ISSN: 0711-2440

Spatial pattern regression for gridded meteorological data: A precipitation and temperature case study

V. Houssou, J. Carreau

G-2025-48 July 2025

La collection Les Cahiers du GERAD est constituée des travaux de recherche menés par nos membres. La plupart de ces documents de travail a été soumis à des revues avec comité de révision. Lorsqu'un document est accepté et publié, le pdf original est retiré si c'est nécessaire et un lien vers l'article publié est ajouté.

Citation suggérée : V. Houssou, J. Carreau (Juillet 2025). Spatial pattern regression for gridded meteorological data: A precipitation and temperature case study, Rapport technique, Les Cahiers du GERAD G– 2025–48, GERAD, HEC Montréal, Canada.

Avant de citer ce rapport technique, veuillez visiter notre site Web (https://www.gerad.ca/fr/papers/G-2025-48) afin de mettre à jour vos données de référence, s'il a été publié dans une revue scientifique

The series Les Cahiers du GERAD consists of working papers carried out by our members. Most of these pre-prints have been submitted to peer-reviewed journals. When accepted and published, if necessary, the original pdf is removed and a link to the published article is added.

Suggested citation: V. Houssou, J. Carreau (July 2025). Spatial pattern regression for gridded meteorological data: A precipitation and temperature case study, Technical report, Les Cahiers du GERAD G–2025–48, GERAD, HEC Montréal, Canada.

Before citing this technical report, please visit our website (https://www.gerad.ca/en/papers/G-2025-48) to update your reference data, if it has been published in a scientific journal.

La publication de ces rapports de recherche est rendue possible grâce au soutien de HEC Montréal, Polytechnique Montréal, Université Mc-Gill, Université du Québec à Montréal, ainsi que du Fonds de recherche du Québec – Nature et technologies.

Dépôt légal – Bibliothèque et Archives nationales du Québec, 2025 – Bibliothèque et Archives Canada, 2025 The publication of these research reports is made possible thanks to the support of HEC Montréal, Polytechnique Montréal, McGill University, Université du Québec à Montréal, as well as the Fonds de recherche du Québec – Nature et technologies.

Legal deposit – Bibliothèque et Archives nationales du Québec, 2025 – Library and Archives Canada, 2025

GERAD HEC Montréal 3000, chemin de la Côte-Sainte-Catherine Montréal (Québec) Canada H3T 2A7 Tél.: 514 340-6053 Téléc.: 514 340-5665 info@gerad.ca www.gerad.ca

Spatial pattern regression for gridded meteorological data: A precipitation and temperature case study

Vihotogbé Houssou ^{a, b} Julie Carreau ^{a, b}

- ^a Département de mathématiques et de génie industriel, Polytechnique Montréal (Qc), Canada
- ^b Groupe d'Études et de Recherche en Analyse des Décisions (GERAD), Montréal (QC), Canada

vihotogbe.houssou@polymtl.ca
julie.carreau@polymtl.ca

July 2025 Les Cahiers du GERAD G-2025-48

Copyright © 2025 Houssou, Carreau

Les textes publiés dans la série des rapports de recherche Les Cahiers du GERAD n'engagent que la responsabilité de leurs auteurs. Les auteurs conservent leur droit d'auteur et leurs droits moraux sur leurs publications et les utilisateurs s'engagent à reconnaître et respecter les exigences légales associées à ces droits. Ainsi, les utilisateurs:

- Peuvent télécharger et imprimer une copie de toute publication du portail public aux fins d'étude ou de recherche privée;
- Ne peuvent pas distribuer le matériel ou l'utiliser pour une activité à but lucratif ou pour un gain commercial;
- Peuvent distribuer gratuitement l'URL identifiant la publication.

Si vous pensez que ce document enfreint le droit d'auteur, contacteznous en fournissant des détails. Nous supprimerons immédiatement l'accès au travail et enquêterons sur votre demande. The authors are exclusively responsible for the content of their research papers published in the series *Les Cahiers du GERAD*. Copyright and moral rights for the publications are retained by the authors and the users must commit themselves to recognize and abide the legal requirements associated with these rights. Thus, users:

- May download and print one copy of any publication from the public portal for the purpose of private study or research;
- May not further distribute the material or use it for any profitmaking activity or commercial gain;
- May freely distribute the URL identifying the publication.

If you believe that this document breaches copyright please contact us providing details, and we will remove access to the work immediately and investigate your claim.

Abstract: We introduce Spatial Pattern Regression (SPR), a method to generate gridded historical meteorological data for climate adaptation. SPR operates in two steps: first extracting spatial structure from high-resolution regional climate model (RCM) simulations as eigenvectors, then using them in linear regression to reconstruct complete gridded fields from station observations at each time step. We compare SPR with standard interpolation methods using data from RCM simulations, where virtual stations are a subset of grid cells and interpolation is done on the rest. Thirty graded case studies are created by varying three factors: region location, size, and network density. Daily precipitation, maximum temperature, and minimum temperature are considered. Results show SPR outperforms standard methods across all three variables for most of the graded case studies. A stress-test with very low network density confirms SPR's robustness. Finally, we systematically assessed how each graded factor affects SPR's performance.

Keywords: Spatial interpolation, principal component analysis, linear regression, high resolution climate simulations, synthetic data framework

1 Introduction

Climate change is intensifying meteorological hazards like heatwaves, extreme weather events, and flooding, triggering an escalating cascade of economic and societal impacts (Warren et al., 2022). To support better adaptation measures, impact studies rely on high-resolution gridded meteorological data that must accurately represent extreme events and spatial heterogeneity. For instance, Lucas-Picher et al. (2020) aimed to reproduce the extreme flood of the Richelieu River in southern Quebec, Canada, which occurred in the spring of 2011. A key factor in achieving this was the use of a recent gridded hydrometeorological dataset—including precipitation, maximum, and minimum temperatures—at a resolution of approximately 7 km, which accounted for orographic precipitation (Livneh et al., 2015). In addition, Lauer et al. (2023) investigated the urban heat island effect during two high-temperature events lasting 2 to 3 days by evaluating various land use–based mitigation strategies. To support their analysis, they relied on high-resolution numerical weather prediction data—including surface air temperature, dew point temperature, and rainfall—at a horizontal resolution of 250 meters.

In most cases, gridded meteorological data are obtained through one of three main approaches: spatial interpolation, physics-based models such as numerical weather prediction or regional climate models, or reanalysis datasets. Spatial interpolation—the first approach—involves statistical techniques that use observed values of meteorological variables at gauged sites to estimate values at ungauged locations. To address the sparsity of gauged locations, auxiliary data available on high-resolution grids and predictive of the meteorological variable of interest are often used. For example, Livneh et al. (2015) used a technique called inverse distance weighting, in which observations from neighboring sites are weighted inversely proportional to their distance from the target location. To enhance the interpolation, orographic scaling was applied using data that incorporate topographic information. Another example is Werner et al. (2019), who used a thin-plate spline interpolation algorithm that incorporates highresolution gridded climatology data. In some cases, impact studies rely on gridded meteorological data produced by numerical weather prediction (NWP) models—the second of the three approaches listed above. NWP uses physics-based models—complex numerical frameworks that integrate multiple physical processes such as atmospheric dynamics, convection, and radiation. These models aim to predict the state of the atmosphere as accurately as possible, supporting both short- to medium-range weather forecasts (from minutes to months) and longer-term climate projections (from months to several decades). In the latter case, the NWP model is typically referred to as a climate model. This type of gridded data was used in the urban heat island study discussed above (Lauer et al., 2023). Another example of such gridded data is Climex, a 50-member ensemble of regional climate simulations with a spatial resolution of approximately 11 km (Leduc et al., 2019). In particular, Faghih and Brissette (2023) used Climex data—after applying a bias correction step—to study the effect of climate change on extreme rainfall and flooding. Lastly, reanalysis datasets—the third approach—are gridded meteorological data produced by numerical weather prediction models that incorporate observational data through data assimilation. Reanalysis aims to address some of the limitations of raw NWP outputs, such as discrepancies with observed weather conditions and systematic biases. An example of a high-resolution reanalysis dataset is CaSR, which provides precipitation and other surface variables at a spatial resolution of 10 km over North America (Gasset et al., 2021). Another prominent example is ERA-Land, a global reanalysis dataset with a spatial resolution of approximately 9 km (Muñoz Sabater et al., 2021).

We present an approach called Spatial Pattern Regression (SPR), which can be viewed as an intermediate between spatial interpolation and reanalysis. SPR relies on information about the spatial structure present in gridded meteorological data produced by a climate model—typically a regional climate model (RCM), which offers sufficiently high resolution. The underlying hypothesis is that the RCM adequately captures the spatial structure of the meteorological variable of interest. This assumption is similar to that made by spatial interpolation methods that rely on RCM-based climatology (see Werner et al. (2019), for instance). SPR operates in two steps. In the first step, the spatial structure is extracted from high-resolution RCM data over the study region. In our implementation, principal com-

ponent analysis (PCA) is used to identify spatial patterns, which are represented by the eigenvectors derived from PCA. In this case, the complete field—i.e., the values of the meteorological variable over the RCM grid—can be represented as a linear combination of the spatial patterns. In the second step, linear regression is used to reconstruct the full field over the study region for a given time step (in our application, a day), using observations of the meteorological variable at a number of sites within the region. Since SPR relies on two statistical techniques—PCA and linear regression—it can be viewed as a spatial interpolation method that blends spatial information derived from RCM data with temporal information from gauged stations. At the same time, because it systematically exploits RCM data by combining it with observations, it also shares similarities with reanalysis. However, unlike reanalysis, the period covered by the RCM data does not need to coincide with the observation period. See § 4 for a detailed explanation of the SPR methodology.

SPR is compared with standard spatial interpolation methods—considered baseline approaches and frequently used in the literature—namely, inverse distance weighting, ordinary kriging, and kriging with external drift. See § 3 for their description. To evaluate and compare the spatial interpolation methods, we introduce graded case studies based on a synthetic data framework, where the rationale is to control the experimental setup and enable systematic comparisons. The field—i.e., high-resolution gridded meteorological data at a given time step—provided by an RCM is considered the ground truth. A subset of grid cell values is treated as "observations," with interpolation aiming to reconstruct the values at the remaining grid cells. This setup allows us to generate a large number of graded case studies, which would not be feasible if ground truth were based on gauged network observations. See § 2 for a detailed description of the graded case studies. We consider three meteorological variables at daily resolution—precipitation, minimum temperature, and maximum temperature—which are commonly used in spatial interpolation for hydrological studies, see for instance Lucas-Picher et al. (2020). All spatial interpolation methods are applied sequentially, processing one day and one meteorological variable at a time.

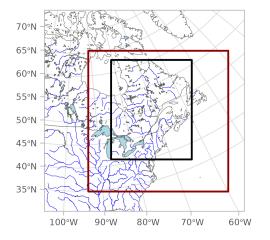
2 Synthetic data framework

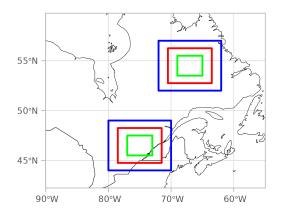
Data generated from a regional climate model (see § 2.1) are used exclusively in order to have total control of the data, eliminating common issues in station data—such as biases, missing observations, sparse coverage in certain regions, instrument failures, inconsistencies in historical records, and discontinuities in spatial and temporal coverage. More precisely, in the proposed synthetic data framework, we assumed that virtual gauged stations, i.e., where synthetic observations are available, correspond to a subset of the grid cells. The remaining cells are used as locations where to carry out the interpolation and thus serve to assess the performance. The graded case studies, see § 2.2, are organized according to different levels of complexity regarding station network density, region size and region location.

2.1 Regional climate model data: Study regions and periods

The ClimEx project investigated the impacts of climate change on extreme meteorological and hydrological events using the Canadian Regional Climate Model over two domains: one in Europe and one in NorthAmerica (Leduc et al., 2019). Among the available meteorological variables, we selected daily precipitation (converted from $kg/m^2/s$ to mm/day), minimum temperature, and maximum temperature (both converted from Kelvin (K) to degrees Celsius (°C)). From the period covered by the ClimEx project (1950–2099), we defined an auxiliary period (1980–2009) and an interpolation period (2000–2009). The auxiliary period is used to extract the auxiliary information—such as the climatology or spatial patterns (see § 3–4)— while the interpolation period is the period over which interpolation is performed. It is important to note that the auxiliary and interpolation periods are not required to overlap; however, overlapping periods do not pose a problem.

The North-American domain spans a regular grid of size 280×280 (78400 cells) with a grid cell size of approximately 11km (see the dark red rectangle in Fig. 1a). Within the subdomain outlined in black





(a) North American domain of the ClimEx project (dark red rectangle) and subdomain (black rectangle) considered in Fig. 1b.

(b) Two study regions—referred to as the south and north regions—within the subdomain outlined in black in Fig. 1a, shown at three different sizes (in different colours).

Figure 1 - Spatial domains used in this study.

in Fig. 1a) we selected two study regions, see Fig. 1b: one in the south and one in the north. Different nested sub-regions were created within each study region to enrich the experimental framework. In the following, we present statistical analyses comparing the northern and southern regions, focusing on the larger area (blue rectangles in Fig. 1b), to assess climatic differences during the auxiliary period (1980–2009).

The first notable difference between the two regions lies in their daily averages. A Student's t-test conducted at a 95% confidence level reveals that the average daily precipitation, minimum temperature, and maximum temperature are all significantly higher in the southern region compared to the northern region. Additionally, climatologies—representing long-term seasonal mean values— are shown in Fig 2. While confirming that average values are higher in the south than in the north, the precipitation climatologies also indicate distinct seasonal behaviours between the two regions.

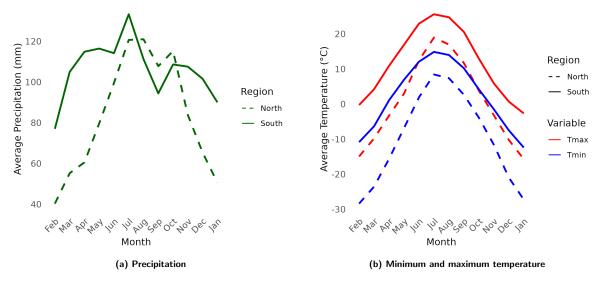


Figure 2 – Climatology of each meteorological variable in the larger study regions (blue rectangles in Fig. 1b) during the auxiliary period (1980–2009).

Moreover, the two regions exhibit different spatial variability. This variability is first assessed by comparing the number of components required in a Principal Component Analysis (PCA) decomposition to capture 90% of the total variance of a given meteorological variable. A greater number of components indicates a higher level of spatial complexity in the data. For precipitation, at least 50 components were needed to capture 90% of the variance in the southern region, compared to only 27 in the northern region, indicating greater spatial complexity in the south. In contrast, for both minimum and maximum temperatures, a single component was sufficient in each region. Furthermore, spatial variability is assessed by comparing semivariance plots (see Fig. 3). Semivariance measures how dissimilar variable values become with distance, with low values indicating strong similarity between nearby points, and high values reflecting greater dissimilarity over larger distances. The plots show that the semivariance in the southern region is significantly higher than in the northern region, with the difference becoming more pronounced as the distance increases. This suggests that spatial variability is greater in the south.

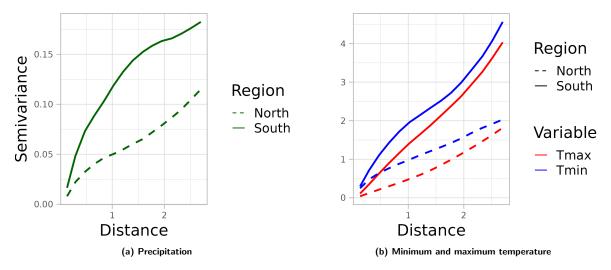


Figure 3 – Semivariance plots for each meteorological variable in the larger study regions (blue rectangles in Fig. 1b) during the auxiliary period (1980–2009).

2.2 Graded case studies

We designed various case studies with controlled complexity for each of the two selected regions (see Fig. 1b) by varying two factors: the size of the region and the density of the virtual gauged station network. The region size can take one of three values—large, medium, or small (corresponding to the nested rectangles in Fig. 1b)—while the network density can take one of five values: 10%, 30%, 50%, 70%, or 90%. In total, 30 case studies were designed, varying by region location, size, and network density (see Table 1).

3 Baseline interpolation methods

Common spatial interpolation methods from the literature serve as baselines to evaluate Spatial Pattern Regression (SPR), the newly proposed method in this work which is described in § 4. Let Z denote the random field of the meteorological variable of interest such that Z(s) represents the random variable at location s. The baseline interpolation methods considered operate by computing a weighted average of the observed values at nearby stations, with the general formula given as follows (Li and

Heap, 2014; Bokke, 2017):

$$\hat{z}_0 = \sum_{i=1}^n w_i z_i \tag{1}$$

with \hat{z}_0 the estimated value of $Z(s_0)$, the random field at the point of interest s_0 ; w_i the weight assigned to point s_i ; z_i the observed value of $Z(s_i)$ at point s_i ; s_i the geographic coordinates of location i (e.g., $s_i = (s_i^{\text{lon}}, s_i^{\text{lat}})$, the longitude and latitude); n the number of neighbouring points used to make the estimate.

Table 1 – Graded case study framework : For each region location (south or north), three region sizes (S, M, or L) and five network density levels (from 90% to 10%) are considered. The number of missing grid cells—where interpolation must be performed—is given by : (100 - density) × total number of grid cells.

	Region South			Region North			
	Size L	Size M	Size S	Size L	Size M	Size S	
Total # of grid cells	3025	1332	342	2970	1332	342	
Network density	# of missing grid cells						
90%	303	133	34	297	133	34	
70%	908	400	103	891	400	103	
50%	1523	666	171	1485	666	171	
30%	2118	932	239	2079	932	239	
10%	2723	1199	308	2673	1199	308	

3.1 Inverse Distance Weighting (IDW)

Inverse Distance Weighting (IDW) is a deterministic and univariate method (Hartkamp et al., 1999; Li and Heap, 2008; Tan and Xu, 2014). It involves weighting the values of neighboring observation stations by the inverse of the distance from the point of interest to these stations (Li and Heap, 2011, 2014; Sokolchuk and Sokac, 2022; Amin Burhanuddin et al., 2015; Pavão et al., 2012; Zimmerman et al., 1999). The closer a station is to the point of interest, the more influence its value will have on the interpolation result (Margaritidis, 2024; Li and Heap, 2008; Bokke, 2017). The formula used for the weights is:

$$w_i = \frac{\frac{1}{d_i^p}}{\sum_{i=1}^n \frac{1}{d_i^p}} \tag{2}$$

with $d_i = \sqrt{(s_i^{\text{lon}} - s_0^{\text{lon}})^2 + (s_i^{\text{lat}} - s_0^{\text{lat}})^2}$ the distance between the point of interest s_0 and station i where $p \geq 1$ controls the relative importance of the distance; n the number of stations in the neighbourhood; W_i the weight associated with the value of each station i.

The value of the power parameter p can have a significant impact on the interpolation results (Hartkamp et al., 1999). These weights, whose sum equals one, are used in IDW interpolation with the general formula (see (1)).

3.2 Ordinary Kriging (OK)

Ordinary Kriging (OK) is the most widely used method within the broader family of Kriging methods (Li and Heap, 2008; Margaritidis, 2024; Snepvangers et al., 2003). It assumes that the spatial correlation between observation points can explain the variability across the entire surface (Sokolchuk and Sokac, 2022). It is a statistical (or geostatistical) method capable of providing a measure of uncertainty related to predictions (Bokke, 2017; Pavão et al., 2012). OK estimates are optimal, unbiased, and have minimum variance. The OK model is

$$Z(s_0) = \mu + \epsilon(s_0) \tag{3}$$

with μ a constant (or global average) and $\epsilon(s_0)$ a stochastic residual with mean zero and spatial autocorrelation modelled using a so-called variogram function (e.g., exponential), which describes how the correlation between observation points decreases as the distance between them increases. The kriging weighting w_i used in (1) are determined by solving a system of linear equations which involves the variogram function. KO assumes stationary autocorrelation, meaning that the amount of correlation depends only on the distance between points, and isotropy, meaning that the autocorrelation behavior is the same in all directions.

3.3 Kriging with External Drift (KED)

Kriging with External Drift (KED) is a version of Kriging where the external drift introduces a trend represented as a linear function of only one auxiliary variable (Hartkamp et al., 1999; Varentsov et al., 2020). The auxiliary variable—in our implementation, we use the climatology—is any variable that is thought to have an influence on the primary variable. The KED model is:

$$Z(s_0) = \mu(s_0) + \epsilon(s_0) \tag{4}$$

with $\mu(s_0) = a + b \ X(s_0)$ the external drift, a function of the auxiliary variable X at point s_0 ; $\epsilon(s_0)$ a stochastic residual with mean zero and spatial auto-correlation modelled using a variogram function as in OK. The prediction formula of KED (Hengl et al., 2003) is equivalent to that of OK, except that the weights W_i used in (1) are determined by taking into account the external drift. As in OK, KED assumes stationary autocorrelation and isotropy.

4 Spatial pattern regression (SPR)

We present the methodology of Spatial Pattern Regression (SPR), the method newly proposed in this work. SPR relies on a representative basis of spatial patterns defined over a high-resolution grid covering the study region. A spatial pattern of a given variable is a recurrent spatial organization of that variable (Carreau and Guinot, 2021). The underlying rationale of SPR is that each field of a given meteorological variable on a given day can be expressed as a function of the spatial patterns in the basis. The first step (see Step 1 in Fig. 4 and § 4.1) thus consists in extracting spatial patterns, which, in our implementation, correspond to the eigenvectors derived from applying PCA to the gridded RCM data. These eigenvectors form a basis that captures the dominant spatial variability present in the data. The second step (see Step 2 in Fig. 4 and § 4.2) consists in integrating temporal information (for a given day, in our case) from values at the observation sites with spatial information captured by the spatial patterns. In our implementation, Step 2 produces a complete field over the study region by performing multiple linear regression to determine the linear combination of the spatial patterns that best matches the observations for that day.

4.1 Identification of a representative basis of spatial patterns

The basis of spatial patterns is identified using PCA, computed via Singular Value Decomposition (SVD) (Link et al., 2019). In what follows, boldface capital letters are used to denote matrices. Let $\mathbf{Z}_{n,p}^{\mathrm{grid}}$ denote the gridded RCM data, where n is the number of time steps (days) in the auxiliary period and p is the number of grid cells in the region considered. Applying SVD to $\mathbf{Z}_{n,p}^{\mathrm{grid}}$ with k singular vectors, we obtain the following decomposition:

$$\mathbf{Z}_{n,p}^{\text{grid}} = \mathbf{U}_{n,k}^{\text{grid}} \; \mathbf{S}_{k,k}^{\text{grid}} \; (\mathbf{V}_{p,k}^{\text{grid}})^T \; + (\bar{Z}_p^{\text{grid}} \; \mathbf{1}_n^T)^T$$
 (5)

where $\mathbf{U}_{n,k}^{\mathrm{grid}}$ is the matrix of left singular vectors, $\mathbf{S}_{k,k}^{\mathrm{grid}}$ is the diagonal matrix of singular values, $\mathbf{V}_{p,k}^{\mathrm{grid}}$ is the matrix of right singular vectors (corresponding to the eigenvectors in PCA), T denotes the transpose operation, 1_n is a vector of ones of length n and $\bar{Z}_p^{\mathrm{grid}}$ is a vector of length p that contains the columnwise average of $\mathbf{Z}_{n,p}^{\mathrm{grid}}$. The basis of spatial patterns, $\mathbf{V}_{p,k}^{\mathrm{grid}}$, are used (see § 4.2) to reconstruct a complete field over the study region by leveraging the observations at gauged stations.

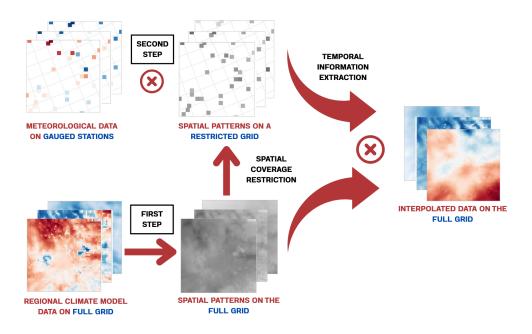


Figure 4 – SPR framework illustration: Step 1 extracts spatial patterns from high-resolution climate model data. Step 2 fits a linear model between observations at gauged stations and the corresponding restricted spatial patterns. The coefficients of the linear model thus contain the appropriate temporal information to reconstruct complete gridded fields using the full spatial patterns.

4.2 Spatio-temporal integration

The purpose of the regression model is to accurately reproduce the observed values at the gauged locations using a linear combination of k spatial patterns. To this end, only the grid cells corresponding to the virtual gauged stations are retained, meaning the spatial patterns are restricted to this reduced set of grid cells. Let $Z_d = (Z(s_1), \ldots, Z(s_d))^T$ denote the vector of random variables representing the meteorological variable at the d gauged stations for a given day. The spatial extent of the spatial patterns is restricted to correspond precisely to the locations of these stations. The spatial patterns $\mathbf{V}_{p,k}^{\mathrm{grid}}$, originally of dimensions $p \times k$, are thus restricted to spatial patterns $\mathbf{V}_{d,k}^{\mathrm{grid}}$, with dimensions $d \times k$, where d < p.

We assume the following linear model to characterize the relationship between the vector of random variables Z_d —centered using $\bar{Z}_d^{\rm grid}$, the mean values from the RCM data—and the restricted spatial patterns:

$$Z_d - \bar{Z}_d^{\text{grid}} = \mathbf{V}_{d,k}^{\text{grid}} \ \beta_k + \epsilon_d \tag{6}$$

with β_k the vector of spatial patterns coefficients of length k and ϵ_d the vector of random errors of length d for the day considered. This model is inspired by the SVD/PCA decomposition where, if $Z_d = Z_p$ for a fixed day $1 \le i \le n$, it corresponds to the i^{th} row of $\mathbf{Z}_{n,p}^{\text{grid}}$. In this setting, $\beta_k = U_{i,k}^{\text{grid}}$ $\mathbf{S}_{k,k}^{\text{grid}}$ as in Eq. (5). The model ensures that the reconstructed values reflect both the temporal variability observed at the stations, captured by the coefficients β_k , and the spatial structure from the RCM, represented by the restricted patterns $\mathbf{V}_{d,k}^{\text{grid}}$, while remaining computationally efficient and interpretable. The coefficients are estimated independently for each day of the interpolation period, yielding a sequence of regression vectors β_k . These coefficients capture local temporal information by describing how the amplitudes of the spatial patterns evolve from day to day. To reconstruct the complete field over the RCM grid for each day, the coefficients β_k are applied to the full basis of spatial patterns $\mathbf{V}_{p,k}^{\text{grid}}$. In other words, Z_p , the vector representing the complete field over the RCM grid for the same day as Z_d

in Eq. (6) is estimated as follows:

$$\hat{Z}_p = \mathbf{V}_{p,k}^{\text{grid}} \hat{\beta}_k + \bar{Z}_p^{\text{grid}}, \tag{7}$$

where $\hat{\beta}_k$ is the vector of fitted regression coefficients from Eq. (6).

5 Model evaluation and hyperparameters selection

For each meteorological variable, interpolation was performed independently for each day and then aggregated over the interpolation period. For precipitation, a transformation of the form $\log(\exp(\cdot)-1)$ was applied prior to interpolation to enforce positivity. To ensure the transformation remained valid, precipitation values were first bounded below by 10^{-5} . After interpolation, the results were backtransformed to the original units using $(\log(\exp(\cdot)-1))$. Additionally, following Werner et al. (2019), interpolated precipitation values below 0.5 after being back-transformed were set to zero; otherwise, the interpolated values were retained.

Since interpolation is performed independently for each day, the training, validation, and test sets are defined in terms of grid cells. The training set consists of the virtual gauged stations, representing a proportion of the total grid cells equal to the network density. The test set includes the remaining grid cells, representing N = 100-density % of the total, see Table 1. A validation set is also defined by randomly selecting N% of the training grid cells. Using the same percentage N% for both the validation and test sets allows for a consistent and comparable assessment of the interpolation models' generalization capacity on each set.

5.1 Hyperparameters selection

For each interpolation method, hyperparameters are selected to maximize performance on the validation set, which is withheld from the training set. In principle, since interpolation is performed day by day, one could select a different set of optimal hyperparameters for each day. However, to simplify the procedure, we instead select hyperparameters based on global performance—i.e., the values that yield the best average performance on the validation set across all days in the interpolation period. While this may not be optimal for every individual day, it is a reasonable and pragmatic choice, as our goal is to identify models that perform well overall across the interpolation period.

For the baseline interpolation methods (OK, KED, and IDW; see § 3), the hyperparameters considered are as follows: for OK and KED, the variogram model (Gaussian, Spherical, Exponential); and for IDW, the weighting power (1, 2, 3, 4, 5). For the proposed method, SPR (see § 4), the number of spatial patterns (k in Eq. (5)) must be selected. Since the maximum number of spatial patterns typically corresponds to the number of grid cells, which varies across regions of different sizes (see Table 1), we define the hyperparameter as a percentage (10% to 90%) of this maximum to ensure a consistent search across all regions. See § 7 for detailed results on the selected hyperparameters.

5.2 Performance evaluation

Each interpolation method is retrained on the full training set using the optimal hyperparameters selected through the training-validation procedure described in 5.1. Performance is then evaluated by comparing the true values z_j with the corresponding interpolated values \hat{z}_j from each method, where j denotes a grid cell in the test set. Two metrics are used for performance evaluation. The first metric is the Root Mean Squared Error (RMSE), a standard measure in regression tasks in general, and in interpolation settings in particular. It is computed for each day in the interpolation period as:

$$\sqrt{\frac{1}{d'} \sum_{j=1}^{d'} (z_j - \hat{z}_j)^2},\tag{8}$$

where d'=p-d, the number of grid cells in the test set; and averaged globally over the interpolation period. The second metric is the Structural Similarity Index Measure (SSIM), commonly used to assess the similarity between a compressed and an original image by considering luminance, contrast, and structure (Wang et al., 2004; Mirbod et al., 2022; Falola et al., 2024). It can be extended to quantify differences in spatial structure between two gridded datasets (Wang et al., 2004). Higher SSIM values indicate better reconstruction, with 1 representing perfect similarity and -1 representing complete dissimilarity. SSIM is computed on interpolated values for each interpolation day as:

$$\frac{(2\mu_z\mu_{\hat{z}} + C_1)(2\sigma_{z\hat{z}} + C_2)}{(\mu_z^2 + \mu_{\hat{z}}^2 + C_1)(\sigma_z^2 + \sigma_{\hat{z}}^2 + C_2)},\tag{9}$$

where all quantities are computed over the test set grid cells indexed by $1 \leq j \leq d'$; specifically, μ_z and $\mu_{\hat{z}}$ are the mean values of the true values z_j and the interpolated values \hat{z}_j , respectively; σ_z^2 and $\sigma_{\hat{z}}^2$ are their respective variances; $\sigma_{z\hat{z}}$ is the covariance between z_j and \hat{z}_j ; and C_1 and C_2 are small constants used for numerical stabilization.

6 Results

We first present an overall comparison of the three baseline methods with SPR across all graded case studies (see § 6.1). We then place the interpolation methods KED—the best-performing baseline method—and SPR under greater challenge by considering a realistic stress-test case study with a network density of 0.1% (see § 6.2). In § 6.3, we assess the effect of individual factors—namely, network density, region size, and region location (see Table 1)—on the interpolation performance of SPR.

6.1 Overall comparison

The three baseline methods and SPR are assessed across all graded case studies by computing daily RMSE (Eq. (8)) and SSIM (Eq. (9)) over the test grid cells. These daily scores are then averaged over the interpolation period (see Fig. 5 for the southern region and Fig. 6 for the northern region). Each panel in the figures corresponds to a specific variable and region size, with each symbol representing the performance of a given interpolation method at a specific network density. The x-axis shows RMSE, and the y-axis shows 1-SSIM; thus, symbols closer to the origin (0,0) indicate better performance.

The proposed method, SPR, shows strong and consistent performance across all three variables — precipitation, minimum temperature, and maximum temperature—in the majority of case studies. In the southern region, for both medium and large sizes (Fig. 5), SPR achieves the best results, with the lowest RMSE and highest SSIM for all three variables. However, KED and OK slightly outperform it for precipitation in the large region at 10% station density. In the small southern region, SPR generally ranks first, although KED and OK perform slightly better for precipitation at 30% density. For minimum and maximum temperatures in the small region with 10% density, SPR and KED perform similarly.

Across methods, OK and IDW are consistently ranked third and fourth in terms of performance. As in the southern region, SPR remains the top-performing method in most northern region case studies (Fig. 6), followed by KED, while OK and IDW trail behind.

From a station density of 50% and above, SPR clearly outperforms all baseline methods, regardless of region size or geographic location. At lower densities (<50%), the competition is mainly between SPR and KED, with similar trends observed in both regions. KED performs slightly better in specific cases: 10% density in large southern regions, 30% density in medium and small southern regions for precipitation, and lowest densities for minimum and maximum temperatures in small southern regions. In some of these situations, KED achieves the lowest RMSE, while SPR maintains better SSIM—for example, in the large northern region for precipitation at 10% density.

This trade-off between RMSE and SSIM generally favors SPR, which tends to better preserve the structural patterns of the virtual observations across most scenarios. Overall, out of 90 synthetic configurations, KED outperforms SPR in just six cases and matches its performance in three.

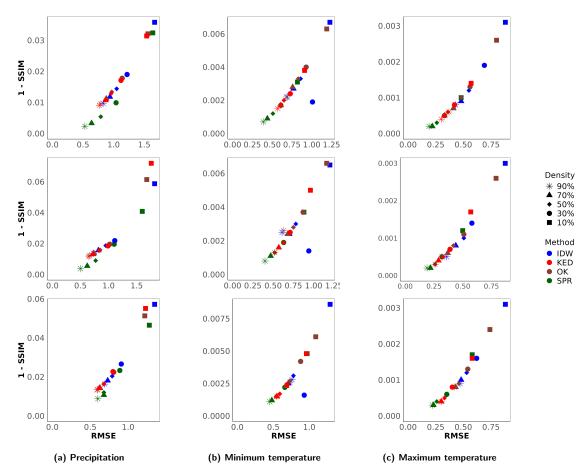


Figure 5 – Comparison of the three baseline methods and SPR in the southern region, in terms of averaged RMSE (x-axis) and averaged 1-SSIM (y-axis). Each column represents a specific meteorological variable, while each row corresponds to a region size—ranging from the largest at the top to the smallest at the bottom. Each color represents a different interpolation method, and each plotting symbol corresponds to a specific network density. The closer a symbol is to the origin, the better the performance.

6.2 Realistic stress-test case study

We designed the following stress-test case study to reflect the main challenge faced by spatial interpolation: the often very low density of station networks. Indeed, especially in remote regions, the number of gauged stations tends to be very low across vast areas. To evaluate the performance of KED and SPR under these typical practical conditions, we selected the larger region in northern Quebec (see Fig. 1b), consisting of 2,970 grid cells. Only three of these grid cells—approximately 0.1% of the total—were randomly selected to serve as virtual stations. The hyperparameters—the variogram model for KED and the number of eigenvectors for SPR—are set to the same values as those selected for the region of the same size with the lowest station network density. The results, reported in Table 2, show that SPR has the lowest average and median RMSE for all variables. It also achieves the highest SSIM for minimum and maximum temperature, except for precipitation, where KED yields a higher SSIM.

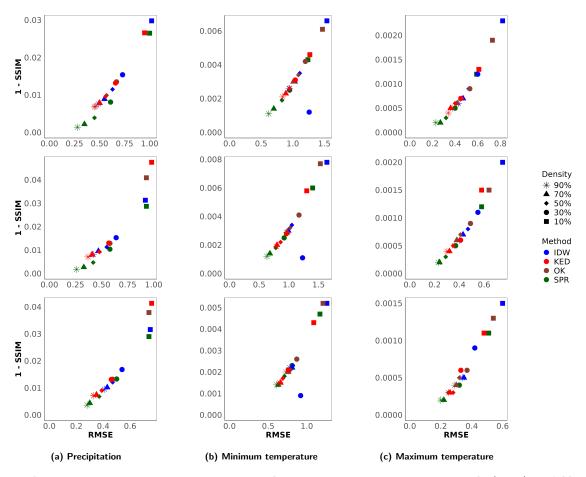


Figure 6 – Comparison of the three baseline methods and SPR in the northern region, in terms of RMSE (x-axis) and 1-SSIM (y-axis). Each column represents a specific meteorological variable, while each row corresponds to a region size—ranging from the largest at the top to the smallest at the bottom. Each color represents a different interpolation method, and each plotting symbol corresponds to a specific network density. The closer a symbol is to the origin, the better the performance.

Table 2 – Realistic stress-case study: daily RMSE statistics (mean, median, standard deviation (SD) and 2.5% and 97.5% quantiles) and average SSIM values for each variable and interpolation method. Lower RMSE and higher SSIM values indicate better performance. The best values for average and median RMSE, as well as SSIM, are shown in bold.

Variable	Method	Mean	I Median	RMSE SD	2.5%	97.5%	SSIM
PR	SPR KED	3.18 3.39	1.97 2.16	3.63 4.06	0.01 0.04	13.17 14.02	0.6654 0.8044
TMIN	SPR KED	3.53 3.99	3.19 3.54	1.8 2.10	1.08 1.24	7.79 9.00	0.9636 0.9541
TMAX	SPR KED	2.59 2.97	2.29 2.59	1.39 1.61	0.86 1.01	6.10 7.09	0.9757 0.9687

To illustrate the differences in the internal mechanisms of SPR and KED, complete interpolated fields along with the corresponding ground truth (i.e., synthetic observations) are shown for three different days for each of the following variables: precipitation, minimum temperature, and maximum temperature (see Fig. 7). In addition, the corresponding spatial RMSE is presented in Fig. 8. KED interpolated fields are more strongly correlated with the climatology than SPR's (average Spearman correlation: 0.9997 for KED and 0.7910 for SPR) across the three variables on the three different days considered in Fig. 7 and Fig. 8. However, SPR interpolated fields are structurally more similar to the observations (average SSIM: 0.4151 for SPR and 0.2504 for KED). This suggests that, when

observational information is limited, KED relies heavily on the auxiliary data provided as external drift. In contrast, SPR is able to generate fields that resemble the climatology, even though the climatology is not explicitly used as input.

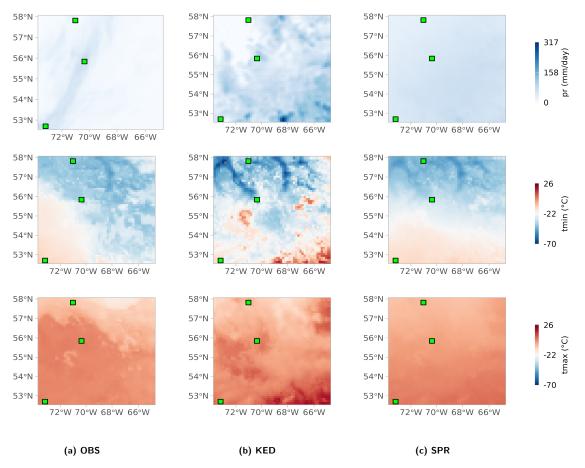


Figure 7 – Realistic stress-case study: interpolated fields of SPR and KED over the larger north region. The observed fields (ground truth) are also included. Green squares indicate the virtual station locations. Solid squares indicate that KED does not interpolate at virtual station locations. Rows correspond to precipitation, minimum temperature, and maximum temperature (top to bottom) on three different days.

6.3 Sensitivity of SPR to individual factors

A factor-wise assessment of SPR is conducted across all three meteorological variables (see Table 1 for the list of factors and their corresponding values).

6.3.1 Effect of station network density

We investigate the influence of network density on the performance of SPR in terms of RMSE, averaged over all days in the interpolation period. In addition, a 95% confidence band is reported, computed as the 2.5% and 97.5% quantiles of the daily RMSE values (see Fig. 9 for the region in the south and Fig. 10 for the region in the north).

We notice that increasing the density of the station network leads to a decrease in average RMSE in both regions, as expected. Beyond 50% density, further increases in network density lead to only marginal improvements in error reduction, regardless of region size or region location. Thus, the higher the station network density, the lower the error and the narrower the confidence interval. Looking at differences across the meteorological variables, we find the following. The average RMSE for precipitation

is higher than that for minimum and maximum temperatures. Additionally, the confidence interval for precipitation is wider, indicating greater uncertainty in interpolating precipitation compared to temperature variables.

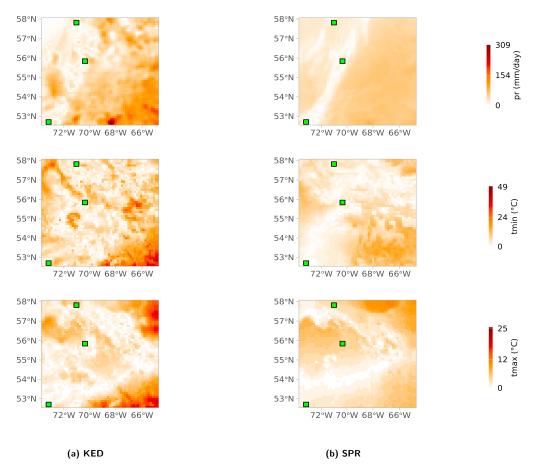


Figure 8 – Realistic stress-case study: comparison of SPR and KED in the larger northern region, based on their spatial RMSE. Each column represents a specific method (SPR or KED), while rows represent a specific variable: precipitation, minimum temperature, temperature, from top to bottom. Green squares indicate the virtual station locations. Solid squares indicate that KED does not interpolate at virtual station locations. Rows correspond to precipitation, minimum temperature, and maximum temperature (top to bottom), shown on the same three days as in Fig. 7.

6.3.2 Effect of region size

We conduct a similar investigation into the influence of the region size on the performance of SPR, based on RMSE averages and confidence bands (see Fig. 11 for the southern region and Fig. 12 for the northern region). Interestingly, the interpolation performance does not decrease steadily with region size—a somewhat surprising result. For each network density and meteorological variable (precipitation, minimum and maximum temperatures), the average RMSE remains relatively stable across region sizes in both study areas. In some cases, larger regions show slightly higher RMSE than smaller ones at the same network density, while smaller regions may exhibit wider confidence bands. However, the confidence interval becomes wider as the region size decreases. These findings suggest that region size has limited impact on the average interpolation performance of SPR.

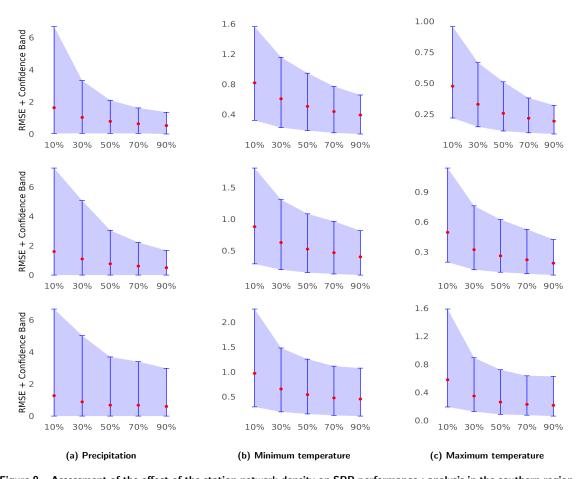


Figure 9 – Assessment of the effect of the station network density on SPR performance: analysis in the southern region, in terms of averaged RMSE (red dot) and 95% confidence band (in blue). Each column represents a specific meteorological variable, while each row corresponds to a region size—ranging from the largest at the top to the smallest at the bottom. For each density (on the x-axis) and variable, a lower RMSE combined with a narrower confidence band indicates better performance.

6.3.3 Effect of region location

Finally, we examine the influence of region location on the performance of SPR, using boxplots of the daily RMSE (see Fig. 13). Indeed, as noted in § 2.1, the northern and southern regions exhibit distinct climatic characteristics, both in terms of magnitude and spatial variability. For each variable and station network density, the distribution of daily average RMSE is different depending on whether it is precipitation or temperature (see Fig. 13). For precipitation, the average RMSE is lower in the north than in the south. This suggests that higher spatial variability—reflected by the need for more PCA components to explain at least 90% of the variance—leads to higher interpolation errors. In contrast, for minimum and maximum temperatures, the average RMSE is higher in the north, despite both regions requiring only one PCA component to reach 90% of the variance. In this case, the difference is likely due to the shape of the spatial variability, which differs between regions, as shown in the semivariance plots (see § 2.1).

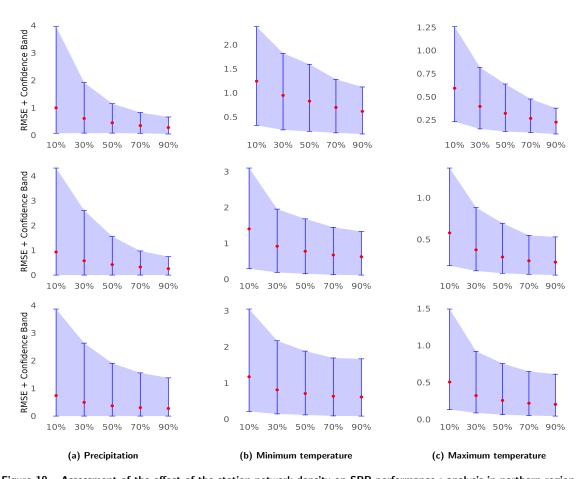


Figure 10 – Assessment of the effect of the station network density on SPR performance: analysis in northern region, in terms of averaged RMSE (red dot) and 95% confidence band (in blue). Each column represents a specific meteorological variable, while each row corresponds to a region size—ranging from the largest at the top to the smallest at the bottom. For each density (on the x-axis) and variable, a lower RMSE combined with a narrower confidence band indicates better performance.

7 Discussion and conclusion

In this work, we introduced Spatial Pattern Regression (SPR), a new method for interpolating meteorological data by leveraging the spatial structure of regional climate model (RCM) simulations. SPR first extracts spatial patterns from these simulations as the eigenvectors obtained from an SVD/PCA decomposition, modeling each meteorological field as a linear combination of these patterns. It then uses multiple linear regression to estimate the coefficients of this linear combination—i.e., the weights that best reconstruct the observed values at gauged stations for a given day. The effectiveness of SPR depends on the representativeness of the extracted spatial patterns and the accuracy of the temporal information captured through regression. Our results show that SPR delivers accurate estimates and outperforms baseline interpolation methods (KED, OK, and IDW) in the majority of the graded case studies within our synthetic data framework. Among the baseline methods, however, KED remains the strongest, consistent with previous findings by Bishop and McBratney (2001) in the context of soil property mapping. Our analysis confirmed that station density and region location significantly influence interpolation performance—findings supported by previous studies such as Stahl et al. (2006); Li and Heap (2014, 2008) and Wagner et al. (2012).

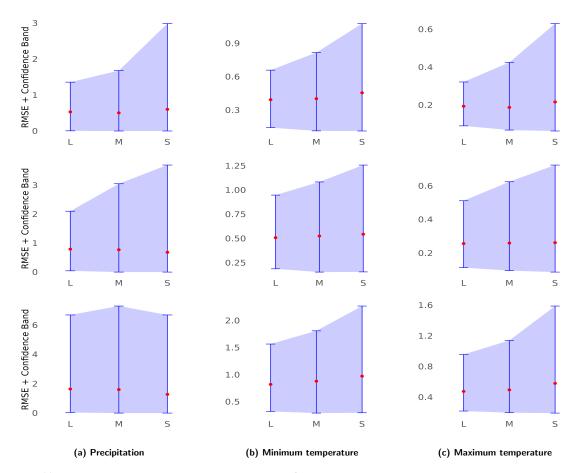


Figure 11 – Assessment of the effect of the region size on SPR performance: analysis in the southern region, in terms of averaged RMSE (red dot) and 95% confidence band (in blue). Each column represents a specific meteorological variable, while rows represent increasing network densities: 90%, 50% and 10% from top to bottom. For each region size (on the x-axis) and variable, a lower RMSE combined with a narrower confidence band indicates better performance.

An interesting feature of SPR lies in the fact that it seeks to reproduce the spatial structure observed in RCM simulations within the interpolated fields. In climate change impact studies—such as in hydrology—interpolated fields are typically used for calibrating models over historical periods, while future RCM simulations are used to assess climate change. By aligning the spatial structures of past interpolated fields with those of future simulations, SPR offers improved consistency between historical calibration and future projection. Unlike traditional interpolation methods, which may incorporate auxiliary information such as elevation or RCM-derived climatology in a somewhat ad hoc manner, SPR offers a systematic and principled approach to integrating such information. The auxiliary period which serves to extract the spatial patterns must be such that key meteorological events and the variability relevant to the targeted meteorological variable are well represented. Apart from this requirement, the auxiliary period can be chosen flexibly and may or may not overlap with the interpolation period when observations are available. Compared to reanalysis products, SPR is simpler to implement for several reasons (Gasset et al., 2021). First, as previously mentioned, the auxiliary RCM data used to extract spatial patterns does not need to cover the same period as the observations. Second, SPR relies on basic statistical tools—SVD/PCA and linear regression—without requiring complex data assimilation procedures. Finally, SPR can incorporate other types of auxiliary gridded data that capture spatial structure, such as elevation, radar, or remote sensing data.

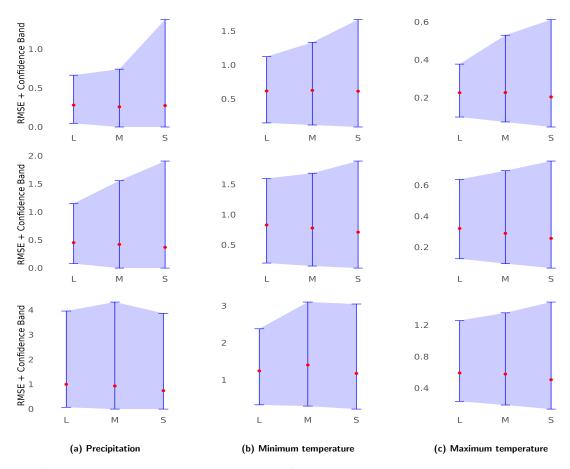


Figure 12 – Assessment of the effect of the region size on SPR performance: analysis in the northern region, in terms of averaged RMSE (red dot) and 95% confidence band (in blue). Each column represents a specific meteorological variable, while rows represent increasing network densities: 90%, 50% and 10% from top to bottom. For each region size (on the x-axis) and variable, a lower RMSE combined with a narrower confidence band indicates better performance.

There are several possible avenues for improving SPR. One potential enhancement is to move beyond selecting a fixed number of leading spatial patterns and instead allow, for each time step, the selection of any subset—not necessarily the top n% based on explained variance. This added flexibility could help capture distinctive features that may appear in lower-ranked eigenvectors, which are often overlooked in the current approach. Additionally, SPR currently fits one regression per day independently, ignoring the temporal continuity of the patterns. A global modeling approach that captures temporal dependencies—given the relatively stable nature of spatial patterns—could improve performance, as suggested by Amato et al. (2020). In summary, SPR offers a promising and efficient alternative for spatial interpolation, with demonstrated accuracy. The current version will serve as a baseline for future developments aimed at incorporating non-linear and temporal dependencies.

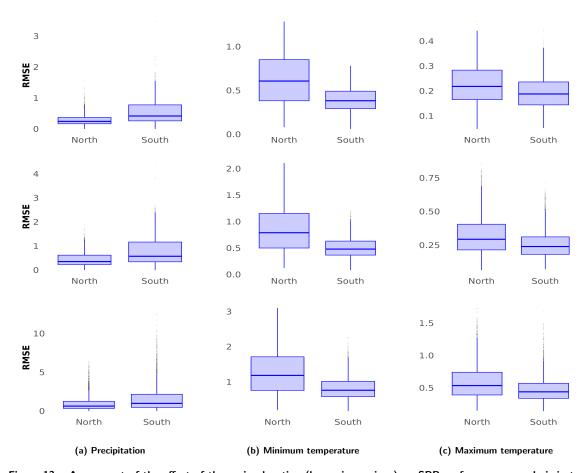


Figure 13 – Assessment of the effect of the region location (large size regions) on SPR performance : analysis in terms of boxplots of daily RMSE. Each column represents a specific meteorological variable, while rows represent increasing network densities : 90%, 50% and 10% from top to bottom.

Software and data availability

Software used: All analyses were conducted using the R programming language (version 4.4.3) within Visual Studio Code (VSCode).

Key R packages: bigmemory, bigstatsr, doParallel, foreach, gstat, sf, sp, terra.

Code availability: The analysis scripts are not publicly shared, but can be made available upon reasonable request to the corresponding author.

Program language: R.

Hardware requirements: PC/Mac/Linux with at least 8 GB RAM recommended for full dataset processing.

Primary data source: Daily climate simulations from the ClimEx project (https://climex-data.srv.lrz.de/Public/). We used one ensemble member (kcj) from the CanESM2-driven simulations, including variables: precipitation, minimum temperature, and maximum temperature.

Data license: Usage of ClimEx data is subject to their terms of use as described on the ClimEx portal.

Author CRediT statement

Vihotogbé Houssou: Conceptualization, Methodology, Software, Validation, Formal analysis, Data curation, Visualization, Writing - original draft, Writing - review & editing.

Julie Carreau: Conceptualization, Methodology, Supervision, Resources, Funding acquisition, Validation, Writing - review & editing.

Appendix: Selected hyperparameters

Table 3 – SPR Optimal proportion of spatial patterns selected for each variable, region (South and North), and region size (L, M, S), at different network densities.

Density	Variable	F	Region South			Region North		
		Size L	Size M	Size S	Size L	Size M	Size S	
90%	Precipitation Minimum Temp. Maximum Temp.	0.45 0.30 0.40	0.45 0.30 0.35	0.45 0.30 0.35	0.50 0.30 0.40	0.60 0.30 0.35	0.50 0.25 0.35	
70%	Precipitation Minimum Temp. Maximum Temp.	0.45 0.30 0.40	0.40 0.25 0.30	0.45 0.30 0.35	0.50 0.30 0.40	0.60 0.30 0.40	0.45 0.30 0.40	
50%	Precipitation Minimum Temp. Maximum Temp.	0.45 0.25 0.30	0.40 0.30 0.30	0.40 0.30 0.40	0.45 0.20 0.30	0.45 0.25 0.30	0.40 0.25 0.30	
30%	Precipitation Minimum Temp. Maximum Temp.	0.40 0.30 0.30	0.20 0.25 0.30	0.30 0.25 0.25	0.40 0.25 0.30	0.30 0.25 0.25	0.30 0.25 0.25	
10%	Precipitation Minimum Temp. Maximum Temp.	0.25 0.40 0.35	0.40 0.45 0.40	0.50 0.50 0.50	0.30 0.30 0.30	0.40 0.50 0.50	0.50 0.50 0.50	

Table 4 – IDW Optimal weighting powers selected for each variable, region (South and North), and region size (L, M, S), at different network densities.

Density	Variable	F	Region Sout	h	Region North			
		Size L	Size M	Size S	Size L	Size M	Size S	
	Precipitation	5	5	5	5	5	5	
90%	Minimum Temp.	5	5	5	5	5	5	
	Maximum Temp.	5	5	5	5	5	5	
	Precipitation	5	5	5	5	5	5	
70%	Minimum Temp.	5	4	5	5	4	5	
	Maximum Temp.	5	4	5	5	5	5	
	Precipitation	4	4	4	4	4	4	
50%	Minimum Temp.	4	4	4	4	4	4	
	Maximum Temp.	4	4	4	4	4	4	
	Precipitation	4	4	4	4	4	4	
30%	Minimum Temp.	3	3	3	3	3	3	
	Maximum Temp.	3	4	4	4	4	4	
	Precipitation	3	2	2	4	3	3	
10%	Minimum Temp.	3	3	2	4	2	2	
	Maximum Temp.	3	3	4	5	3	3	

Table 5 – OK Optimal kriging models selected for each variable, region (South and North), and region size (L, M, S), at different network densities.

Density	Variable	F	Region Sout	h	Region North			
		Size L	Size M	Size S	Size L	Size M	Size S	
90%	Precipitation Minimum Temp. Maximum Temp.	Sph Sph Exp	Sph Sph Sph	Sph Sph Sph	Sph Sph Sph	Sph Sph Sph	Sph Sph Sph	
70%	Precipitation Minimum Temp. Maximum Temp.	Sph Sph Exp	Sph Sph Sph	Sph Sph Sph	Sph Sph Sph	Sph Sph Sph	Sph Sph Sph	
50%	Precipitation Minimum Temp. Maximum Temp.	Sph Sph Sph	Sph Sph Sph	Sph Sph Sph	Sph Sph Sph	Sph Sph Sph	Sph Sph Sph	
30%	Precipitation Minimum Temp. Maximum Temp.	Sph Sph Sph	Sph Sph Sph	Sph Exp Exp	Sph Sph Sph	Sph Sph Sph	Sph Sph Exp	
10%	Precipitation Minimum Temp. Maximum Temp.	Sph Sph Sph	Exp Exp Exp	Exp Exp Exp	Sph Sph Exp	Exp Exp Exp	Exp Exp Exp	

Table 6 – KED Optimal kriging models selected for each variable, region (South and North), and region size (L, M, S), at different network densities.

Density	Variable	F	Region South			Region North		
		Size L	Size M	Size S	Size L	Size M	Size S	
90%	Precipitation	Sph	Sph	Sph	Sph	Sph	Sph	
	Minimum Temp.	Sph	Sph	Sph	Sph	Sph	Sph	
	Maximum Temp.	Sph	Sph	Sph	Sph	Sph	Sph	
70%	Precipitation	Sph	Sph	Sph	Sph	Sph	Sph	
	Minimum Temp.	Sph	Sph	Sph	Sph	Sph	Sph	
	Maximum Temp.	Sph	Sph	Sph	Sph	Sph	Sph	
50%	Precipitation	Sph	Sph	Sph	Sph	Sph	Sph	
	Minimum Temp.	Sph	Sph	Sph	Sph	Sph	Sph	
	Maximum Temp.	Sph	Sph	Sph	Sph	Sph	Sph	
30%	Precipitation	Sph	Sph	Sph	Sph	Sph	Sph	
	Minimum Temp.	Sph	Sph	Sph	Sph	Sph	Sph	
	Maximum Temp.	Sph	Sph	Exp	Sph	Sph	Exp	
10%	Precipitation	Sph	Exp	Exp	Sph	Exp	Exp	
	Minimum Temp.	Sph	Exp	Exp	Sph	Exp	Exp	
	Maximum Temp.	Sph	Exp	Exp	Sph	Exp	Exp	

Références

Amato, F., Guignard, F., Robert, S., Kanevski, M., 2020. A novel framework for spatio-temporal prediction of environmental data using deep learning. Scientifics Reports 10. doi:10.1038/s41598-020-79148-7.

Amin Burhanuddin, S.N.Z., Deni, S., Mohamed Ramli, N., 2015. Geometric median for missing rainfall data imputation, in: The 2nd ISM International Statistical Conference 2014 (ISM–II): Empowering the Applications of Statistical and Mathematical Sciences, pp. 113–119. doi:10.1063/1.4907433.

Bishop, T., McBratney, A., 2001. A comparison of prediction methods for the creation of field-extent soil property maps. Geoderma 103, 149–160. URL: https://www.sciencedirect.com/science/article/pii/S001670610100074X, doi:https://doi.org/10.1016/S0016-7061(01)00074-X. estimating uncertainty in soil models.

Bokke, A., 2017. Comparative evaluation of spatial interpolation methods for estimation of missing meteorological variables over ethiopia. Journal of Water Resource and Protection 09, 945–959. doi:10.4236/jwarp.2017.98063.

- Carreau, J., Guinot, V., 2021. A pca spatial pattern based artificial neural network downscaling model for urban flood hazard assessment. Advances in Water Resources 147, 103821. URL: https://www.sciencedirect.com/science/article/pii/S0309170820307107, doi:https://doi.org/10.1016/j.advwatres.2020.103821.
- Faghih, M., Brissette, F., 2023. Temporal and spatial amplification of extreme rainfall and extreme floods in a warmer climate. Journal of Hydrometeorology 24, 1331-1347. URL: https://journals.ametsoc.org/view/journals/hydr/24/8/JHM-D-22-0224.1.xml, doi:10.1175/JHM-D-22-0224.1.
- Falola, Y., Churilova, P., Liu, R., Huang, C.K., Delgado, J.F., Misra, S., 2024. Generating extremely low-dimensional representation of subsurface earth models using vector quantization and deep autoencoder. Petroleum Research URL: https://www.sciencedirect.com/science/article/pii/S2096249524000619, doi:https://doi.org/10.1016/j.ptlrs.2024.07.001.
- Gasset, N., Fortin, V., Dimitrijevic, M., Carrera, M., Bilodeau, B., Muncaster, R., Gaborit, E., Roy, G., Pentcheva, N., Bulat, M., Wang, X., Pavlovic, R., Lespinas, F., Khedhaouiria, D., Mai, J., 2021. A 10 km north american precipitation and land-surface reanalysis based on the gem atmospheric model. Hydrology and Earth System Sciences 25, 4917–4945. URL: https://hess.copernicus.org/articles/25/4917/2021/, doi:10.5194/hess-25-4917-2021.
- Hartkamp, A., de Beurs, K., Stein, A., White, J., 1999. Interpolation techniques for climate variables. Geographic Information Systems Series 99–01. International Maize and Wheat Improvement Center (CIMMYT), Mexico 1999. ISSN: 1405–7484.
- Hengl, T., Heuvelink, G., Stein, A., 2003. Comparison of kriging with external drift and regression-kriging. Technical Note.
- Lauer, A., Pausata, F.S.R., Leroyer, S., Argueso, D., 2023. Effect of urban heat island mitigation strategies on precipitation and temperature in montreal, canada: Case studies. PLOS Climate 2, e0000196. URL: https://doi.org/10.1371/journal.pclm.0000196, doi:10.1371/journal.pclm.0000196.
- Leduc, M., Mailhot, A., Frigon, A., et al., 2019. The ClimEx project: A 50-member ensemble of climate change projections at 12-km resolution over europe and northeastern north america with the canadian regional climate model (CRCM5). J. of App. Meteo. & Clim. 58, 663–693. doi:https://doi.org/10.1175/JAMC-D-18-0021.1.
- Li, J., Heap, A.D., 2008. A review of spatial interpolation methods for environmental scientists. Geoscience Australia.
- Li, J., Heap, A.D., 2011. A review of comparative studies of spatial interpolation methods in environmental sciences: Performance and impact factors. Ecological Informatics 6, 228–241. doi:https://doi.org/10.1016/j.ecoinf.2010.12.003.
- Li, J., Heap, A.D., 2014. Spatial interpolation methods applied in the environmental sciences: A review. Environmental Modelling & Software 53, 173–189. doi:https://doi.org/10.1016/j.envsoft.2013.12.008.
- Link, R., Snyder, A., Lynch, C., Hartin, C., Kravitz, B., Bond-Lamberty, B., 2019. Fldgen v1.0: an emulator with internal variability and space—time correlation for earth system models. Geoscientific Model Development 12, 1477—1489. doi:10.5194/gmd-12-1477-2019.
- Livneh, B., Bohn, T.J., Pierce, D.W., Munoz-Arriola, F., Nijssen, B., Vose, R., Cayan, D.R., Brekke, L., 2015. A spatially comprehensive, hydrometeorological data set for mexico, the u.s., and southern canada 1950–2013. Scientific Data 2, 150042. URL: https://doi.org/10.1038/sdata.2015.42, doi:10.1038/sdata.2015.42.
- Lucas-Picher, P., Arsenault, R., Poulin, A., Ricard, S., Lachance-Cloutier, S., Turcotte, R., 2020. Application of a high-resolution distributed hydrological model on a U.S.-Canada transboundary basin: Simulation of the multiyear mean annualhydrograph and 2011 flood of therichelieu river basin. Journal of Advances in Modeling Earth Systems 12, e2019MS001709. URL: https://agupubs.onlinelibrary.wiley.com/doi/abs/10.1029/2019MS001709, doi:https://doi.org/10.1029/2019MS001709, arXiv:https://agupubs.onlinelibrary.wiley.com/doi/pdf/10.1029/2019MS001709. e2019MS001709 2019MS001709.
- Margaritidis, A., 2024. Comparison of spatial interpolation methods of precipitation data in central macedonia, greece. Computational Water, Energy, and Environmental Engineering 13, 13–37. doi:10.4236/cweee.2024.131002.
- Mirbod, M., Ghatari, A.R., Saati, S., Shoar, M., 2022. Industrial parts change recognition model using machine vision, image processing in the framework of industrial information integration. Journal of Industrial Information Integration 26, 100277. URL: https://www.sciencedirect.com/science/article/pii/S2452414X21000741, doi:https://doi.org/10.1016/j.jii.2021.100277.

- Pavão, C., França, G., Marotta, G., Mnezes, P.H.B., Neto, G., Roig, H., 2012. Spatial interpolation applied a crustal thickness in brazil. Journal of Geographic Information System 2151-1969 4, 142–152.
- Muñoz Sabater, J., Dutra, E., Agustí-Panareda, A., Albergel, C., Arduini, G., Balsamo, G., Boussetta, S., Choulga, M., Harrigan, S., Hersbach, H., Martens, B., Miralles, D.G., Piles, M., Rodríguez-Fernández, N.J., Zsoter, E., Buontempo, C., Thépaut, J.N., 2021. Era5-land: a state-of-the-art global reanalysis dataset for land applications. Earth System Science Data 13, 4349–4383. URL: https://essd.copernicus.org/articles/13/4349/2021/, doi:10.5194/essd-13-4349-2021.
- Snepvangers, J., Heuvelink, G., Huisman, J., 2003. Soil water content interpolation using spatio-temporal kriging with external drift. Geoderma 112, 253–271. doi:https://doi.org/10.1016/S0016-7061(02)00310-5. pedometrics 2001.
- Sokolchuk, K., Sokac, M., 2022. Comparison of spatial interpolation methods of hydrological data on example of the pripyat river basin (within ukraine). Acta Hydrologica Slovaca .
- Stahl, K., Moore, R., Floyer, J., Asplin, M., McKendry, I., 2006. Comparison of approaches for spatial interpolation of daily air temperature in a large region with complex topography and highly variable station density. Agricultural and Forest Meteorology 139, 224-236. URL: https://www.sciencedirect.com/science/article/pii/S0168192306001638, doi:https://doi.org/10.1016/j.agrformet.2006.07.004.
- Tan, Q., Xu, X., 2014. Comparative analysis of spatial interpolation methods: an experimental study. Sensors and Transducers 165, 155–163.
- Varentsov, M., Esau, I., Wolf, T., 2020. High-resolution temperature mapping by geostatistical kriging with external drift from large-eddy simulations. Monthly Weather Review 148, 1029–1048. URL: http://dx.doi.org/10.1175/MWR-D-19-0196.1, doi:10.1175/MWR-D-19-0196.1.
- Wagner, P.D., Fiener, P., Wilken, F., Kumar, S., Schneider, K., 2012. Comparison and evaluation of spatial interpolation schemes for daily rainfall in data scarce regions. Journal of Hydrology 464-465, 388-400. URL: https://www.sciencedirect.com/science/article/pii/S0022169412006270, doi:https://doi.org/10.1016/j.jhydrol.2012.07.026.
- Wang, Z., Bovik, A., Sheikh, H., Simoncelli, E., 2004. Image quality assessment: from error visibility to structural similarity. IEEE Transactions on Image Processing 13, 600–612. doi:10.1109/TIP.2003.819861.
- Warren, F., et al., 2022. Canada in a Changing Climate: Regional Perspectives Report. Government of Canada, Ottawa, ON. Accessed: 2025-04-07 at https://changingclimate.ca/regional-perspectives/.
- Werner, A., Schnorbus, M., Shrestha, R., Cannon, A., Zwiers, F., Dayon, G., Anslow, F., 2019. A long-term, temporally consistent, gridded daily meteorological dataset for northwestern north america. Scientific Data 6, 1–16. doi:https://doi.org/10.1038/sdata.2018.299.
- Zimmerman, D., Pavlik, C., Ruggles, A., Armstrong, M.P., 1999. An experimental comparison of ordinary and universal kriging and inverse distance weighting. Mathematical Geology, 375–390doi:10.1023/A:1007586507433.