

Presentation of an Algorithm for the Automated Estimation and Gap-Filling of Missing Data in Daily Weather Records

J.S. Gosselin^{a,*}, R. Martel^a, C. Rivard^b

^a*Institut national de la recherche scientifique, Centre Eau Terre Environnement, 490 rue de la Couronne, Quebec City, Quebec, Canada*

^b*Geological Survey of Canada, Quebec Division, 490 rue de la Couronne, Quebec City, Quebec, Canada*

Abstract

Daily weather data are useful in several areas of Earth sciences, including hydrology, hydrogeology and agronomy. However, weather datasets are often incomplete. The estimation of missing data can be a complex and tedious task. This is particularly the case for daily precipitation because of their high spatial and temporal variability. A user friendly, menu-driven, and interactive computer program for rapid and automatic completion of daily climatological series has been developed. Missing data for a given weather station are estimated using a multiple linear regression model, generated using data from nearby stations. For daily precipitation, it is possible to activate an option that forces the algorithm to preserve the probability distribution of data. This is an advantage over conventional approaches that tend to overestimate the number of wet days and underestimate the high intensity precipitation events. The software also allows downloading and automatic formatting of raw data available on the Environment Canada website. The software is demonstrated for two weather station located in Monteregie Est region, southern Quebec. Cross-validation was used to check the method and to define the optimal parameters to minimize the error in estimating missing daily precipitation.

Keywords: heat transport, recharge assessment, uncertainty analysis, subsurface temperature time series

1. Introduction

Climate data are useful in several fields of Earth sciences, including hydrology, hydrogeology and agronomy. However, weather datasets are, most of the time, incomplete. This can represent a major hindrance in various applications, such as for the use of hydrological or hydrogeological models that heavily depend on these data. Filling the gaps in weather datasets can quickly become a tedious task as the size of the data records and the number of stations increase. Moreover, it can also be quite complex when aspects such as time-efficiency of the method and accuracy of the estimated missing values are taken into account. This is particularly true for the estimation of missing daily precipitation data because of their high spatial and temporal variability (Simolo et al., 2010). Although there are various

*Corresponding authors

Email address: jnsebgosselin@gmail.com (J.S. Gosselin)

methods to estimate missing daily weather data that are well covered in textbooks and technical papers, few tools to perform this task efficiently and automatically are available.

This paper presents an open source algorithm, written in the Python programming language, that can be used to automatically fill the gaps in daily weather datasets and to assess the uncertainty on the estimated values. An application of the method, using the WHAT software, is also presented for the Montérégie Est study area, located in southern Quebec, Canada.

for filling the gaps in the daily weather datasets of a given weather station (hereafter called the target station) using data from the neighboring stations.

Can also be used to compute daily potential evapotranspiration.

The algorithm is available for free at : . In addition, it is included as part of the WHAT software, which blablabla.

Secondly, the program also includes an automated, robust, and efficient method to quickly and easily fill the gaps in the daily weather datasets downloaded from the CDCD. WHAT also includes a cross-validation resampling algorithm to conveniently validate and assess the uncertainty of the estimated missing values.

In addition the algorithm can also b

A guide for the operation of the software Gosselin (2015) is available for download at this web address: <https://github.com/jnsebgosselin/WHAT>.

2. Theory

The algorithm described in this paper is based on the implementation of the classical MLR (Multiple Linear Regression) method presented in Eischeid et al. (2000). The MLR method is a robust and well known spatial interpolation technique that can indirectly account for local effects, such as topography, land cover, land use and surface water. While creating serially complete daily datasets of air temperature and total precipitation for the western U.S., Eischeid et al. (2000) found that the MLR method consistently outperformed the other classical methods tested (normal ratio, inverse distance, optimal interpolation, and single best estimator). The same result was also found by Xia et al. (1999) for a study in Bavaria, Germany. Moreover, in a study conducted in Iran for different climate conditions (dry to extra humid conditions), Kashani and Dinpashoh (2011) found that the estimation obtained with the MLR method compared well with those obtained with more recent methods, more specifically the artificial neural network (reference) and the genetic programming (references) techniques.

Figure 1 presents a flowchart of the gap-filling algorithm. The algorithm consists of two nested loops: the external 'Loop A' iterates over the weather variables contained in the dataset of the target station (min, max, and mean air temperature and total precipitation), while the inner 'Loop B' iterates over the missing values in the data series of the current weather variable in 'Loop A'. Each missing value is estimated independently with a two-step procedure in 'Loop B': the first step consists in the selection of the neighboring stations, while the second step consists in building a MLR model, estimating the missing value, and filling the corresponding gap in the data series.

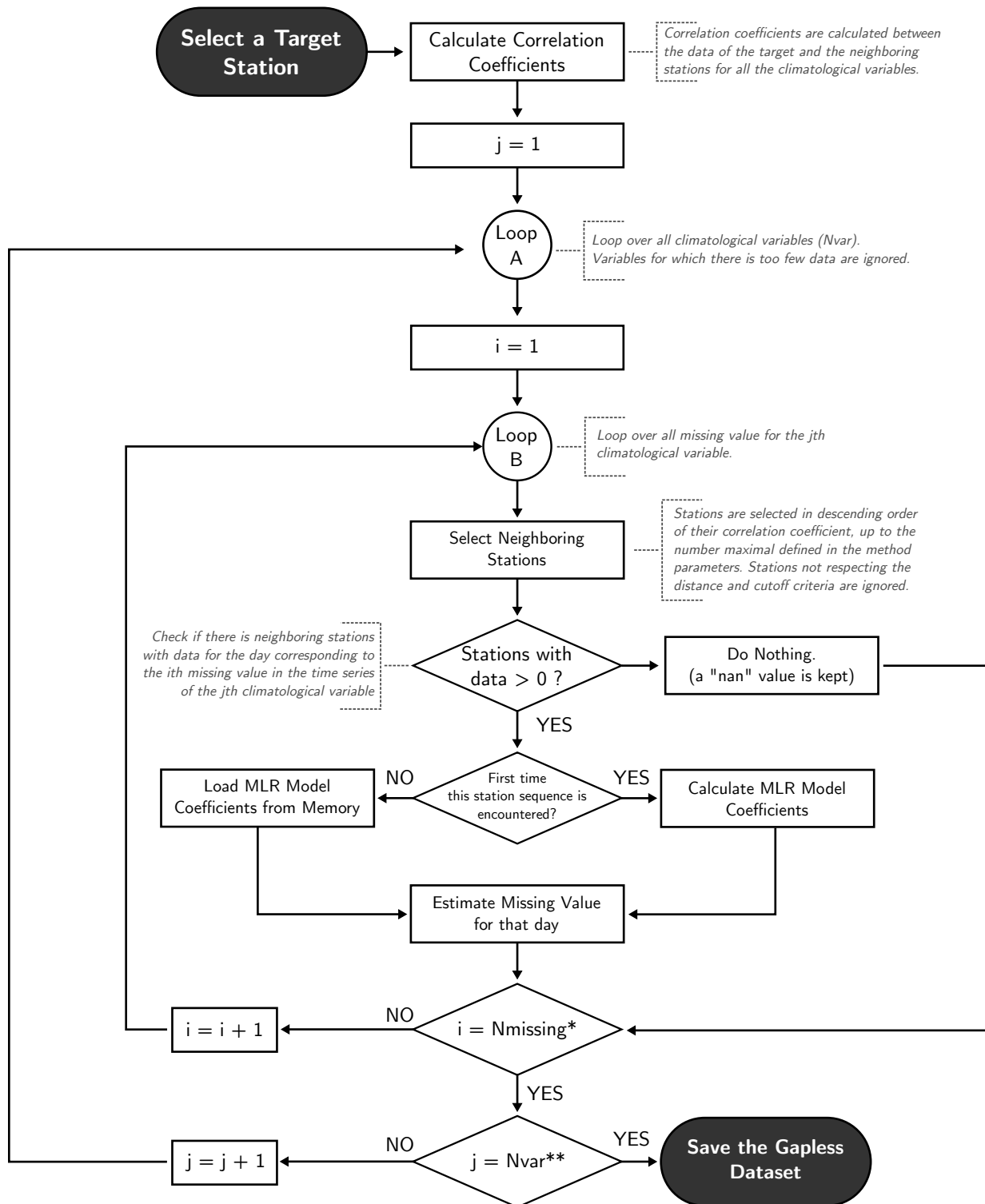


Figure 1

2.1. Correlation Coefficients Calculations

Correlation coefficients are first calculated between the available data of the target station and those of the neighboring stations for the j^{th} weather variable contained in the dataset. The correlation coefficients are calculated individually for each neighboring station using all the available data. Neighboring stations that have less than 182 days (half a year) of synchronous data with the target station or with a correlation coefficient below a value of 0.35 will not be used to fill the gaps in the data series of the current weather variable. The 0.35 threshold defined for the correlation coefficient is based on the value used by Eischeid et al. (2000) in their application of the method.

Moreover, the correlation between the data of the neighboring and target stations will generally decrease as the horizontal distance and elevation difference between them increase. Therefore, it is possible to discard neighboring stations that are located further of the target station than specified thresholds, either in the horizontal or the vertical direction. The default values are set to 100 km and 350 m for the horizontal and vertical distance respectively based on the literature Tronci et al. (1986); Xia et al. (1999); Simolo et al. (2010).

2.2. Selection of the Neighboring Stations

As stated by Eischeid et al. (2000), the selection of neighboring stations is critically important for the accurate estimation of missing weather data. Problems arise though because the list of neighboring stations with available data can vary from one day to the other. Therefore, it is not possible to use a single MLR model to estimate the missing values in the dataset of the target station all at once. The selection of the neighboring stations and the generation of a MLR model must be done instead individually for each day with a missing value in the dataset of the target stations.

Neighboring stations with available data are selected in descending order of their correlation coefficient, up to the maximal number of stations that is specified as a parameter of the algorithm. The default value for the maximal number of neighboring station used for the generation of the MLR models is 4. Tests run by Eischeid et al. (2000) showed that using more than 4 neighboring stations did not significantly improve, and may even have degraded, the accuracy of the estimate. If for a given day with a missing value, no neighboring stations have a measured value, no calculation is done and a 'NaN' value is kept in the dataset.

2.3. Generation of the Multiple Linear Regression Model

Each time a MLR model is generated for a given sequence of neighboring stations, the resulting model parameters are stored into memory. Therefore, after the neighboring stations have been selected for a given day with a missing data (section 2.2), the program checks if this sequence of selected stations has already been encountered before for the current weather variable. If so, the stored MLR parameters will be used directly to estimate the missing data for the current day. Otherwise, the model will generate a new MLR model and will store the results into memory. Since a MLR model is generated only one time for a given sequence of neighboring stations, the algorithm becomes faster with time.

The MLR model can be generated using either an Ordinary Least Square (OLS) or a Least Absolute Deviations (LAD) criteria. For the OLS criteria, the MLR model is obtained by solving the linear matrix equation $\mathbf{X}\mathbf{a} = \mathbf{Y}$ by computing the $N \times 1$ parameter vector \mathbf{a} that minimizes the Euclidean L2-norm $\|\mathbf{Y} - \mathbf{X}\mathbf{a}\|_2$, where \mathbf{Y} is a $M \times 1$ vector containing the M daily data of the target station and \mathbf{X} is a $M \times N$ matrix containing the M synchronous daily data of the N selected neighboring stations. (REFERENCE: Numpy Documentation)

However, daily precipitation series are generally characterized by long-tailed, positively skewed, distributions. In this case, the generation of the MLR model using the more robust LAD method is a more appropriate approach that will yield more reliable results than using the OLS method as described above Menke (1989); Eischeid et al. (2000). Resolution of the MLR model with the LAD method is achieved in the gap-filling algorithm using an iterative reweighted least-squares method. The downside in using the LAD method compared to using the OLS is an increase in computation time by about a factor 10.

2.4. Estimating Missing Daily Values

Once the parameters of the MLR model are known, the missing value for the corresponding day at time t_i can be estimated as follows:

$$Y_p|_{t_i} = a_0 + \sum_{k=1}^N a_k \cdot X_k(t_i) \quad (1)$$

where $Y_p(t_i)$ is the value estimated at time t_i for the j_i th weather variable in the dataset of the target station, $X_k(t_i)$ is the synchronous available data of the k_i th neighboring stations, a_k are the regression coefficients, and N is the total number of selected neighboring stations used for the regression.

When all the missing values in the dataset of target station have been estimated and filled, the resulting gapless time series is saved in a file with a '.out' extension. Moreover, detailed information about the estimated values are also saved in an accompanying '.log' file. The outputs of the gap-filling algorithm are discussed in more details in section 3.3.

2.5. Uncertainty of the estimated values

Each time a MLR model is generated from a new sequence of neighboring stations for a given weather variable, the resulting model is also used to estimate values for the days in the data series of the target station for which there exists a measured value. The accuracy of the new MLR model is then approximated by computing a Root-Mean-Square Error (RMSE) between the values estimated with the MLR model and the respective measured values. The RMSE thus calculated is saved, along with the estimated value, in the '.log' file.

The algorithm also includes a cross-validation resampling procedure to estimate the accuracy of the method, in addition to fill the gaps in the dataset. The procedure to enable this functionality in the algorithm is presented in section 3.2.

More specifically, the procedure consists in estimating alternately a weather data value for each day of the selected station's dataset, even for days for which data are not missing. In other words, the loop B in the flowchart of Figure

8.1 will iterate over all the days of the dataset and not only over days for which there is a missing data. In addition, the memory feature will be deactivated and a MLR model will be estimated for each day independently. If a measured value is present for the current day being estimated, this value will be temporarily discarded from the data series to avoid self-influence of the observation on the estimation procedure. Before estimating a value for a given day, the corresponding measured data in the dataset of the target station is temporarily discarded to avoid self-influence of this observation on the generation of the MLR model. The model is then generated and used to estimate a value on this given day and the corresponding observed data is put back in the dataset of the selected station. When a value for every day of the dataset has thus been estimated, the estimated values are saved in a tsv (tabular-separated values) file in the *Output* folder with the extension “.err”, along with the “.log” and “.out” files described in ???. The accuracy of the method can then be estimated by computing the RMSE between the estimated weather data and the respective non-missing observations in the original dataset of the selected station. Activating this feature will significantly increase the computation time of the gap filling procedure, especially if the least absolute deviation regression model is selected, but can provide interesting insights on the performance of the procedure for the specific datasets used for a given project.

Moreover, the graphs that are presented in Section blabla are automatically generated when the gap-filling routine is completed. There is currently no tool provided in WHAT to directly analyze the results from the Jackknife procedure. However, all the source code that has been written for the production of the figures of Section 8.4 can be downloaded freely on GitHub at (<https://github.com/jnsebgosselin/WHAT>).

3. Operation

3.1. Input Data

3.2. Parameters

3.3. Output

4. Application: Monteregie Est Case Study

The Monteregie Est region is located in southern Quebec, Canada, on the south shore of the St. Lawrence River. It covers a total area of 9032 km², from the St. Lawrence River at its northern limit to the border of the United States (states of New York and Vermont) at its southern limit (see Figure X).

This region has been the subject of an extensive characterization project within the ‘Programme d’acquisition de connaissances sur les eaux souterraines du Québec’ (PACES) whose main objective was to prepare a realistic and concrete picture of the groundwater resources for the region (?).

135 4.1. Study Area

The climate is characterized by significant seasonal differences in temperature, resulting in warm summers and cold winters. Precipitation, as rain or snow, are distributed rather evenly throughout the year.

Among all the weather stations for which data were available in and around the study area, a total of 32 was selected based on the availability and continuity of the weather data between 1980 and 2014. A list of these selected
140 stations is presented in Table X with their coordinates, altitude, total time periods for which data were available, mean annual cumulative precipitation, and mean annual air temperature. Most of these information are generated automatically by WHAT in the file “weather_datasets_summary.log” (see Section ??).

4.2. Materials and Method

For this purpose, the Canadian Daily Climate Database (CDCD), owned by the Government of Canada, contains
145 daily data for air temperature and precipitation dating back to 1840 to the present for about 8450 stations distributed across Canada. Data can be downloaded manually on the Government of Canada website (www.climate.weather.gc.ca) for each year individually and saved in a csv file. This process involves a lot of repetitive manipulations and is a time consuming task. Moreover, the re-organization of the individual data files, saved for each year separately, in a more convenient format can also represent a tedious task when done manually. Alternatively, it is possible to order
150 a DVD containing the entire database for a small fee. This option has the disadvantage of only providing an image in time as data cannot be updated.

WHAT (Well Hydrograph Analysis Toolbox) is a computer program that addresses the aforementioned issues (Gosselin et al., 2015). Firstly, it provides a graphical interface to the online CDCD that allows to search for stations interactively using location coordinates, download the available data for the selected weather stations, and automati-
155 cally organize the data in a more convenient format.

Pour l'ensemble de ces stations, les précipitations totales annuelles sont d'environ 1100 mm/y en moyenne. Les précipitations totales les plus élevées sont observées à la station de Brome (~1280 mm/y et les plus faibles à la station de Sorel (~960 mm/year). La température annuelle moyenne dans la région d'étude est de 5.9 °C, variant de 4.3 to 6.7 °C tandis que les températures mensuelles moyennes fluctuent entre -12 to 21 °C. Les températures mensuelles
160 minimales sont observées en janvier (-17.1 to -13.6 °C) tandis que les températures mensuelles maximales sont observées en juillet (24 to 26.7 °C). Les températures les plus élevées sont généralement observées aux stations de Philipsburg et Saint-Bernard-de-Lacolle (température annuelle moyenne de 6.7 °C) et les plus faibles à la station de Bonsecours (température annuelle moyenne de 4.3 °C)

Les figures 1.3 et 1.4 illustrent respectivement les variations spatiales des précipitations totales annuelles et de la
165 température moyenne annuelle pour la période de 1970-2000. Les valeurs présentées sur ces figures ont été interpolées par krigeage ordinaire sur une grille de 250 x 250 m, à partir des valeurs rapportées pour les 16 stations actives mentionnées ci-haut. Dans la région d'étude, la tendance générale indique que les précipitations annuelles totales diminuent du sud-sud-est vers le nord-nord-ouest et que les températures annuelles moyennes diminuent du sud-ouest

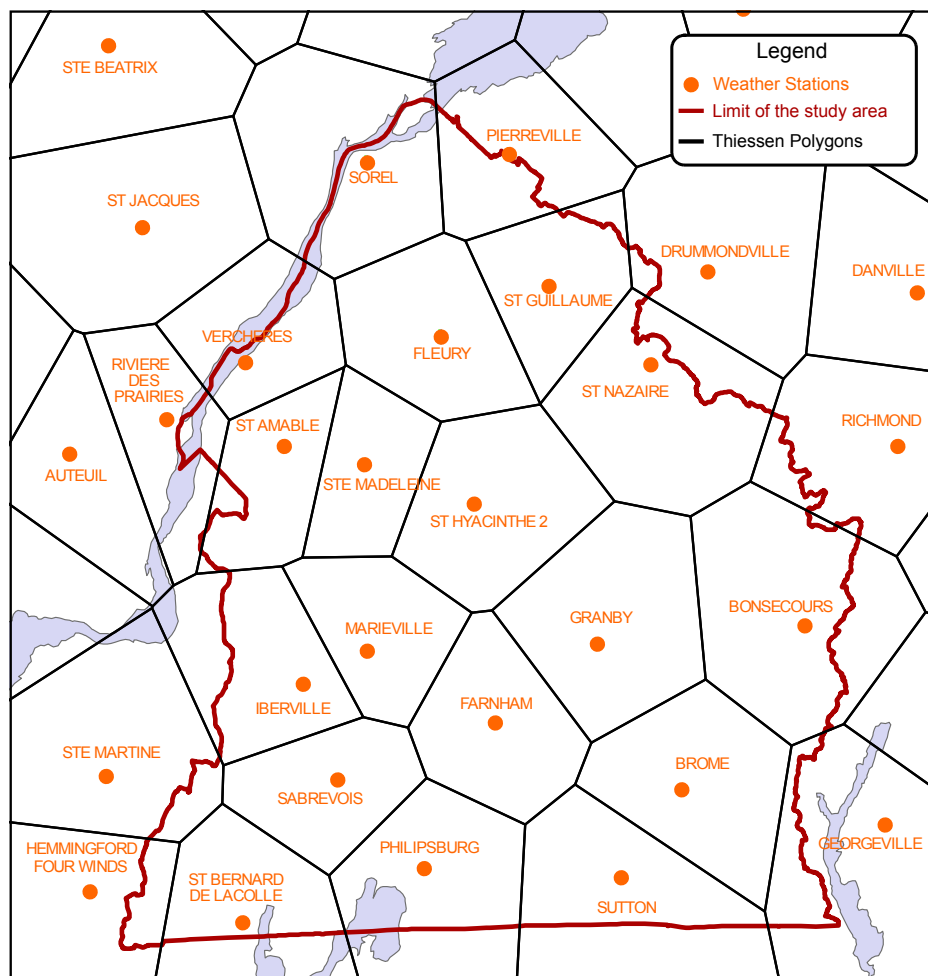


Figure 2: Locations of the weather stations in the Montérégie Est area.

vers le nord-est. Outre l'influence de la latitude, la température de la région est également influencée par la présence
 170 des Appalaches au sud-est et du fleuve Saint-Laurent au nord-ouest.

The weather network of the Montérégie Est region, located in the province of Quebec, Canada, has been used to
 test the method. This region feature strongly variable topography and land cover conditions. The network is presented
 in figure X. Also, stations from bordering states were extracted to improve the spatial distribution of sites surrounding
 target stations located near state borders.

175 Daily weather data for 32 weather stations in and around the Montérégie Est area were also retrieved from the
 Canadian Daily Climate Database (CDCD) with the software WHAT for the years 2000 to 2012. Missing values
 in the weather time series were also estimated with WHAT to produce gapless meteorological records of daily air
 temperature and precipitation.

Tests have shown that inclusion of more than four stations does not significantly improve the interpolation and
 180 may in fact degrade the estimate.

4.3. Results and Discussion

The quality of the estimates is strongly affected by seasonality. Stations at higher elevations are difficult to estimate accurately, in large part because of the topographical diversity of the surrounding stations leading to degradation of spatial coherence among stations.

The tendency for all of the methods to have a negative bias is indicative of the nature of precipitation distributions to be positively skewed (interpolated values will tend to cluster about the median error rather than the mean).

According to Xia et al. (1999), the two most important factors in climatology are the inter-correlations in the station network, and the seasonal variations in the relations between the stations.

5. Discussion

However, weighing and regression-based techniques, including the MLR method, all tend to overestimate the number of rainy days, while heavy precipitation events are systematically underestimated. Therefore, the rainfall probability distribution is usually not preserved with these techniques). However, Simolo et al. (2010) have proposed a two-step procedure to modify the MLR method to address these issues.

An alternative approach would have been to calculate the correlation coefficient with a subset of data from the target series centered around the missing value, as it was done in Simolo et al. (2010) for instance. This approach allows for a better representation of the seasonal variations in the relationships between the stations. The downsides include a more complex algorithm to implement and a reduction of the method robustness and efficiency.

References

- Eischeid, J.K., Pasteris, P.A., Diaz, H.F., Plantico, M.S., Lott, N.J., 2000. Creating a serially complete, national daily time series of temperature and precipitation for the Western United States. *Journal of Applied Meteorology* 39, 1580–1591. doi:1520-0450(2000)039<1580:CASCND>2.0.CO;2. d181.
- Gosselin, J.S., 2015. WHAT (Well Hydrograph Analysis Toolbox). URL: <https://github.com/jnsebgosselin/WHAT>.
- Gosselin, J.S., Rivard, C., Martel, R., 2015. User Manual for WHAT. Document written for software version 4.1.7-beta. Technical Report. INRS-ETE, Quebec City, Qc, Can. URL: <https://github.com/jnsebgosselin/WHAT/raw/master/WHATMANUAL/WHATMANUAL.pdf>.
- Kashani, M.H., Dinpashoh, Y., 2011. Evaluation of efficiency of different estimation methods for missing climatological data. *Stochastic Environmental Research and Risk Assessment* 26, 59–71. URL: <http://link.springer.com/article/10.1007/s00477-011-0536-y>, doi:10.1007/s00477-011-0536-y. d265.
- Menke, W., 1989. *Geophysical Data Analysis: Discrete Inverse Theory*. 1 edition ed., Academic Press, San Diego.
- Simolo, C., Brunetti, M., Maugeri, M., Nanni, T., 2010. Improving estimation of missing values in daily precipitation series by a probability density function-preserving approach. *International Journal of Climatology* 30, 1564–1576. doi:10.1002/joc.1992. d184.
- Tronci, N., Molteni, F., Bozzini, M., 1986. A comparison of local approximation methods for the analysis of meteorological data. *Archives for Meteorology, Geophysics, and Bioclimatology, Series B* 36, 189–211. URL: <http://link.springer.com/article/10.1007/BF02278328>, doi:10.1007/BF02278328. d218.
- Xia, Y., Fabian, P., Stohl, A., Winterhalter, M., 1999. Forest climatology: estimation of missing values for Bavaria, Germany. *Agricultural and Forest Meteorology* 96, 131–144. doi:S0168-1923(99)00056-8. d182.