



## **Data-Driven Estimation of Groundwater Level Time-Series at Unmonitored Sites Using Comparative Regional Analysis**

Downloaded from: <https://research.chalmers.se>, 2023-09-08 04:53 UTC

Citation for the original published paper (version of record):

Haaf, E., Giese, M., Reimann, T. et al (2023). Data-Driven Estimation of Groundwater Level Time-Series at Unmonitored Sites Using Comparative Regional Analysis. *Water Resources Research*, 59(7). <http://dx.doi.org/10.1029/2022WR033470>

N.B. When citing this work, cite the original published paper.

# Water Resources Research



## RESEARCH ARTICLE

10.1029/2022WR033470

## Data-Driven Estimation of Groundwater Level Time-Series at Unmonitored Sites Using Comparative Regional Analysis

E. Haaf<sup>1</sup> , M. Giese<sup>2</sup> , T. Reimann<sup>3</sup> , and R. Barthel<sup>2</sup> 

<sup>1</sup>Department of Architecture and Civil Engineering, Chalmers University of Technology, Gothenburg, Sweden, <sup>2</sup>Department of Earth Sciences, University of Gothenburg, Gothenburg, Sweden, <sup>3</sup>Institute for Groundwater Management, TU Dresden, Dresden, Germany

### Key Points:

- Daily groundwater levels at unmonitored sites are estimated through transfer of head duration curves based on the similarity of site characteristics at monitored sites
- Nonlinearity of controls on groundwater levels favors using of Machine Learning (e.g., regression trees) over multiple linear regression for prediction
- The dynamic nature of controls on groundwater levels can be disentangled, which is central for studies of recharge seasonality, droughts, and floods

### Supporting Information:

Supporting Information may be found in the online version of this article.

### Correspondence to:

E. Haaf,  
[ezra.haaf@chalmers.se](mailto:ezra.haaf@chalmers.se)

### Citation:

Haaf, E., Giese, M., Reimann, T., & Barthel, R. (2023). Data-driven estimation of groundwater level time-series at unmonitored sites using comparative regional analysis. *Water Resources Research*, 59, e2022WR033470. <https://doi.org/10.1029/2022WR033470>

Received 16 AUG 2022

Accepted 1 JUN 2023

### Author Contributions:

**Conceptualization:** E. Haaf, M. Giese  
**Data curation:** E. Haaf  
**Formal analysis:** E. Haaf  
**Funding acquisition:** E. Haaf, R. Barthel  
**Investigation:** E. Haaf  
**Methodology:** E. Haaf  
**Project Administration:** E. Haaf  
**Resources:** E. Haaf  
**Software:** E. Haaf  
**Supervision:** T. Reimann, R. Barthel  
**Validation:** E. Haaf, M. Giese  
**Visualization:** E. Haaf

© 2023. The Authors.

This is an open access article under the terms of the [Creative Commons Attribution License](https://creativecommons.org/licenses/by/4.0/), which permits use, distribution and reproduction in any medium, provided the original work is properly cited.

**Abstract** A new method is presented to efficiently estimate daily groundwater level time series at unmonitored sites by linking groundwater dynamics to local hydrogeological system controls. The proposed approach is based on the concept of comparative regional analysis, an approach widely used in surface water hydrology, but uncommon in hydrogeology. Using physiographic and climatic site descriptors, the method utilizes regression analysis to estimate cumulative frequency distributions of groundwater levels (groundwater head duration curves, HDC) at unmonitored locations. The HDC is then used to construct a groundwater hydrograph using time series from distance-weighted neighboring monitored (donor) locations. For estimating time series at unmonitored sites, in essence, spatio-temporal interpolation, stepwise multiple linear regression (MLR), extreme gradient boosting (XGB), and nearest neighbors are compared. The methods were applied to 10-year daily groundwater level time series at 157 sites in unconfined alluvial aquifers in Southern Germany. Models of HDCs were physically plausible and showed that physiographic and climatic controls on groundwater level fluctuations are nonlinear and dynamic, varying in significance from “wet” to “dry” aquifer conditions. XGB yielded a significantly higher predictive skill than nearest neighbor and MLR. However, donor site selection is of key importance. The study presents a novel approach for regionalization and infilling of groundwater level time series that also aids conceptual understanding of controls on groundwater dynamics, both central tasks for water resources managers.

## 1. Introduction

Groundwater head observations are the basis for most investigations in hydrogeology. However, boreholes for groundwater observation as well as corresponding groundwater level time series are often scarce and unevenly distributed in both space and time. This is a disadvantage for effective management of groundwater resources at the regional scale (Butler et al., 2021), where water managers assess the current and future status of groundwater resources (Lóaiciga & Leipnik, 2001). In consequence, methods are needed to estimate groundwater head time series at ungauged sites.

Two main approaches are commonly used by hydrogeologists to predict temporal changes in groundwater head at a given site, numerical and statistical models. The typical approach is to implement a process-based, numerical groundwater flow model. However, numerical models typically require large amounts of data and effort, while investigators commonly are confronted with a lack of comprehensive description and documentation of the subsurface. This results in significant uncertainty, both regarding conceptualization and parametrization (e.g., Enemark et al., 2019). Dealing with this uncertainty leads to a tedious and time-consuming process to construct, calibrate, and run these process-based models (Bakker & Schaars, 2019). Additionally, models for meaningful local projections at large spatial scales are not yet available (Berg & Sudicky, 2019). An alternative to regional scale modeling with less need for detailed subsurface description are lumped (rainfall-runoff) hydrological models with a groundwater component (Barthel & Banzhaf, 2016). However, these models are problematic as they usually imply oversimplification of the groundwater component, disregarding the local descriptors of hydrogeological systems and their three-dimensional setup (Barthel & Banzhaf, 2016; Butler et al., 2021). Generally, lumped models may provide adequate descriptions of groundwater systems only for simple hydrogeological situations such as shallow, unconfined aquifers, but not for more complex systems, such as deep and confined aquifers.

A different type of approach requiring only measured groundwater level data for groundwater time series estimation are parametric or data-driven methods. This type of approach requires few data on local system descriptors,

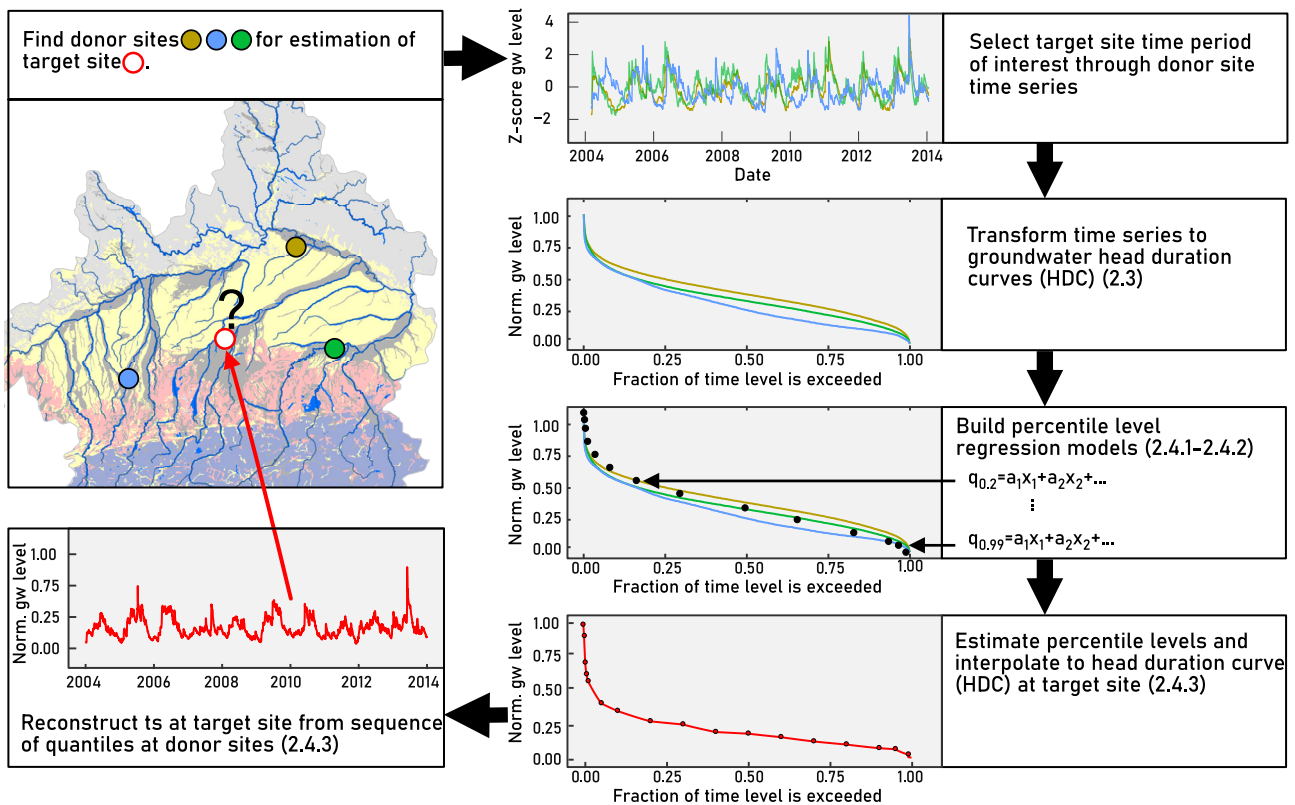
Writing – original draft: E. Haaf  
Writing – review & editing: E. Haaf, M. Giese, T. Reimann, R. Barthel

while often long and measurement-dense series of input signal and groundwater measurements are necessary to achieve good calibrations. In contrast to groundwater-gradient driven methods, data-driven methods either use spatio-temporal geostatistics (e.g., Ruybal et al., 2019; Varouchakis et al., 2022) or transfer net precipitation input into groundwater level changes (Chen et al., 2002). However, available methods predict groundwater level only at monthly or annual resolution and consequently do not capture the large intra-annual and intra-monthly variability of groundwater dynamics (e.g., Heudorfer et al., 2019). An approach to predict time series at higher temporal scales are transfer functions, that can be used to yearly, monthly and daily temporal resolutions, such as impulse-response functions (e.g., Collenteur et al., 2019; Marchant & Bloomfield, 2018; Von Asmuth, 2012) or artificial neural networks (cf. Rajaei et al., 2019; Wunsch et al., 2022). However, no formal method is known to transfer information from such models from monitored to unmonitored aquifers, although recently attempted in streamflow (Kratzert et al., 2019). This means that these methods can currently only make predictions when sufficient local time series data are available (e.g., 10 years weekly data, Wunsch et al., 2021).

In summary, neither numerical models nor the currently available data-driven tools provide a straightforward approach to estimate daily groundwater levels at unmonitored sites to aid regional scale management. Therefore, new and complementary methodologies are required to overcome scarcity and patchy data distribution. Such approaches should be less data-hungry than numerical models, yet account for local hydrogeological conditions and allow prediction at high temporal resolution despite limited local data availability. In surface-water-oriented hydrology, data scarcity has been countered with approaches of classification and similarity analysis, embraced by the hydrological community, particularly within the PUB initiative (Predictions in Ungauged Basins; Blöschl et al., 2013; Hrachowitz et al., 2013; McDonnell & Woods, 2004; Sivakumar & Singh, 2012; Wagener et al., 2007). These concepts attempt to systematically link the physical form and structure of catchments to their functioning by comparative analysis. Such links can then be used to transfer information to similar systems for prediction, that is, regionalization or spatio-temporal interpolation. However, such approaches are rarely considered in groundwater research, which is pointed out by various authors, for example, Barthel et al. (2021), de Marsily et al. (2005), Green et al. (2011), and Voss (2005). Recently, a number of studies initiated the implementation of such regionalization approaches in groundwater, quantitatively connecting groundwater response to physiographic and climatic descriptors (Boutt, 2017; Giese et al., 2020; Haaf & Barthel, 2018; Haaf et al., 2020; Heudorfer et al., 2019; Rinderer et al., 2014, 2016, 2017, 2019). These approaches, however, have not yet been exploited to predict daily groundwater level time series at unmonitored sites.

When looking for methodological inspiration in the body of literature within the surface water community, and more specifically the PUB initiative, a large majority of approaches use regionalization mainly as a tool to calibrate lumped rainfall-runoff models at unmonitored sites (He et al., 2011; Hrachowitz et al., 2013). As mentioned above, such lumped models are often not useful for describing groundwater dynamics and, when available, are time-consuming to set up and calibrate (Jackson et al., 2016; Mackay et al., 2014). Simpler statistical methods for regionalization of streamflow time series, however, have been proposed by for example, Shu and Ouarda (2012) based on Hughes and Smakhtin (1996). These methods make use of the characteristic relationship between flow duration curve (FDC; cumulative frequency of time where a flow is equaled or exceeded) and physiographic and climatic site descriptors, a relationship that is well investigated (Yokoo & Sivapalan, 2011). FDCs in surface water hydrology are commonly used to study the flow regime throughout the range of discharges and integrate effects of climate, topography, geology, and also anthropogenic activity (Ridolfi et al., 2020; Sugiyama et al., 2003; Vogel & Fennessey, 1995). This implies that the shape of a specific FDC is theoretically inferable from site descriptors. The technique evaluated in this study takes advantage of this through estimation of duration curves at unmonitored (target) sites based on similarity to neighboring donor sites. Then, from the estimated duration curve, time series are reconstructed at the target site into a daily time series (Hughes & Smakhtin, 1996; Mohamoud, 2010; Shu & Ouarda, 2012; Smakhtin, 1999).

Cumulative frequency or duration curves of groundwater heads are not as broadly used for studying groundwater resources, except when for example, analyzing the relative state of groundwater storage (e.g., Maxe, 2013). Giese et al. (2020) estimated aggregates (indices) of head duration curves (HDC) and linked differences in shapes to local, intermediate, and regional groundwater flow patterns. Haaf et al. (2020) found correlations between HDC indices and map-derivable physiographic and climatic site descriptors. These are indications that alike streamflow, system controls are integrated in groundwater head regimes and may be exploited by analysis of duration curves.



**Figure 1.** Principle steps to estimate groundwater head time series at unmonitored sites using the head duration curve methodology.

Regionalization and subsequent estimation of daily time series at unmonitored sites through duration curves of groundwater head is evaluated in this paper. The approach is based on the methodology proposed by Shu and Ouada (2012) for streamflow. It is adapted to groundwater, where groundwater HDCs as well as groundwater-relevant and map-derivable site descriptors are used. Within surface-water, this method has only been tested using stepwise multiple linear regression (MLR). In this study, besides MLR, a comparison is carried out with estimation through averaging of the nearest neighbor sites (NN) and extreme gradient boosting (XGB). XGB can represent nonlinear relationships between groundwater dynamics and site descriptors and has shown to be powerful in for example, recharge studies (Naghbi et al., 2020). In summary, a method is evaluated that may be used when an aquifer is partially monitored or unmonitored, but time series data at a particular site of interest are incomplete or unavailable. The regionalization approach is applied to unconfined, alluvial aquifers in a humid climate in Southern Germany at unmonitored sites using solely map-derivable site descriptors and data from neighboring locations.

## 2. Method and Data

### 2.1. General Strategy

The methodology of estimating groundwater head time series at an unmonitored site, is based on information from donor sites and requires the steps as explained in Figure 1. The principle is as follows: donor sites are selected from a regional data set with a time series period that is of interest for target site estimation. Next, time series are transformed to HDCs, and at 15 fixed percentile levels, percentile level models are constructed based on multiple regression analysis and gradient boosted regression trees, creating a regional HDC model (Sections 2.4.1 and 2.4.2). The regional model can then be used to predict the 15 fixed percentiles at locations across the entire domain and reconstruct site-specific HDCs with logarithmic inter- and extrapolation between percentile levels. Finally, from the site-specific HDC time series are estimated at unmonitored sites, with a distance-based weighting method using the sequence of records from donor sites (Section 2.4.3). For performance comparison, time series are also evaluated using only a distance-based average of time series from donor sites, further called



**Table 1**  
*Descriptive Statistics of Physiographic and Climatic Descriptors, Discussed in the Paper*

Variable	Description	Range		Unit
		Minimum	Maximum	
dist_stream (B)	Estimated distance from well to nearest stream (main rivers)	6	10,958	m
well_elevation (B)	Estimated Elevation of well	310	839	m asl.
P_avg (C)	Mean annual precipitation	675	1,613	mm
T_avg (C)	Mean annual temperature	6.4	9.3	°C
SI (C)	Seasonality index of precipitation	0.11	0.31	–
A_thickness (G)	Average thickness of saturated zone	1	50.1	m
A_Depth (G)	Bottom of formation	3	110	m
Depth_to_GW (G)	Average depth to water table	0.3	39.8	m
Broadleaved_forest (L)	% of 3 km buffer occupied by broadleaved forest	0	44.5	%
Coniferous_forest (L)	% of 3 km buffer occupied by coniferous forest	0	93.5	%
Urban (L)	% of 3 km buffer occupied by urban fabric	0	74.9	%
slp_sk (M) <sup>a</sup>	Mean slope	0/–0.1	1.95/2.6	–
twi (M)	Mean value of topographic wetness index	5.8	8.9	–

*Note.* Class of variable in parenthesis: (G) Geology, (M) Morphology, (L) Land cover, (B) Boundaries, and (C) Climate.

<sup>a</sup>Skewness was calculated for local and regional scale respectively. For these, the ranges are given separated by a slash //r.

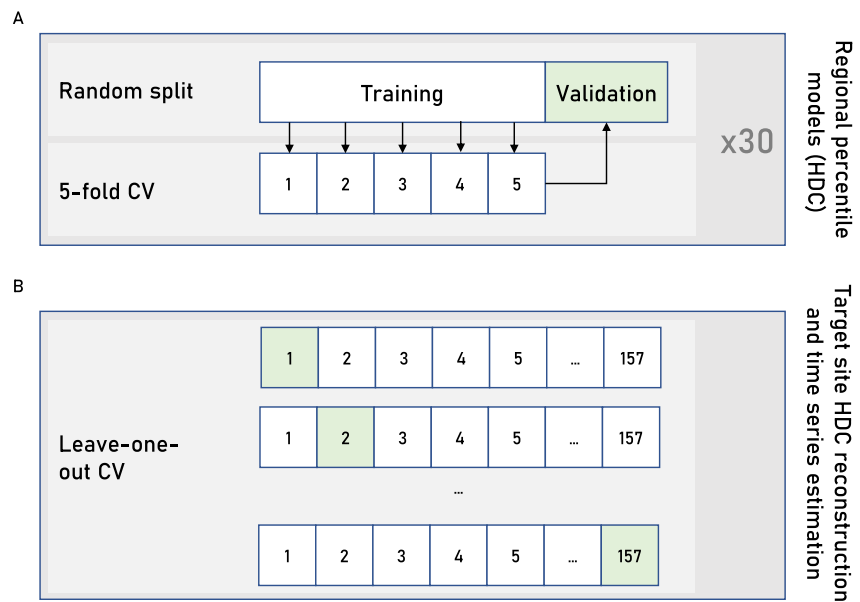
Nearest-Neighbor (NN). Then, the number of neighbors and the performance of daily groundwater head estimations at target sites are evaluated using leave-one-out cross-validation (Section 2.5). The models that are used for estimation of time series are then checked for plausibility (Section 2.6). In Section 2.7 the case data set is described, which is further analyzed using cluster analysis to understand results with regard to different groundwater regimes and systems. All data analysis was carried out by using the programming language R (R Development Core Team, 2022).

## 2.2. Data Selection and Processing

Groundwater head time series are selected from a data set described by Haaf et al. (2020). The data set contains groundwater head time series from the Upper Danube catchment in Bavaria, Southern Germany, with available geological information and absence of patterns of direct anthropogenic impact (for a more detailed explanation refer to Heudorfer et al. (2019)). From this data set observation wells were selected that come (a) with continuous daily time series and at least 10 year record length, (b) less than 1% missing data, which are (c) concurrent with a record period 2004–2014. The resulting set of 157 observation wells are mostly located in shallow, quaternary sediments in river valleys and fluvial sand as well as in gravel deposits, with a few boreholes located in deeper tertiary sediments. All wells are classified as penetrating unconfined aquifers. Then, at each site, 47 physiographical and meteorological descriptors were derived, described in detail in Haaf et al. (2020). In addition to Haaf et al. (2020), the percentage of land cover within a 3 km radius of each site was derived from the CORINE land cover data set (Bossard et al., 2000). Here, the classes were partially reclassified, such that artificial surfaces (class 1 were divided into Urban (1.1) and Other (1.2–1.4) surfaces) and forests was kept at level 3 (only broadleaved, coniferous and mixed forest classes were present). Agriculture, wetlands and water were kept at level 1. Table 1 shows selected descriptors that are most important for models on this study and therefore discussed in more detail. Remaining descriptors can be found in Table S1 in Supporting Information S1. Descriptors are called predictors when in context of regression models.

## 2.3. Transformation to Head Duration Curves (HDCs)

In a first step, groundwater head time series were normalized (on a 0–1 scale). Subsequently, duration curves of groundwater heads were calculated at each site. This was done, by first ranking all  $n$  observed, normalized



**Figure 2.** (a) Validation strategy for regional percentile models. The data set is randomly split and models are fit on 5-fold CV. Model fits are evaluated on the validation data. (b) When building the local head duration curve and estimating the time series, leave-one-out cross validation is carried out, removing the target site from the model data during model fitting but retaining all other (donor) sites. Model fits are evaluated on the left-out target site.

groundwater heads  $l_i$ ,  $i = 1, 2, \dots, n$  in descending order, where  $i$  is the rank of an observation. The head duration curve (HDC) is then constructed following the Weibull plotting formula (Sugiyama et al., 2003):

$$p_i = P(L \geq l_i) = \frac{i}{n + 1}, \quad (1)$$

where  $p_i$  is the percentile where a given groundwater head  $l_i$  is equaled or exceeded. Groundwater HDCs are subsequently created by plotting the percentage levels  $p_i$  against the corresponding heads  $l_i$  (as seen in Figure 1).

#### 2.4. Regression Analysis for Regional Percentile Models

To estimate the duration curve at an ungauged site, forward stepwise regression (MLR, see Section 2.4.1) and XGB (see Section 2.4.2) were applied to build regional models from physiographic and climatic predictors at a selected percentage level (0.1%, 0.5%, 1%, 5%, 10%, 20%, 30%, 40%, 50%, 60%, 70%, 80%, 90%, 95%, 99%). The validation strategy is shown in Figure 2a. The 15 percentile models are fit with 5-fold cross-validation on 80% of the training set using a random split. After this, predictions are made on the 20% validation data and compared to observed values. Since the data set is relatively small this procedure is iterated 30 times to be able to see how different selections of sites for the regional model impacts prediction. Each percentile model and iteration is evaluated using the coefficient of determination ( $R^2$ ).

##### 2.4.1. Construction of Percentile Models With MLR

MLR models at selected percentage levels are built using a selective inference framework. Selective inference adjusts  $p$ -values for the effect of sequential selection of variables (Taylor & Tibshirani, 2015). This is necessary since conventional stepwise regression leads to an overestimation of the strength of apparent relations. The consequence of conventional models is therefore selection of non-significant predictors and therefore overfitting (Taylor & Tibshirani, 2015). Instead of using  $p$ -values based on the  $t$ -test for forward selection, the procedure is here stopped based on the false discovery rate (exceeding 0.1; G'Sell et al., 2016). The selected variables are then used to build a regression relationship for the training data set with  $m$  observations (from well locations) and percentage levels,  $p = 1, 2, \dots, m$ , where  $H_p$  is the percentile of the normalized head  $H$  and  $x_p$  the selected climatic and physiographic descriptors with the following form:

$$H_p = \beta_0 + \sum_j x_{pj} \beta_j + \epsilon_p, \quad (2)$$

errors  $\epsilon_p$  being independent and normally distributed and where  $\beta$  is a vector of model parameters that are estimated.

#### 2.4.2. Construction of Percentile Models With XGB

Alternative models for each percentile were constructed using XGB, an implementation of boosted regression trees (Friedman, 2001). Hereby, the *xgb.train* function from the XGBoost R package (Chen & Guestrin, 2016) was used to predict  $H_p$  based on the entire set of climatic and physiographic descriptors. To optimize the model fit the XGB ensemble is stopped at the number of decision trees, where the difference between training and evaluation error reaches a minimum. Hyperparameters for the fitted XGB model can be found in the Table S2 in Supporting Information S1.

#### 2.4.3. From Regional Percentile Models to Locally Estimated Time Series

Once regional percentile models are built, percentile levels can be predicted for a given target site using XGB and MLR models using leave-one-out cross validation (LOOCV) (Figure 2b). To reconstruct a continuous HDC from the 15 percentile levels, logarithmic interpolation is used to estimate percentiles of groundwater heads between the percentage levels. The percentile to be estimated is found by identifying the closest (modeled) fixed percentage levels  $p_i$  above and  $p_{i-1}$  below and their corresponding groundwater heads  $H_i$  and  $H_{i-1}$ . The groundwater head  $H$  can then be found using the following equation:

$$\ln(H) = \ln(H_i) + \frac{\ln(H_{i-1}) - \ln(H_i)}{p_{i-1} - p_i} \times (p - p_i) \quad (3)$$

In cases where percentiles are estimated that are larger than the highest percentage point or lower than the lowest (modeled) percentage point, logarithmic extrapolation is used. Hereby, the closest two percentage points are found ( $p_{n1}$ ,  $p_{n2}$ ) and the corresponding groundwater heads ( $H_{n1}$ ,  $H_{n2}$ ). Extrapolating to the percentile  $p$  is done using the equation below.

$$\ln(H) = \ln(H_{n1}) + \frac{\ln(H_{n1}) - \ln(H_{n2})}{p_{n1} - p_{n2}} \times (p - p_{n2}) \quad (4)$$

Reconstruction of the groundwater head time series from interpolated duration curves can then be carried out following the principle given by Smakhtin (1999) for streamflow estimation. Groundwater heads  $H_t$  at the target site are estimated by looking up the donor site's percentile of the duration curve at the first date to be estimated. Then the same percentile is found in the reconstructed target site's duration curve and the corresponding groundwater head is chosen as the estimated level at the particular date. This process is repeated for all dates available within the record of the donor sites. However, not all donor sites are given the same weight for estimation at the target site. The estimated series of groundwater heads at the target site  $H_t$  are rather put together (Equation 5) by weighting each donor site's contribution based on the Euclidean distance  $d_t$  to the target.

$$H_t = \sum_{j=1}^n w_j H_{sj} / \sum_{j=1}^n w_j \quad (5)$$

The weights are calculated based on a dissimilarity measure:

$$w_j = \frac{1/d_t}{\sum_{j=1}^n 1/d_t} \quad (6)$$

Groundwater heads are also estimated at each target site using a straightforward NN method. Here, NN means that no duration curve is reconstructed but only the actual time series of each donor site  $L_{ij}$  is used, however, weighted according to Equations 5 and 6.

### 2.5. Evaluation of Time Series Estimation

The performance of the daily groundwater head prediction was evaluated using LOOCV as performed by Shu and Ouarda (2012). Using a LOOCV procedure means that one (target) site is considered unmonitored and thus left out from the data set (Figure 2b). With the remaining data set ( $n - 1$  sites), the groundwater head time series are

estimated at the target site. Here, a maximum of  $n = 20$  sites were allowed as donor sites. Then, the performance at that site is evaluated by calculating the Kling-Gupta Efficiency (KGE), Pearson correlation coefficient ( $R$ ), and root-mean-square error (RMSE) as goodness of fit measures between estimated and observed time series. These steps are repeated at each of the  $n$  sites and the average (cross-validated) estimate is found by aggregating the goodness of fit-estimates from each sub-sample.

## 2.6. Plausibility Analysis of Models

To examine the plausibility of models used to predict percentile points along the HDC, the impact on model output is analyzed using standardized regression coefficients (MLR) and Shapley Additive Explanations values (SHAP) for XGB (Lundberg et al., 2020) using the  $R$  package *SHAPforxgboost* (Liu & Just, 2021). SHAP values quantify how much individual predictors, across the predictor's value range, contribute to the output variable (here the percentile point). More specifically, the SHAP value gives the difference in the model output depending on if the model is fit with or without the predictor. Using scatterplots, SHAP values can then be interpreted locally which allows understanding of the dependence structure within each model for each predictor. Further, mean absolute SHAP of all data points for each model is estimated, yielding global feature importance across each percentile. This supports understanding of the dynamic changes of importance of controls across different aquifer states and allows qualitative comparison to standardized regression coefficients of MLR models.

## 2.7. Cluster Analysis

In order to get a better understanding of the data set, regarding similarities in dynamics and subsequently site descriptors, hierarchical cluster analysis was performed with a strategy described in Giese et al. (2020). Prior to cluster analysis, the selected groundwater head time series are transformed to  $z$ -scores. As input into the clustering algorithm, Euclidean pairwise distances between time series were computed. Subsequently, hierarchical cluster analysis using Ward linkage is performed on the matrix of pairwise distances. The hierarchical relationship between the series can then be displayed in a dendrogram. From the dendrograms a scree plot is constructed, by sorting the heights of the dendrograms branches and plotting these against the number of nodes. The inflection point of the scree plot is then identified to select the number of clusters that sufficiently describes the patterns of member time series, while still generalizing the data set to a manageable level.

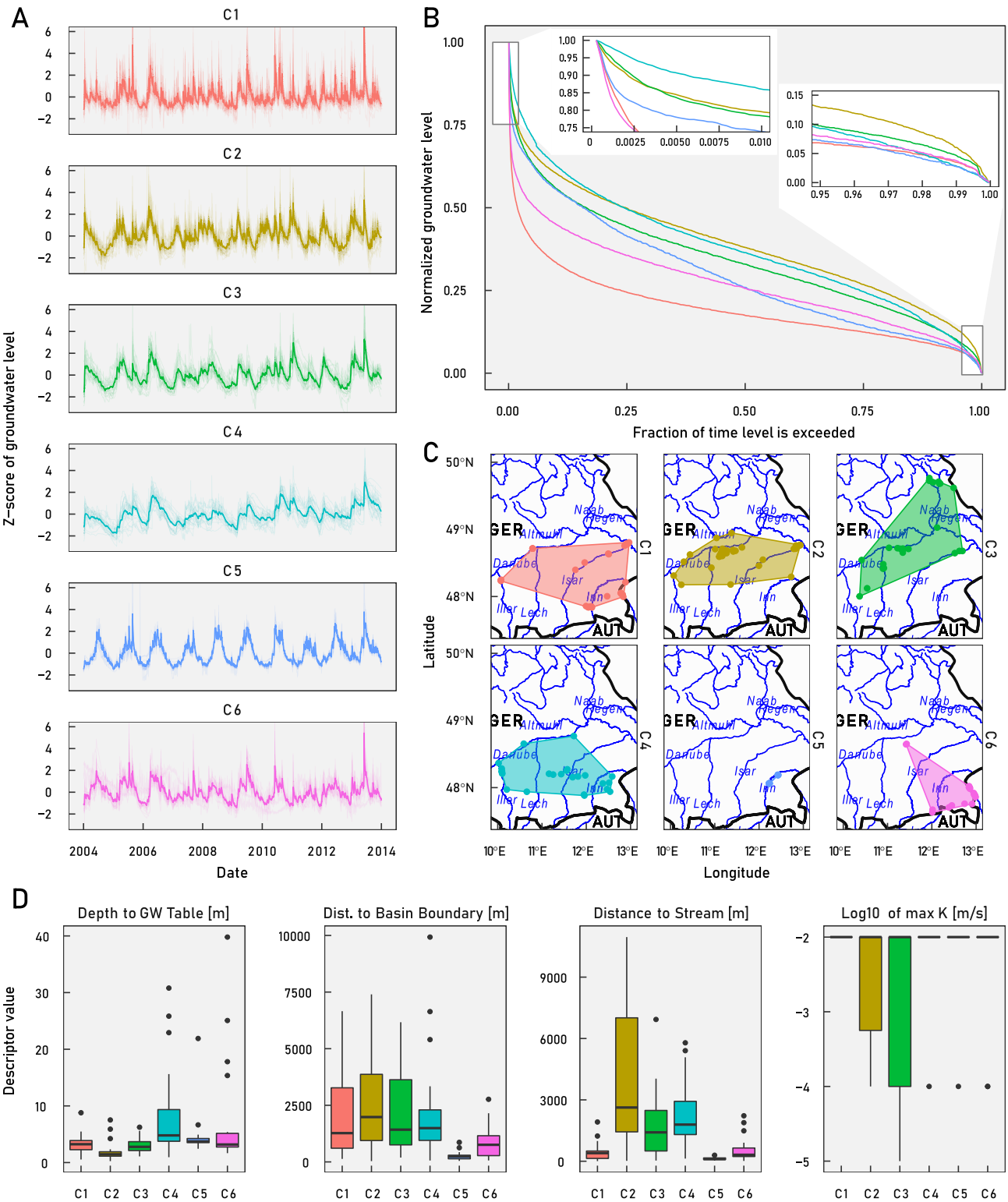
# 3. Results and Discussion

## 3.1. Hydrogeological Description of Clusters

Cluster analysis of the data set based on similarity of groundwater head time series results in hydrogeologically meaningful groups. The six identified clusters (see Figures S1 and S2 in Supporting Information S1) are either made up of wells exclusively located in alluvial deposits or in alluvial deposits and outwash plains. Further, cluster separation can be linked to differences in distance to stream, depth to water table, size of aquifer, local hydrology and geographical location.

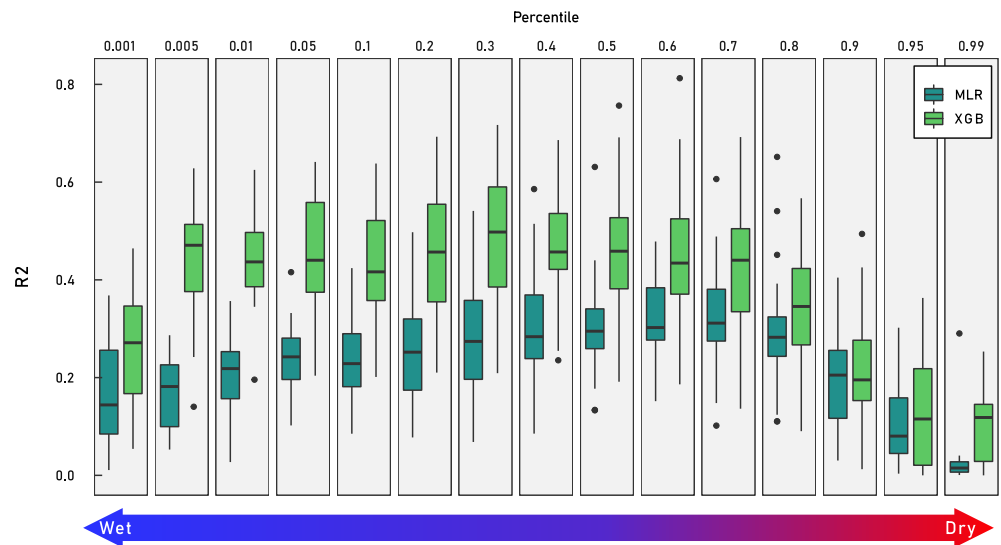
Figures 3a and 3b show that groundwater head time series in clusters C1 and C6 have similar groundwater regimes. Time series in C1 show a relatively fast response (flashy) and overprinting of high peaks to varying degree, which is seen to a slightly lesser degree in C6. Inter- and intra-annual patterns are mostly absent. Groundwater heads in these two clusters are shallow ( $75\% < 5$  m) and with the wells relatively close to groundwater basin boundaries and streams in medium size aquifers (Figure 3d). Presumably, these clusters represent wells tapping mainly local groundwater flow systems (Giese et al., 2020). The pronounced flashiness is linked to interaction with streams (Haaf et al., 2020) and can also be seen in the low percentiles of the duration curves that are significantly steeper in the flashier C1 and C6 than other clusters (Figure 3b). Differences between C1 and C6 can be attributed to the different geographical areas, with C1 located in more extensive aquifers far downstream of the headwater catchment in the South and C6 located mainly in smaller alluvial aquifers in the Salzach and Inn catchments at the foot of the Alps (Figure 3c and Figure S3 in Supporting Information S1).

Flashiness in cluster C2 is like C6, however, exhibiting intra-annual variations and weak inter-annual seasonality. Like C1 and C6, C2 is characterized as local flow due to the very shallow wells, however, wells are in intermediate locations in large aquifers. Therefore, dynamics are not closely coupled to the major rivers, which are at larger



**Figure 3.** (a) Time series within each cluster. (b) Mean of groundwater head duration curve of color related to cluster in (a). (c) Location of cluster members with convex hull and stream network, ISO 3166-1 alpha-3 country codes. (d) Hydrogeological descriptors of sites within each cluster.





**Figure 4.** Performance measured in coefficient of determination of cross-validated percentile regression models.

distances, but presumably to (unmapped) smaller creeks and to vegetation considering the shallow groundwater table.

C3 is less flashy than C2, but shows a similar inter- and intra-annual pattern, which can also be seen in the similarity of the two cluster's HDCs (Figure 3b). C3 wells are, similar to C2, located in larger aquifers, but are deeper and closer to streams, likely representing local and intermediate flow systems.

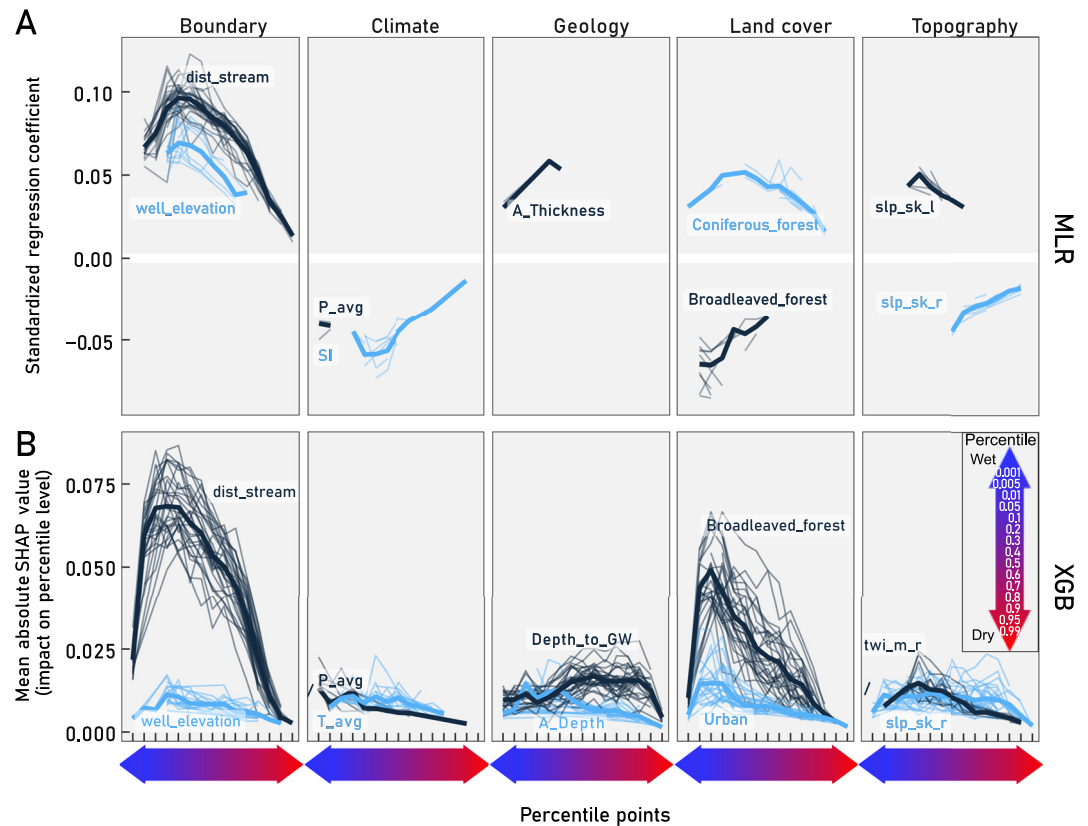
C4 has dominant inter-annual variability, which is linked to the larger distance to groundwater head and streams (Haaf et al., 2020). The larger inter annual variability in C4 is also seen in the less steep lower percentiles of the duration curves (Figure 3b) and is linked to mainly intermediate and regional flow systems.

Groundwater hydrographs in cluster C5 show a very distinct pattern compared to the remaining clusters. The HDC falls steeply at lower percentiles, following the flashier C1 and C6, until stabilizing and resembling more the weakly intra-annual dominated HDCs of C2 and C3, before crossing back to C1 and C6 at higher percentiles, due to cluster's weak intra-annual periodicity. The distinct pattern and in-group similarity of the 14 wells in C5 is explained by their locations, concentrated near the Inn, which is regulated by run-of-the-river hydroelectric plants with pondage (Figure 3c).

### 3.2. Performance of HDC Reconstruction

After regression analysis, models were found for all 15 fixed percentage points. Regression models fitted on 30 different sets of hold-out data resulted in a distribution of results that are robust with regard to central tendency. Median XGB model performance on hold-out data expressed as  $R^2$  is around 0.5, except for the lowest and upper percentiles (0.1%, 80%–99%), that is, wet and dry states, where goodness-of-fit declines (Figure 4). XGB models perform significantly better than MLR models that show a similar behavior across percentiles but with lower goodness-of-fit (median  $R^2$ : 0.3). Figure 4 also shows that the range of  $R^2$  is large, which is very likely related to the size of the data set. The consequence of small data sets, when using hold-out data is that the evaluation data (here,  $n = 32$ ) may not be representative of the training data across sets of hold-out data. Further, when running models on the entire data set (training + evaluation), both XGB and MLR models show around 100% and 70% performance improvement from median  $R^2$ . Performance loss across hold-out data and against the entire data set indicates that generalization from the training set is moderate.

When comparing results to studies using an analogous methodology in streamflow, model results of  $R^2$  between 0.72 and 0.99 are reported and analogous lower values in the extremes (Mohamoud, 2010; Shu & Ouarda, 2012). However, these studies used neither hold-out data, cross-validation methods, or p-value adjustment for step-wise MLR. This means that the higher  $R^2$  values of models presented in these studies are likely overfitting and



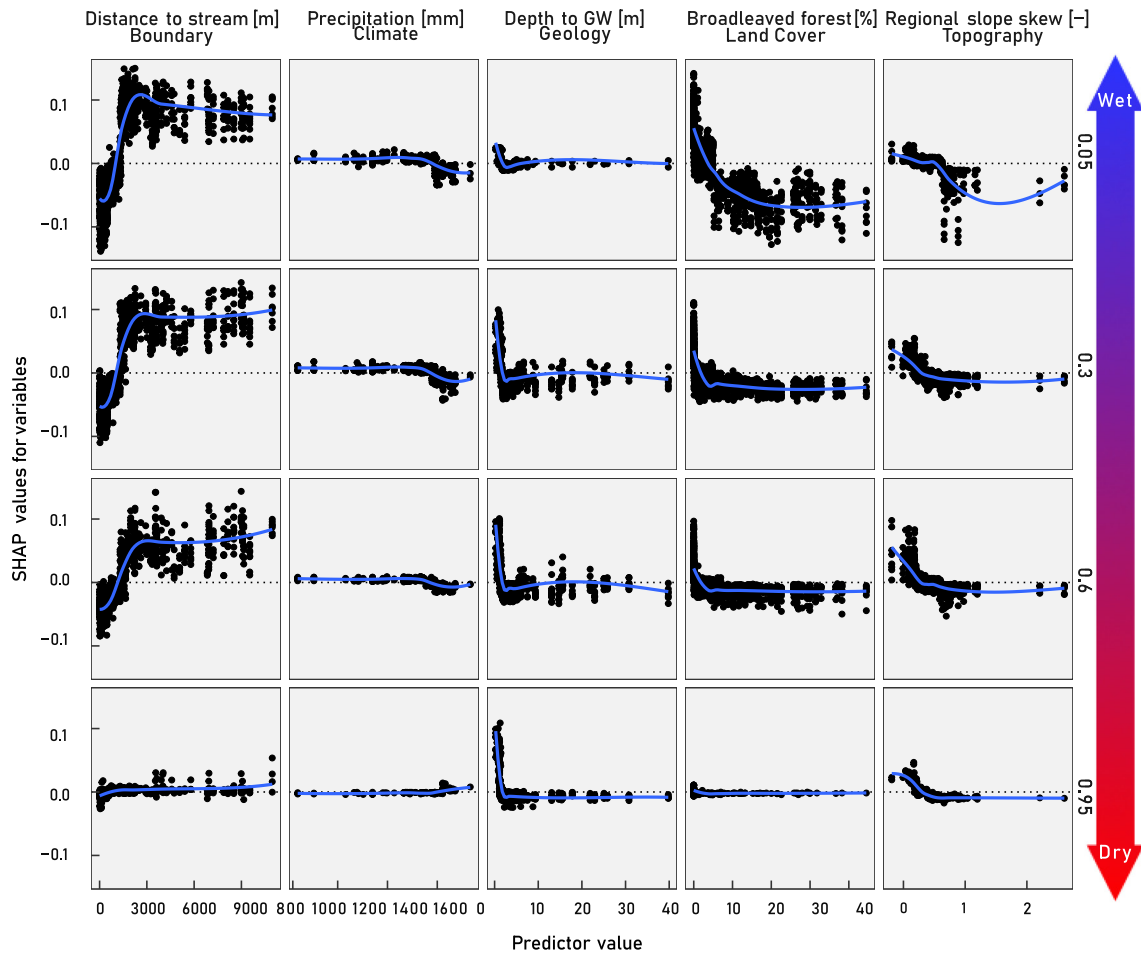
**Figure 5.** Relative predictor importance across percentage point models stratified by predictor class for multiple linear regression and extreme gradient boosting models (scales not comparable). Data from all hold-out data sets are plotted. If lines do not cover all percentiles, the predictor is not statistically significant (a) or zero (b), that is, has no impact on the prediction. In addition the mean of hold-out data sets is plotted to emphasize the central behavior of the data. (a) Standardized regression coefficients show both relative predictor importance and direction of relationship between predictor and model output. (b) Mean absolute Shapley Additive Explanations value shows relative importance through impact on the output variable.

generalization outside of the data set could be questioned. The performance achieved on evaluation + training data by XGB and MLR models in this study would thus be more comparable and are in fact in parity with performance reported in streamflow studies.

### 3.3. Dynamic Controls on Groundwater Heads

Relative predictor importance across percentage point models stratified by predictor class for MLR and XGB models respectively is shown in Figure 5. Standardized regression coefficients in MLR give both relative predictor importance (higher absolute value) but also the direction of the relationship between predictor and output variable (percentile level of HDC) through the sign of the coefficient (Figure 5a). Mean absolute SHAP value on the other hand, shows only relative predictor importance (Figure 5b). Further, for clarity of presentation, only the most salient variables are shown (MLR: variables are shown that are selected in at least 30% of hold-out data sets; XGB: only the top two predictors are shown per predictor class based on overall mean absolute SHAP value).

The main result is that the importance of predictors varies across percentiles. This implies that different site (or system) descriptors to varying extents control the groundwater dynamics when the aquifer is moving from “wet” to “dry” states and vice versa. An example is distance to stream that is important through all aquifer states but dominating in wet states (both MLR and XGB, Figures 5a and 5b). Depth to the groundwater table, on the other hand, becomes more dominant when the aquifer is in dry states (only XGB, Figure 5b). A pattern that can be seen across all variables is that predictor strength declines significantly (approaches zero) at higher percentiles, which is also connected to lower goodness-of-fit at these percentiles (Figure 4). Consequently, predictability of percentiles coupled to groundwater drought is lower.



**Figure 6.** Relationship between feature value and impact on prediction for five selected variables across four percentiles (all percentiles can be found in Figure S4 in Supporting Information S1). Each point represents an observation of the predictor variable and its Shapley Additive Explanations value. Data from all hold-out data sets are plotted and fitted with a local polynomial regression to emphasize the central behavior of the data.

Another important finding is that many of the most important predictors are consistently selected across both MLR and XGB as well as show a similar importance progression across percentiles (e.g., distance to stream, well elevation, average annual precipitation, broadleaved forest and regional slope skewness). This means that many of the important variables have a sufficiently linear relationship with percentiles of groundwater HDCs so that it can be picked up by MLR. For instance, MLR models show that percentage points of the HDC increases with distance to stream. Consistently high values across percentage points means flat HDCs, which is associated to lower flashiness (compare e.g., HDC and time series of C3 and C4 to C1 in Figures 3a and 3b). This is plausible and expected, since streams are the aquifer's given drainage boundary and known through previous regional scale empirical studies (e.g., Boutt, 2017; Giese et al., 2020; Haaf et al., 2020; Vidon, 2012). However, SHAP values of individual data points related to XGB prediction allows us to look more closely at linearity of relationships between HDC and predictor value ranges (Figure 6). The SHAP values reveal a more complex relationship, where the relationship between distance to stream and dynamics is constant up to about 500 m distance, turning into a linear relationship, where groundwater dynamics become less flashy with distance until reaching a plateau at about 3,000 m distance. Here, presumably groundwater is no longer strongly connected to the stream and a constant contribution to the HDC is reached (Figure 6). This effect is consistent across aquifer states, however weakens, when the groundwater head drops into dry states. The nonlinearity of relationships with threshold effects is common, as described below for variables selected in Figure 6:

- Average annual precipitation has relatively low impact on the HDC, which is also true for other climate predictors in this study. However, precipitation below approximately 800 mm leads to slightly less flashy

dynamics in wet states. This can be coupled to less infiltration and recharge events. At higher precipitation rates, no systematic impact on HDC can be seen.

- Depth to groundwater table only affects the HDC when very shallow, approximately 2 m and above. Shallow water tables increase the percentile level accordingly, meaning that less flashiness may be expected. Sites, where groundwater heads are very shallow may be coupled to discharge zones. Here the aquifer is continuously replenished through recharge from uplands with significant upward hydraulic gradients (Gribovski et al., 2010; Winter 2001). Generally, this effect increases in importance at higher percentiles, that is, in a drier aquifer state.
- If the percentage of broadleaved forests exceeds approximately 10%, groundwater heads become flashier in wet states, which can be linked to higher soil moisture, preferential flow and recharge than other land cover types, reducing surface runoff (Brinkmann et al., 2019; Dubois et al., 2021).
- If regional slopes are right skewed, sites are located in alluvial valley bottoms at the fringes of higher hill ranges (Haaf et al., 2020; Montgomery, 2001). In these locations amplitudes are expected to be higher due to front slope flow and mountain block recharge, which is also seen here particularly in wet aquifer states with lower SHAP values at higher slope skewness. Low slope skewness ( $<0.0.3$ ) on the other hand contributes to less flashy groundwater dynamics.

Overall, the progression of controls have implications not only for prediction but also conceptual understanding of groundwater dynamics in this region. The nonlinear relationships of groundwater dynamics and controls and the alternating dominance of these controls throughout different aquifer states are likely of interest, for example when studying for example, vulnerability to drought events and climate change. Certainly, there is a need for a dedicated analysis of the dependence of controls on aquifer states, which was outside of the scope in this study.

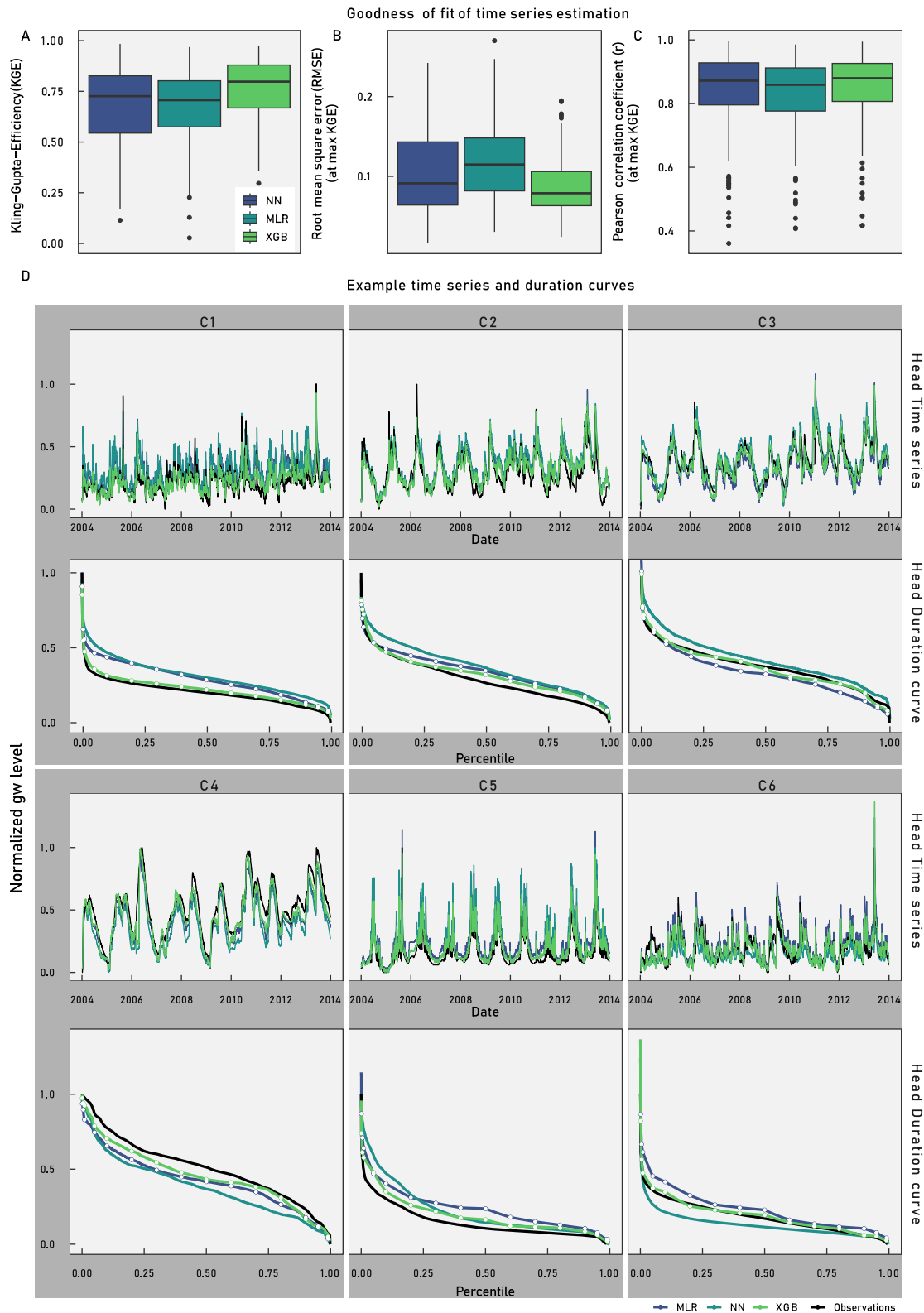
### 3.4. Performance of Estimation Techniques

Daily groundwater head time series were estimated at target sites, using trained models from each of MLR and XGB models as well as using the NN method. The XGB model had a higher KGE than NN at 120 of 157 (76%) sites, and a higher KGE than MLR at 136 of 157 (87%) sites. In consequence KGE is also considerably higher for XGB than NN and MLR (Figure 7a). Interestingly, MLR has a lower median KGE than NN, (slightly higher performance at the lower quartiles) which means that HDC modeling in the case of MLR deteriorates estimation on average, compared to the simple NN approach.

The higher performance of XGB can almost entirely be attributed to smaller amplitude errors between simulated and observed time series. Amplitude errors are expressed by the RMSE component of KGE, which is much improved when using XGB compared to NN and MLR (Figure 7b). The correlation component of the KGE on the other hand shows no significant differences between methods, meaning that timing errors between observed and simulated time series are not significantly improved through XGB or MLR (Figure 7c). As discussed by Mohamoud (2010), timing errors are coupled to the mismatch of time sequence in hydrograph events (here, e.g., recharge events) at donor and target sites. The difference in magnitude and timing errors among the methods can be seen in example duration curves and time series from each cluster in Figure 7d. While the time series estimated with different methods follow the time sequence pattern of the observed data quite similarly, a larger difference can be seen in how well the magnitude is matched. This is more easily visible in duration curves, where all methods show some systematic departures from the observed data. XGB estimates, however, follow the true duration curve most closely. In conclusion, from a water resources management perspective, the HDC estimation approach using XGB implies better estimation of the quantitative status of groundwater resources through reduced amplitude errors.

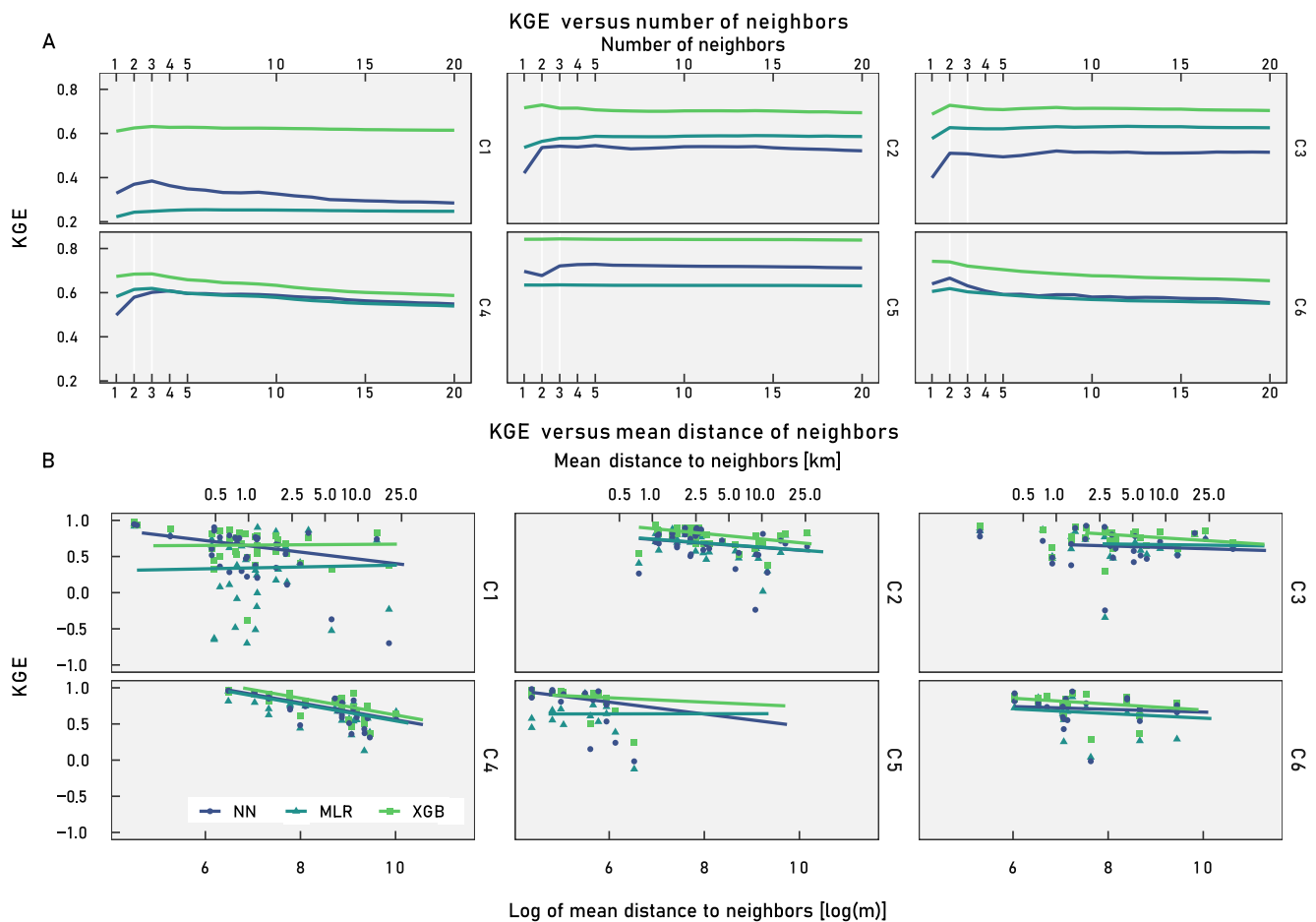
Figure 8a shows that an optimal number of donor sites (neighbors) is generally reached with only 1–3 neighbors, as expressed by the maximum KGE. Sourcing more neighbors generally results in plateauing or even decrease of estimation performance across different groundwater regimes, as expressed by clusters C1–C6. Although the number of optimal donor sites is consistent, C4 and C6 exhibit a sharp decline, when more than three or two donor sites respectively are added. A possible reason for this is that these two clusters contain sites with significantly deeper groundwater tables (Figure 3d). This means that donor sites with for example, more shallow water table and therefore deviating groundwater response will be weighted in and cause a mismatch of time sequence, decreasing the quality of the predicted groundwater head time series at the target site.

Not only hydrogeological suitability of donor sites is important, but also proximity (Figure 8b). Performance decreases approximately with the natural logarithm of mean distance of neighbors. However, even at large mean



**Figure 7.** (a–c) Performance of estimation of daily groundwater head time series for the three approaches across all unmonitored sites, measured as Kling-Gupta Efficiency (a), root-mean-square error (b), Pearson's  $r$  (c). (d) Example of estimated duration curves (white dots are estimates from percentile models) and time series at a single target site in each cluster with approximate median performance.





**Figure 8.** (a) Mean performance—measured as Kling-Gupta Efficiency (KGE)—of the three estimation methods plotted against number of included neighboring sites, stratified by cluster. (b) Performance of all sites—measured by KGE—plotted versus mean distance to neighbors, stratified by estimation method and cluster.

distances to donor sites (e.g., >5 km), estimation performance at many sites may remain high. This is particularly the case for cluster C2 and C3. These clusters show higher performances by both HDC-based estimation techniques MLR and XGB. On the other hand, at sites with sufficient neighbors nearby (<5 km), NN is preferred over MLR. Overall, however, XGB yields best performance independently of mean distance to neighbors.

### 3.5. Hydrogeological Controls and Plausibility of Models

From a hydrogeological perspective, there are obviously missing descriptors to describe groundwater heads, such as aquifer properties, transmissivity and storativity. These are often not consistently available at the scale of this study (regional scale), or only with a low level of certainty at the level of 1–2 orders of magnitude (e.g., hydraulic conductivity in this study). However, it can be argued that the importance of storativity in this study is reduced, since normalization on a 0–1 scale of groundwater head time series reduce the importance of amplitude. Regarding hydraulic conductivity a relatively homogenous selection of sites is used (Figure 3d). When assuming order of magnitude similarity of hydraulic conductivity, the predictor aquifer thickness ( $A_{thickness}$ ) may be considered a rough proxy for transmissivity. With these simplifications and proxy variables, model fits are acceptable, but still contain significant uncertainty, resulting in lower quality of time series prediction. Adding hydraulic properties, that is, storativity values and less uncertainty regarding hydraulic conductivity to the set of predictors would likely improve the fit of regression models. It would further allow for use of more heterogeneous data sets. Different strategies to extract such hydraulic properties at wells from groundwater head time series of unconfined aquifers was recently proposed using transfer function noise models (Peterson & Fulton, 2019) and spectral analysis (Houben et al., 2022).

Apart from the missing hydraulic properties, other factors likely play a role in explaining the moderate goodness-of-fit of the HDC models. Some of the uncertainty may be due to different hydraulic properties stratified within the zone of fluctuation. This is the case at only a few sites according to the borehole logs. Other sources of uncertainty may be found in data (groundwater head measurements, spatial resolution of DEM and climate data) or method of estimating physiographic and climatic descriptors. Other reasons may be found in the overrepresentation of relatively shallow alluvial aquifers, particularly in the north-east of the study area. Using mean squared error as a loss function, regression models tend to better represent the bulk of the sites within the data set, which are mainly lowland riverine aquifers with shallow groundwater heads (local groundwater flow) and less so the peri-alpine river valleys in the north-east. A functional stratification of the data prior to HDC model building by for example, the dominating predictor distance to stream, or more conceptually-based, using the hydrological landscape concept (Winter 2001) may improve the predictive performance of the HDC models for sites that are less well represented. Using these functional pre-classifications should also improve transferability of methods to other study domains. For such an exercise, however, a data set would be necessary with sufficient data points that ensures robust models in each functional stratum.

### 3.6. Improvement of Donor Selection

The bias of the models toward well-represented hydrogeological settings as described above, also has consequences on donor-based reconstruction of time series at unmonitored sites. As discussed in Section 3.4, differences in timing error between the three methods, NN, MLR, and XGB, are very small and related to the similarity of time sequences between target and donor sites. A mismatch occurs, when inadequate donor sites are selected, which can be seen for example, in cluster C4 and C6 (Figure 8a). Performance in these clusters declines with each additional donor and is presumably related to donors for intermediate/regional flow (C4) target sites being selected from (C6) sites that are located near rivers. In other words, donor sites have hydrological responses that differ from the target sites. Similar responses at sites with intermediate and regional flow systems can however be expected even at larger distances (Giese et al., 2020; Haaf & Barthel, 2018). In consequence, careful selection of donor sites is crucial to the performance of the method (also pointed out by authors applying the approach to streamflow: e.g., Hughes & Smakhtin, 1996; Shu & Ouarda, 2012; Smakhtin, 1999) and geographical proximity should not always be the main or sole selection criteria for donor sites.

Likely, a cleverer approach than solely proximity for donor site selection, would surely improve the performance of the presented approach. Such a strategy could be based on a distance metric that uses physiographic and climatic site descriptors for quantification of similarity between sites, as proposed for streamflow by Shu and Ouarda (2012). However, after studying the nonlinearity of relationships between site descriptors and groundwater dynamics, a non-continuous approach may be more useful. Often, step changes could be seen, which indicates that a discrete classification approach may provide a more optimal pool of donor sites. Such classes of similar responses could be developed from the SHAP values in Figure 6, for example, that neighbors must be within the same distance to stream, that is, within one of three classes (1–500, 500–1,500, >1500 m). For many of the sites, however, nearby sites still provide the most adequate timing of events. Therefore, any of the donor selection strategies discussed above must be combined with an approach that applies weights to donors within the similar class based on proximity.

## 4. Conclusions

Using the presented method, groundwater head duration curves can be transferred based on comparative regional analysis of map-derivable site descriptors from monitored to unmonitored sites. Neighboring donor sites can then be used to successfully reconstruct the daily groundwater head time series based on the transferred duration curve. Apart from time series estimation at unmonitored sites—in essence spatio-temporal interpolation—the modeling approach aided by physiographic and climatic descriptors also gives insight into hydrological processes through identification of significant controls. Specifically, at the study site, controls on groundwater dynamics were nonlinear, which favors use of Machine Learning (i.e., gradient boosted regression trees) over MLR and therefore makes possible improved conceptual hydrogeological understanding as well as higher predictive skill. The method and results were robust as tested through nested cross-validation, however, require thorough testing with larger data sets for application in other hydrogeological settings.

The study also showed that only 1–3 neighboring donor sites are generally necessary to optimally reconstruct time series of unmonitored sites. Further, it could be shown that inclusion of physiographic and climatic controls

add information to time series predictions. However, the fewer nearby donor sites are available, the more benefit can be drawn from these controls with the proposed comparative regional analysis approach, compared to NN averaging of time series. Importantly, the selection of donor sites was identified as a key factor to improve predictive skill and should be expanded on using a combination of geographical proximity and functional classes of groundwater sites from which to draw appropriate neighbors. Finally, the study shows ways forward to investigate the dynamic nature of controls on groundwater heads, which may provide valuable insight to studies of recharge seasonality, droughts and floods.

## Data Availability Statement

Groundwater time series cannot be provided publicly by the authors based on the data usage agreement with the LfU, but can be downloaded from <https://www.gkd.bayern.de/en/groundwater/upper-layer> and <https://www.gkd.bayern.de/en/groundwater/deeper-layer>. The selected station names are provided in Supporting Information S1. Processed data for regional duration curve models can be found at <https://doi.org/10.5281/zenodo.8046570>. Code for reproduction of results can be obtained from the corresponding author. All the analysis was performed in the statistical language R (R Development Core Team, 2022) using apart from the packages mentioned in the body “tidyverse” (Wickham et al., 2019), “lubridate,” “rsample” (Frick et al., 2022), and “selectiveInference.” The authors thank the contributors of all these packages.

## Acknowledgments

The authors would like to thank the German federal state agency Bayerisches Landesamt für Umwelt (LfU, <https://www.lfu.bayern.de>) for the provision of data and Supporting Information S1. The work partly contributes to the Swedish Research Council Formas project 2016-00513 and the Swedish Transport Administration project TRV2019/45670. Big thanks to Lars Rosén and the reviewers for valuable comments.

## References

- Bakker, M., & Schaars, F. (2019). Solving groundwater flow problems with time series analysis: You may not even need another model. *Ground Water*, 57(6), 826–833. <https://doi.org/10.1111/gwat.12927>
- Barthel, R., & Banzhaf, S. (2016). Groundwater and surface water interaction at the regional-scale—A review with focus on regional integrated models. *Water Resources Management*, 30(1), 1–32. <https://doi.org/10.1007/s11269-015-1163-z>
- Barthel, R., Haaf, E., Giese, M., Nygren, M., Heudorfer, B., & Stahl, K. (2021). Similarity-based approaches in hydrogeology: Proposal of a new concept for data-scarce groundwater resource characterization and prediction. *Hydrogeology Journal*, 29(5), 1693–1709. <https://doi.org/10.1007/s10040-021-02358-4>
- Berg, S. J., & Sudicky, E. A. (2019). Toward large-scale integrated surface and subsurface modeling. *Ground Water*, 57(1), 1–2. <https://doi.org/10.1111/gwat.12844>
- Blöschl, G., Sivapalan, M., Wagener, T., Viglione, A., & Savenije, H. (2013). *Runoff prediction in ungauged basins: Synthesis across processes, places and scales*. Cambridge University Press.
- Bossard, M., Feranec, J., & Otaheh, J. (2000). CORINE land cover technical guide: Addendum 2000.
- Bout, D. F. (2017). Assessing hydrogeologic controls on dynamic groundwater storage using long-term instrumental records of water table levels. *Hydrological Processes*, 31(7), 1479–1497. <https://doi.org/10.1002/hyp.11119>
- Brinkmann, N., Eugster, W., Buchmann, N., & Kahmen, A. (2019). Species-specific differences in water uptake depth of mature temperate trees vary with water availability in the soil. *Plant Biology*, 21(1), 71–81. <https://doi.org/10.1111/plb.12907>
- Butler, J. J., Knobbe, S., Reboulet, E. C., Whittemore, D., Wilson, B. B., & Bohling, G. C. (2021). Water well hydrographs: An underutilized resource for characterizing subsurface conditions. *Groundwater*, 59(6), 808–818. <https://doi.org/10.1111/gwat.13119>
- Chen, T., & Guestrin, C. (2016). XGBoost: A scalable tree boosting system. <https://doi.org/10.1145/2939672.2939785>
- Chen, Z., Grasby, S. E., & Osadetz, K. G. (2002). Predicting average annual groundwater levels from climatic variables: An empirical model. *Journal of Hydrology*, 260(1), 102–117. [https://doi.org/10.1016/S0022-1694\(01\)00606-0](https://doi.org/10.1016/S0022-1694(01)00606-0)
- Collentour, R. A., Bakker, M., Calje, R., Klop, S. A., & Schaars, F. (2019). Pastas: Open source software for the analysis of groundwater time series. *Ground Water*, 57(6), 877–885. <https://doi.org/10.1111/gwat.12925>
- de Marsily, G., Delay, F., Gonçalves, J., Renard, P., Teles, V., & Violette, S. (2005). Dealing with spatial heterogeneity. *Hydrogeology Journal*, 13(1), 161–183. <https://doi.org/10.1007/s10040-004-0432-3>
- Dubois, E., Larocque, M., Gagné, S., & Meyzonat, G. (2021). Simulation of long-term spatiotemporal variations in regional-scale groundwater recharge: Contributions of a water budget approach in cold and humid climates. *Hydrology and Earth System Sciences*, 25(12), 6567–6589. <https://doi.org/10.5194/hess-25-6567-2021>
- Enemark, T., Peeters, L. J. M., Mallants, D., & Batelaan, O. (2019). Hydrogeological conceptual model building and testing: A review. *Journal of Hydrology*, 569, 310–329. <https://doi.org/10.1016/j.jhydrol.2018.12.007>
- Frick, H., Chow, F., Kuhn, M., Mahoney, M., Silge, J., & Wickham, H. (2022). rsample: General Resampling Infrastructure. Retrieved from <https://rsample.tidymodels.org>
- Friedman, J. H. (2001). Greedy function approximation: A gradient boosting machine. *Annals of Statistics*, 29(5), 1189–1232. <https://doi.org/10.1214/aos/1013203451>
- Giese, M., Haaf, E., Heudorfer, B., & Barthel, R. (2020). Comparative hydrogeology—Reference analysis of groundwater dynamics from neighbouring observation wells. *Hydrological Sciences Journal*, 65(10), 1685–1706. <https://doi.org/10.1080/02626667.2020.1762888>
- Green, T. R., Taniguchi, M., Kooi, H., Gurdak, J. J., Allen, D. M., Hiscock, K. M., et al. (2011). Beneath the surface of global change: Impacts of climate change on groundwater. *Journal of Hydrology*, 405(3–4), 532–560. <https://doi.org/10.1016/j.jhydrol.2011.05.002>
- Gribovski, Z., Szilágyi, J., & Kalicz, P. (2010). Diurnal fluctuations in shallow groundwater levels and streamflow rates and their interpretation—A review. *Journal of Hydrology*, 385(1–4), 371–383. <https://doi.org/10.1016/j.jhydrol.2010.02.001>
- G'Sell, M. G., Wager, S., Chouldechova, A., & Tibshirani, R. (2016). Sequential selection procedures and false discovery rate control. *Journal of the Royal Statistical Society Series B: Statistical Methodology*, 78(2), 423–444. <https://doi.org/10.1111/rssb.12122>
- Haaf, E., & Barthel, R. (2018). An inter-comparison of similarity-based methods for organisation and classification of groundwater hydrographs. *Journal of Hydrology*, 559, 222–237. <https://doi.org/10.1016/j.jhydrol.2018.02.035>

- Haaf, E., Giese, M., Heudorfer, B., Stahl, K., & Barthel, R. (2020). Physiographic and climatic controls on regional groundwater dynamics. *Water Resources Research*, 56(10), e2019WR026545. <https://doi.org/10.1029/2019wr026545>
- He, Y., Bárdossy, A., & Zehe, E. (2011). A review of regionalisation for continuous streamflow simulation. *Hydrology and Earth System Sciences*, 15(11), 3539–3553. <https://doi.org/10.5194/hess-15-3539-2011>
- Heudorfer, B., Haaf, E., Stahl, K., & Barthel, R. (2019). Index-based characterization and quantification of groundwater dynamics. *Water Resources Research*, 55(7), 5575–5592. <https://doi.org/10.1029/2018WR024418>
- Houben, T., Pujades, E., Kalbacher, T., Dietrich, P., & Attinger, S. (2022). From dynamic groundwater level measurements to regional aquifer parameters—Assessing the power of spectral analysis. *Water Resources Research*, 58(5), e2021WR031289. <https://doi.org/10.1029/2021wr031289>
- Hrachowitz, M., Savenije, H. H. G., Blöschl, G., McDonnell, J. J., Sivapalan, M., Pomeroy, J. W., et al. (2013). A decade of predictions in ungauged basins (PUB)—A review. *Hydrological Sciences Journal*, 58(6), 1198–1255. <https://doi.org/10.1080/02626667.2013.803183>
- Hughes, D. A., & Smakhtin, V. (1996). Daily flow time series patching or extension: A spatial interpolation approach based on flow duration curves. *Hydrological Sciences Journal*, 41(6), 851–871. <https://doi.org/10.1080/02626669609491555>
- Jackson, C. R., Wang, L., Pachocka, M., Mackay, J. D., & Bloomfield, J. P. (2016). Reconstruction of multi-decadal groundwater level time-series using a lumped conceptual model. *Hydrological Processes*, 30, 30–18. <https://doi.org/10.1002/hyp.10850>
- Kratzert, F., Klotz, D., Herrnegger, M., Sampson, A. K., Hochreiter, S., & Nearing, G. S. (2019). Toward improved predictions in ungauged basins: Exploiting the power of machine learning. *Water Resources Research*, 55(12), 11344–11354. <https://doi.org/10.1029/2019WR026065>
- Liu, Y., & Just, A. C. (2021). SHAPforxgboost: SHAP plots for 'XGBoost', R package version 0.1.1. Retrieved from <https://CRAN.R-project.org/package=SHAPforxgboost>
- Lóaiciga, H. A., & Leipnik, R. B. (2001). Theory of sustainable groundwater management: An urban case study. *Urban Water*, 3(3), 217–228. [https://doi.org/10.1016/s1462-0758\(01\)00040-1](https://doi.org/10.1016/s1462-0758(01)00040-1)
- Lundberg, S. M., Erion, G., Chen, H., Degraeve, A., Prutkin, J. M., Nair, B., et al. (2020). From local explanations to global understanding with explainable AI for trees. *Nature Machine Intelligence*, 2(1), 56–67. <https://doi.org/10.1038/s42256-019-0138-9>
- Mackay, J., Jackson, C. R., & Wang, L. (2014). A lumped conceptual model to simulate groundwater level time-series. *Environmental Modelling & Software*, 61, 229–245. <https://doi.org/10.1016/j.envsoft.2014.06.003>
- Marchant, B. P., & Bloomfield, J. P. (2018). Spatio-temporal modelling of the status of groundwater droughts. *Journal of Hydrology*, 564, 397–413. <https://doi.org/10.1016/j.jhydrol.2018.07.009>
- Maxe, L. (2013). Bedömningsgrunder för grundvatten. In *Sveriges geologiska undersökning SGU-rapport 2013, 1*.
- McDonnell, J. J., & Woods, R. (2004). On the need for catchment classification. *Journal of Hydrology*, 299(1), 2–3. [https://doi.org/10.1016/s0022-1694\(04\)00421-4](https://doi.org/10.1016/s0022-1694(04)00421-4)
- Mohamoud, Y. M. (2010). Prediction of daily flow duration curves and streamflow for ungauged catchments using regional flow duration curves. *Hydrological Sciences Journal*, 53(4), 706–724. <https://doi.org/10.1623/hysj.53.4.706>
- Montgomery, D. (2001). Slope distributions, threshold hillslopes, and steady-state topography. *American Journal of Science*, 301(4–5), 432–454. <https://doi.org/10.2475/ajs.301.4-5.432>
- Naghibi, S. A., Hashemi, H., Berndtsson, R., & Lee, S. (2020). Application of extreme gradient boosting and parallel random forest algorithms for assessing groundwater spring potential using DEM-derived factors. *Journal of Hydrology*, 589, 125197. <https://doi.org/10.1016/j.jhydrol.2020.125197>
- Peterson, T. J., & Fulton, S. (2019). Joint estimation of gross recharge, groundwater usage, and hydraulic properties within HydroSight. *Groundwater*, 57(6), 860–876. <https://doi.org/10.1111/gwat.12946>
- Rajaei, T., Ebrahimi, H., & Nourani, V. (2019). A review of the artificial intelligence methods in groundwater level modeling. *Journal of Hydrology*, 572, 336–351. <https://doi.org/10.1016/j.jhydrol.2018.12.037>
- R Development Core Team. (2022). *R: A language and environment for statistical computing*. R Foundation for Statistical Computing. Retrieved from <http://www.R-project.org>
- Ridolfi, E., Kumar, H., & Bárdossy, A. (2020). A methodology to estimate flow duration curves at partially ungauged basins. *Hydrology and Earth System Sciences*, 24(4), 2043–2060. <https://doi.org/10.5194/hess-24-2043-2020>
- Rinderer, M., McGlynn, B. L., & van Meerveld, H. J. (2017). Groundwater similarity across a watershed derived from time-warped and flow-corrected time series. *Water Resources Research*, 53(5), 3921–3940. <https://doi.org/10.1002/2016wr019856>
- Rinderer, M., Meerveld, H. J., & McGlynn, B. L. (2019). From points to patterns: Using groundwater time series clustering to investigate subsurface hydrological connectivity and runoff source area dynamics. *Water Resources Research*, 55(7), 5784–5806. <https://doi.org/10.1029/2018wr023886>
- Rinderer, M., van Meerveld, H. J., & Seibert, J. (2014). Topographic controls on shallow groundwater levels in a steep, prealpine catchment: When are the TWI assumptions valid? *Water Resources Research*, 50(7), 6067–6080. <https://doi.org/10.1002/2013wr015009>
- Rinderer, M., van Meerveld, I., Stähli, M., & Seibert, J. (2016). Is groundwater response timing in a pre-alpine catchment controlled more by topography or by rainfall? *Hydrological Processes*, 30(7), 1036–1051. <https://doi.org/10.1002/hyp.10634>
- Ruybal, C. J., Hogue, T. S., & McCray, J. E. (2019). Evaluation of groundwater levels in the Arapahoe aquifer using spatiotemporal regression kriging. *Water Resources Research*, 55(4), 2820–2837. <https://doi.org/10.1029/2018wr023437>
- Shu, C., & Ouara, T. B. M. J. (2012). Improved methods for daily streamflow estimates at ungauged sites. *Water Resources Research*, 48(2), 2523. <https://doi.org/10.1029/2011wr011501>
- Sivakumar, B., & Singh, V. P. (2012). Hydrologic system complexity and nonlinear dynamic concepts for a catchment classification framework. *Hydrology and Earth System Sciences*, 16(11), 4119–4131. <https://doi.org/10.5194/hess-16-4119-2012>
- Smakhtin, V. Y. (1999). Generation of natural daily flow time-series in regulated rivers using a non-linear spatial interpolation technique. *Regulated Rivers: Research & Management*, 15(4), 311–323. [https://doi.org/10.1002/\(sici\)1099-1646\(199907/08\)15:4<311::aid-rrr544>3.0.co;2-w](https://doi.org/10.1002/(sici)1099-1646(199907/08)15:4<311::aid-rrr544>3.0.co;2-w)
- Sugiyama, H., Vudhivanich, V., Whitaker, A. C., & Lorsirirak, K. (2003). Stochastic flow duration curves for evaluation of flow regimes in rivers. *JAWRA Journal of the American Water Resources Association*, 39(1), 47–58. <https://doi.org/10.1111/j.1752-1688.2003.tb01560.x>
- Taylor, J., & Tibshirani, R. J. (2015). Statistical learning and selective inference. *Proceedings of the National Academy of Sciences*, 112(25), 7629–7634. <https://doi.org/10.1073/pnas.1507583112>
- Varouchakis, E. A., Guardiola-Albert, C., & Karatzas, G. P. (2022). Spatiotemporal geostatistical analysis of groundwater level in aquifer systems of complex hydrogeology. *Water Resources Research*, 58(3), e2021WR029988. <https://doi.org/10.1029/2021wr029988>
- Vidon, P. (2012). Towards a better understanding of riparian zone water table response to precipitation: Surface water infiltration, hillslope contribution or pressure wave processes? *Hydrological Processes*, 26(21), 3207–3215. <https://doi.org/10.1002/hyp.8258>
- Vogel, R. M., & Fennessy, N. M. (1995). Flow duration curves II: A review of applications in water resources planning. *JAWRA Journal of the American Water Resources Association*, 31(6), 1029–1039. <https://doi.org/10.1111/j.1752-1688.1995.tb03419.x>

- Von Asmuth, J. R. (2012). Groundwater system identification through time series analysis.
- Voss, C. I. (2005). The future of hydrogeology. *Hydrogeology Journal*, *13*(1), 1–6. <https://doi.org/10.1007/s10040-005-0435-8>
- Wagener, T., Sivapalan, M., Troch, P., & Woods, R. (2007). Catchment classification and hydrologic similarity. *Geography Compass*, *1*(4), 901–931. <https://doi.org/10.1111/j.1749-8198.2007.00039.x>
- Wickham, H., Averick, M., Bryan, J., Chang, W., McGowan, L., François, R., et al. (2019). Welcome to the tidyverse. *Journal of Open Source Software*, *4*(43), 1686. <https://doi.org/10.21105/joss.01686>
- Winter, T. C. (2001). The concept of hydrologic landscapes. *Journal of the American Water Resources Association*, *37*(2), 335–349. <https://doi.org/10.1111/j.1752-1688.2001.tb00973.x>
- Wunsch, A., Liesch, T., & Broda, S. (2021). Groundwater level forecasting with artificial neural networks: A comparison of long short-term memory (LSTM), convolutional neural networks (CNNs), and non-linear autoregressive networks with exogenous input (NARX). *Hydrology and Earth System Sciences*, *25*(3), 1671–1687. <https://doi.org/10.5194/hess-25-1671-2021>
- Wunsch, A., Liesch, T., & Broda, S. (2022). Deep learning shows declining groundwater levels in Germany until 2100 due to climate change. *Nature Communications*, *13*(1), 1221. <https://doi.org/10.1038/s41467-022-28770-2>
- Yokoo, Y., & Sivapalan, M. (2011). Towards reconstruction of the flow duration curve: Development of a conceptual framework with a physical basis. *Hydrology and Earth System Sciences*, *15*(9), 2805–2819. <https://doi.org/10.5194/hess-15-2805-2011>