow, I., Bengio, Y., & Courville, A. (2016). *Deep learning*. MIT Press.

Haanpää, L., Af Ursin, P., Nermes, M., Kaljonen, A., & Isolauri, E. (2018). Association of allergic diseases with children's life satisfaction: Population-based study in Finland. *British Medical Journal Open*. https://doi.org/10.1136/bmjopen-2017-019281

Harvey, A., & Scott, A. (1994). Seasonality in dynamic regression models. *The Economic Journal, 1994*(104), 1324–1345.

Howard, L. E., & Levetin, E. (2014). Ambrosia pollen in Tulsa, Oklahoma: Aerobiology, trends, and forecasting model development. *Annals of Allergy, Asthma & Immunology: Official Publication of the American College of Allergy, Asthma, & Immunology,*. https://doi.org/10.1016/j.anai.2014.08.019

Hyndman, R. J., & Athanasopoulos, G. (2018). *Forecasting: Principles and practice* (2nd ed.). Otexts: Victoria. https://otexts.org/fpp2/

Iglesias-Otero, M. A., Fernández-González, M., Rodríguez-Caride, D., Astray, G., Mejuto, J. C., & Rodríguez-Rajo, F. J. (2015). A model to forecast the risk periods of Plantago pollen allergy by using the ANN methodology. *Aerobiologia*. https://doi.org/10.1007/s10453-014-9357-z

Karatzas, K., Tsiamis, A., Charalampopoulos, A., Damialis, A., & Vokou, D. (2019). Pollen season identification for three pollen taxa in Thessaloniki, Greece: A 30-year retrospective analysis. *Aerobiologia*. https://doi.org/10.1007/s10453-019-09605-y

Kawashima, S., Thibaudon, M., Matsuda, S., Fujita, T., Lemonis, N., Clot, B., & Oliver, G. (2017). Automated pollen monitoring system using laser optics for observing seasonal changes in the concentration of total airborne pollen. *Aerobiologia*. https://doi.org/10.1007/s10453-017-9474-6

Kmenta, M., Bastl, K., Jäger, S., & Berger, U. (2014). Development of personal pollen information—the next generation of pollen information and a step forward for hay fever sufferers. *International Journal of Biometeorology*. https://doi.org/10.1007/s00484-013-0776-2

Makra, L., Matyasovszky, I., Thibaudon, M., & Bonini, M. (2011). Forecasting ragweed pollen characteristics with nonparametric regression methods over the most polluted areas in Europe. *International Journal of Biometeorology*. https://doi.org/10.1007/s00484-010-0346-9

Muzalyova, A., Brunner, J. O., Traidl-Hoffmann, C., & Damialis, A. (2019). Pollen allergy and health behavior: Patients trivializing their disease. *Aerobiologia*. https://doi.org/10.1007/s10453-019-09563-5

Myszkowska, D., & Majewska, R. (2014). Pollen grains as allergenic environmental factors–new approach to the forecasting of the pollen concentration during the season. *Annals of Agricultural and Environmental Medicine: AAEM*. https://doi.org/10.5604/12321966.1129914

Nakao, A., Nakamura, Y., & Shibata, S. (2015). The circadian clock functions as a potent regulator of allergic reaction. *Allergy*. https://doi.org/10.1111/all.12596

Nowosad, J., Stach, A., Kasprzyk, I., Chłopek, K., Dąbrowska-Zapart, K., Grewling, Ł., Latałowa, M., Pędziszewska, A., Majkowska-Wojciechowska, B., Myszkowska, D., Piotrowska-Weryszko, K., Weryszko-Chmielewska, E., Puc, M., Rapiejko, P., & Stosik, T. (2018). Statistical techniques

for modeling of Corylus, Alnus, and Betula pollen concentration in the air. *Aerobiologia*. https://doi.org/10.1007/s10453-018-9514-x

Ocana-Peinado, F., Valderrama, M. J., & Aguilera, A. M. (2008). A dynamic regression model for air pollen concentration. *Stochastic Environmental Research and Risk Assessment*. https://doi.org/10.1007/s00477-007-0153-y

Oteros, J., Pusch, G., Weichenmeier, I., Heimann, U., Möller, R., Röseler, S., Traidl-Hoffmann, C., Schmidt-Weber, C., & Buters, J. (2015). Automatic and online pollen monitoring. *International Archives of Allergy and Immunology*. https://doi.org/10.1159/000436968

Oteros, J., Sofiev, M., Smith, M., Clot, B., Damialis, A., Prank, M., Werchan, M., Wachter, R., Weber, A., Kutzora, S., Heinze, S., Herr, C., Menzel, A., Bergmann, K., Traidl-Hoffmann, C., Schmidt-Weber, C., & Buters, J. (2019). Building an automatic pollen monitoring network (ePIN): Selection of optimal sites by clustering pollen stations. *The Science of the Total Environment*. https://doi.org/10.1016/j.scitotenv.2019.06.131

Pankratz, A. (2012). *Forecasting with Dynamic Regression Models (Wiley Series in Probability and Statistics, v.935)*. Hoboken: Wiley.

Passali, D., Cingi, C., Staffa, P., Passali, F., Muluk, N. B., & Bellussi, M. L. (2018). The international study of the allergic rhinitis survey: Outcomes from 4 geographical regions. *Asia Pacific allergy*. https://doi.org/10.5415/apallergy.2018.8.e7

Pawankar, R. (2014). Allergic diseases and asthma: A global public health concern and a call to action. *The World Allergy Organization Journal*. https://doi.org/10.1186/1939-4551-7-12

Piotrowska, K. (2012). Forecasting the Poaceae pollen season in eastern Poland. *Grana*. https://doi.org/10.1080/00173134.2012.659204

Puc, M. (2012). Artificial neural network model of the relationship between Betula pollen and meteorological factors in Szczecin (Poland). *International Journal of Biometeorology*. https://doi.org/10.1007/s00484-011-0446-1

Ritenberga, O., Sofiev, M., Kirillova, V., Kalnina, L., & Genikhovich, E. (2016). Statistical modelling of non-stationary processes of atmospheric pollution from natural sources: Example of birch pollen. *Agricultural and Forest Meteorology*. https://doi.org/10.1016/j.agrformet.2016.05.016

Rodríguez-Rajo, F. J., Valencia-Barrera, R. M., Vega-Maray, A. M., Suárez, F. J., Fernández-González, D., & Jato, V. (2006). Prediction of airborne Alnus concentration by using ARIMA models. *Annals of Agricultural and Environmental Medicine: AAEM, 2006*(13), 25–32.

Ríos, B., Torres-Jardón, R., Ramírez-Arriaga, E., Martínez-Bernal, A., & Rosas, I. (2016). Diurnal variations of airborne pollen concentration and the effect of ambient temperature in three sites of Mexico City. *International Journal of Biometeorology, 60,* 771–787. https://doi.org/10.1007/s00484-015-1061-3

Schiele, J.,Rabe F., SchmittGlaser, M., Haring Brunner, J. O.Bauer, B.Schuller, B.Traidl-Hoffmann, C., Damialis A. (2019). Automated Classification of Airborne Pollen using Neural Networks. *Annual International Conference of the IEEE Engineering in Medicine and Biology Society. IEEE Engineering in Medicine and Biology Society. Annual International Conference*, doi: https://doi.org/10.1109/EMBC.2019.8856910.

Schmitz, R., Thamm, M., Ellert, U., Kalcklösch, M., & Schlaud, M. (2014). Verbreitung häufiger Allergien bei Kindern und Jugendlichen in Deutschland: Ergebnisse der KiGGS-Studie - Erste Folgebefragung (KiGGS Welle 1). *Bundesgesundheitsblatt, Gesundheitsforschung, Gesundheitsschutz*. https://doi.org/10.1007/s00103-014-1975-7

Simoleit, A., Gauger, U., Mücke, H.-G., Werchan, M., Obstová, B., Zuberbier, T., & Bergmann, K.-C. (2016). Intradiurnal patterns of allergenic airborne pollen near a city motorway in Berlin, Germany. *Aerobiologia, 32,* 199–209. https://doi.org/10.1007/s10453-015-9390-6

Sofiev, M. (2019). On possibilities of assimilation of near-real-time pollen data by atmospheric composition models. *Aerobiologia, 35,* 523–531. https://doi.org/10.1007/s10453-019-09583-1

Sofiev, M., & Bergmann, K.-C. (Eds.). (2013). *Allergenic pollen: A review of the production, release, distribution and health impacts*. Springer.

Sánchez Mesa, J. A., Galán, C., & Hervás, C. (2005). The use of discriminant analysis and neural networks to forecast the severity of the Poaceae pollen season in a region with a typical Mediterranean climate. *International Journal of Biometeorology*. https://doi.org/10.1007/s00484-005-0260-8

Taskaya-Temizel, T., & Casey, M. C. (2005). A comparative study of autoregressive neural network hybrids. *Neural Networks: The Official Journal of the International Neural Network Society*. https://doi.org/10.1016/j.neunet.2005.06.003

Twomey, J. M., & Smith, A. E. (1995). Performance measures, consistency, and power for artificial neural network models. *Mathematical and Computer Modelling*. https://doi.org/10.1016/0895-7177(94)00207-5

Valencia, J. A., Astray, G., Fernández-González, M., Aira, M. J., & Rodríguez-Rajo, F. J. (2019). Assessment of neural networks and time series analysis to forecast airborne Parietaria pollen presence in the Atlantic coastal regions. *International Journal of Biometeorology*. https://doi.org/10.1007/s00484-019-01688-z

Weger, L. A., Bergmann, K. C., Rantio-Lehtimäki, A., Dahl, A., Buters, J., Déchamp, C., Belmonte, J., Thibaudon, M., Cecchi, L., Besancenot, J. P., Galán, C., & Waisel, Y. (2013). Impact of pollen. In M. Sofiev & K. C. Bergmann (Eds.), *Allergenic pollen: A review of the production, release, distribution and health impacts*. Dordrecht: Springer.

Zewdie, G. K., Liu, X., Wu, D., Lary, D. J., & Levetin, E. (2019). Applying machine learning to forecast daily Ambrosia pollen using environmental and NEXRAD parameters. *Environmental Monitoring and Assessment*. https://doi.org/10.1007/s10661-019-7428-x

Ziello, C., Sparks, T. H., Estrella, N., Belmonte, J., Bergmann, K. C., Bucher, E., Brighetti, M. A., Damialis, A., Detandt, M., Galán, C., Gehrig, R., Grewling, L., Bustillo, A. M. G., Hallsdóttir, M., Kockhans-Bieda, M. C., Linares, C., Myszkowska, D., Pàldy, A., Sánchez, A., … Thibaudon, M. (2012). Changes to airborne pollen counts across Europe. *PloS one*. https://doi.org/10.1371/journal.pone.0034076

Ziska, L. H., Makra, L., Harry, S. K., Bruffaerts, N., Hendrickx, M., Coates, F., Saarto, A., Thibaudon, M., Oliver, G., Damialis, A., Charalampopoulos, A., Vokou, D., Heiđ‗marsson, S., Guđjohnsen, E., Bonini, M., Oh, J., Sullivan, K., Ford, L., Brooks, G. D., Myszkowska, D., et al. (2019). Temperature-related changes in airborne allergenic pollen abundance and seasonality across the northern hemisphere: A retrospective data analysis. *The Lancet Planetary Health*. https://doi.org/10.1016/S2542-5196(19)30015-4

Ščevková, J., Dušička, J., Mičieta, K., & Somorčík, J. (2015). Diurnal variation in airborne pollen concentration of six allergenic tree taxa and its relationship with meteorological parameters. *Aerobiologia*. https://doi.org/10.1007/s10453-015-9379-1