

Statistical Modeling of Water Quality Data on River Networks

2025 Workshop on Spatial Statistics and Related Fields

2025-07-24

Seoncheol Park

Department of Mathematics

Hanyang University

 [pscstat at hanyang.ac.kr](mailto:pscstat@hanyang.ac.kr)

Joint work with Joonpyo Kim (Sejoing Univ.), Hyungryul Park (Yonsei Univ.),
Yeonje Lee (Sejong Univ.), and Seungyeon Lim (Hanyang Univ.)



Contents

Topics

- Water Quality Data on River Networks
- Expectile-based Probabilistic Forecasting
- Adaptive Boosting on River Networks

Related Works

- Park, H., Kim, J., & **Park, S.** (2025+). Expectile-based Probabilistic Forecasting for Spatio-Temporal River Network Data. *Under Revision*.
- Lim, S. & **Park, S.** (2025+). Adaptive Boosting on Linear Networks. *In Preparation*.



Water Quality Data on River Networks

Geum River Networks

- Images from (Park & Oh, 2022)

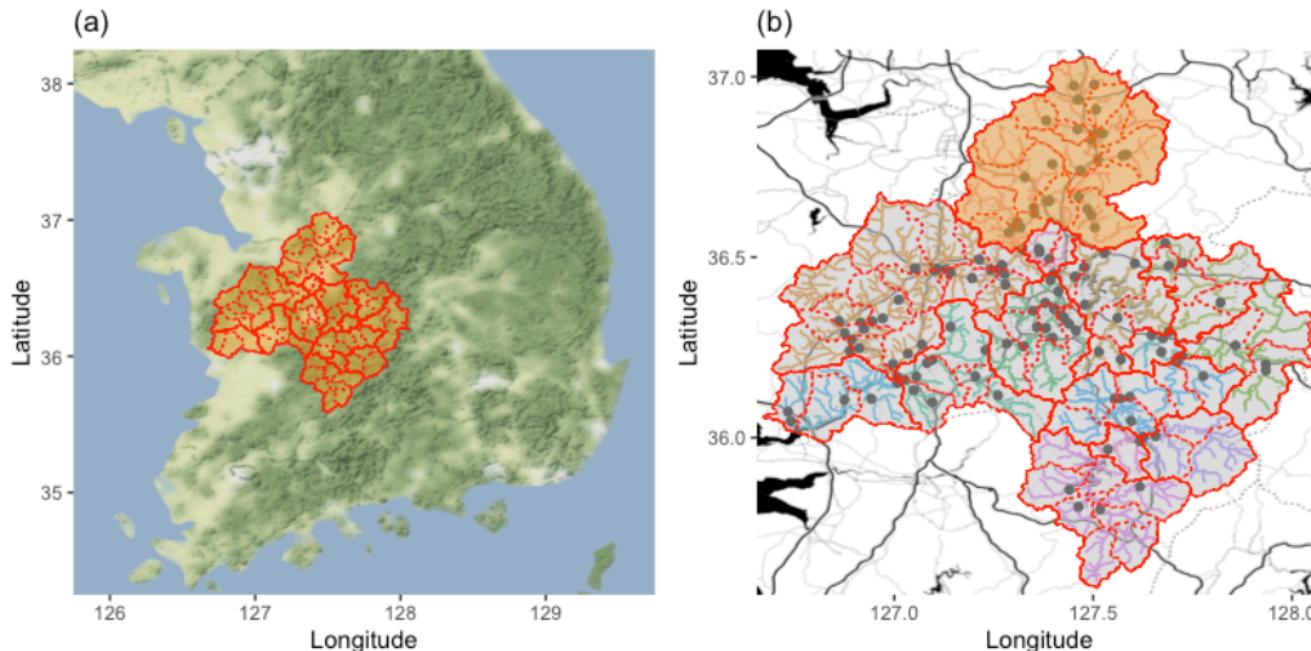


Figure 1: (a) The Geum River basin and (b) Data locations in the basin.

Water Quality Data

내년부터 유기물질 측정지표 바뀐다

환경정책 설명회 개최…총유기탄소 지표 도입

2019.11.11 11:12

정현섭 객원기자



환경부는 내년부터 수질 유기물질 측정지표로서 총유기탄소(TOC, Total Organic Carbon)를 도입할 계획이다.

지난 8일 라마다 서울 신도림 호텔에서 개최된 '환경정책설명회 및 최신기술 발표·전시회'에서 환경부는 이같이 밝혔다.

이날 행사에서 환경부 및 환경 기관, 민간기업 등의 실무자들은 정책 현안과 변화에 대해 설명하고, 환경산업체들은 최신 공법 및 기술에 대해 소개했다.

수질의 유기물질 측정지표, 내년부터 TOC로 전환

하수나 폐수에 포함된 다량의 유기물질들이 처리되지 않은 채로 방류되면 공공수역의 수질을 악화시키기 때문에 지속적인 측정을 통한 점검이 필요하다.

대표적인 유기물질 측정지표로는 BOD(생화학적 산소요구량), COD(화학적 산소요구량으로 망간(Mn)이나 크롬(Cr)을 산화제로 이용), TOC 등이 사용된다.

현재 국내의 물환경보전법에서는 BOD와 COD(Mn)를 적용하고 있는데, 환경부는 2020년부터 유기물질 측정지표 COD(Mn)을 TOC로 전환할 계획이다.

Figure 2: The Korean Ministry of Environment introduced Total Organic Carbon (TOC) as an indicator for measuring organic substances in water quality starting in 2020.

- The Korean Water Environment Information System

- Water quality index for organic compounds:

- Biochemical Oxygen Demand (BOD)
- Chemical Oxygen Demand (COD)
- Total Organic Carbon (TOC)

- Water quality index for *algal bloom*:

- Total Nitrogen (TN)
- Total Phosphorus (TP)



Water Quality Data (cont.)

- Seasonality, irregularly observed time series, outliers or extreme values, ...

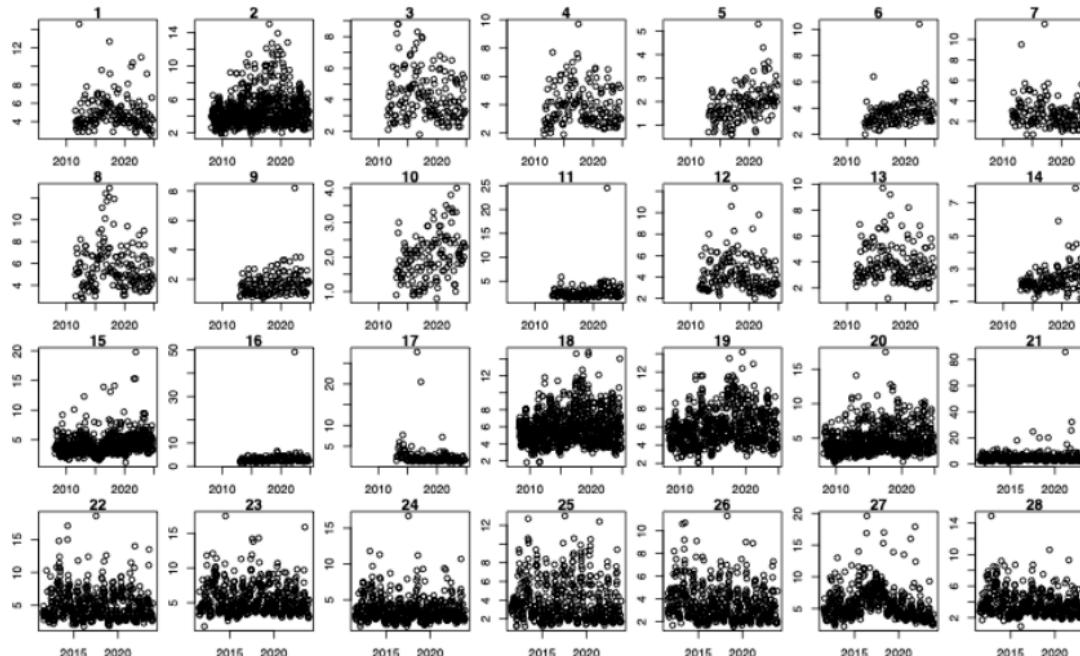


Figure 3: Scatterplot of TOC time series (2005~2024) at the Miho River.

Expectile-based Probabilistic Forecasting

Data: Miho River

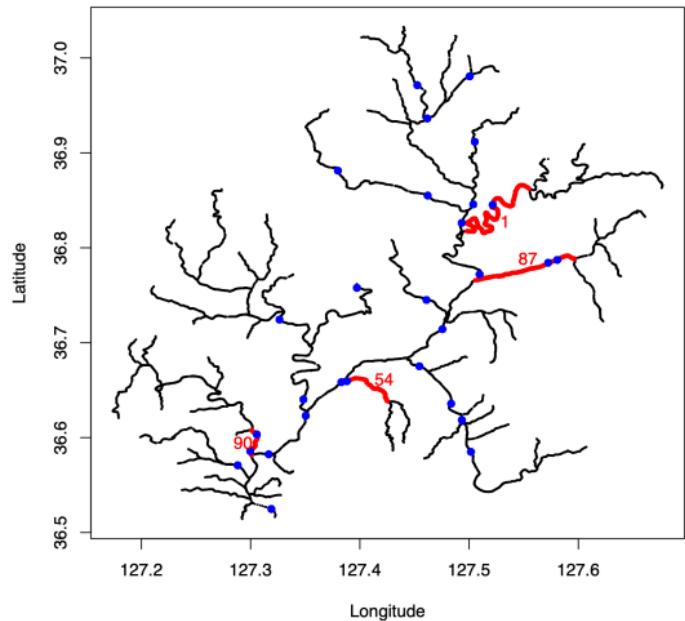


Figure 4: Miho River network structure.

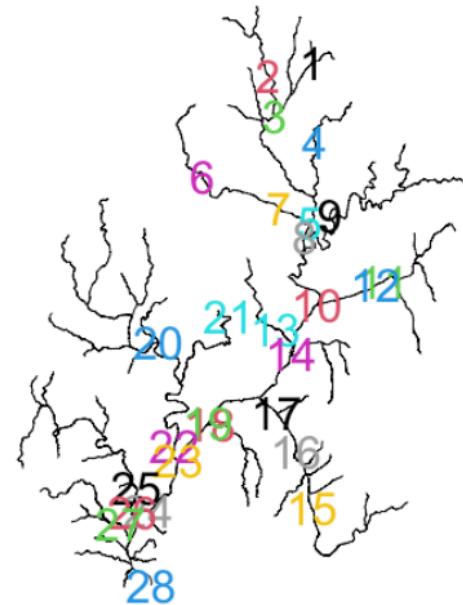


Figure 5: Water Quality Observation Sites in the Miho River.

Data: Miho River (cont.)

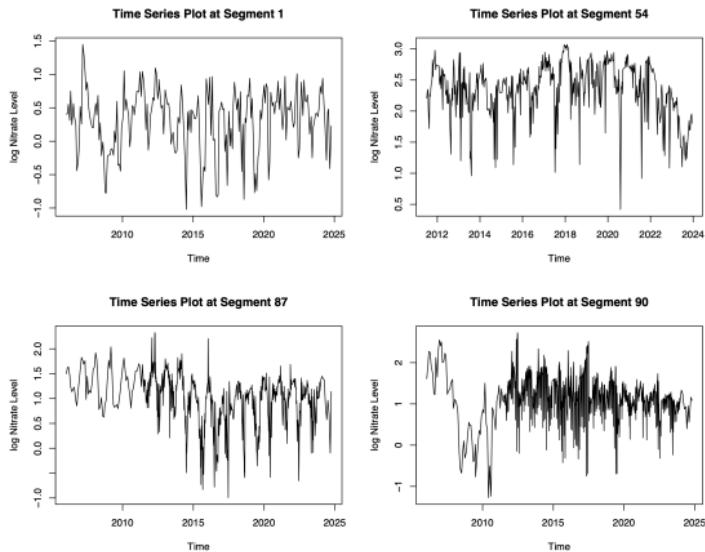


Figure 6: Time series plots of the log total nitrogen (TN) at selected segments.

- The dataset has following characteristics, making the forecasting problem more challenging:
 - ▶ Observation time points are irregular and their ranges differ across sites.
 - ▶ The observed time series show heteroskedasticity and seasonality.
 - ▶ Their variability differs across regions, and some outliers are present.

Related Works

- (O'Donnell et al., 2014) proposed a flexible regression model which extends a spatio-temporal geoadditive model with spatial components defined as functions of stream segments. (**TN**)
- (Gallacher et al., 2017) and (Kim et al., 2022) suggested a flow-adaptive principal component analysis for river network data. (**TN, TOC**)
- (Park & Oh, 2022) proposed a nonparametric regression model based on a lifting scheme for the statistical modeling of TOC data in the Geum River network. (**TOC**)
- (Santos-Fernandez et al., 2022) explored some Bayesian models for stream networks. (**Water temperature**)

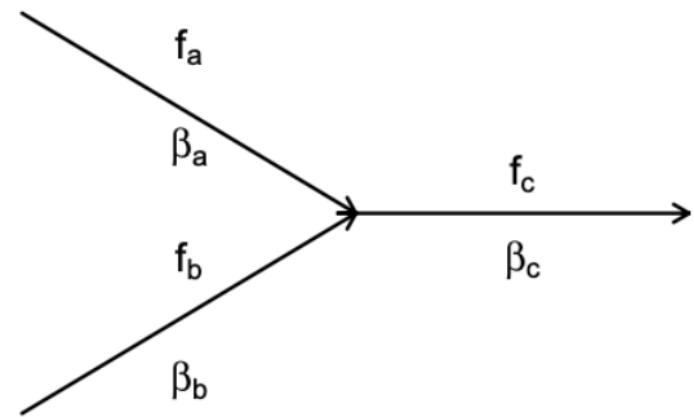


Figure 7: A stochastic representation of a confluence, with model parameters (β_a, β_b) , flows (f_a, f_b) and the corresponding outgoing versions (β_c, f_c) .

Quantile and Expectile Regression

- For $\tau \in (0, 1)$, τ -th quantile $q_\tau(Y)$ of a real random variable Y can be obtained by

$$q_\tau(Y) = \operatorname{argmin}_q E[\rho_\tau(Y - q)], \quad (1)$$

solving the above optimization problem, where $\rho_\tau(x) = x(\tau - I(x < 0))$ is a check function. In this view, (Koenker & Bassett, 1978) suggested a **quantile regression** to represent a conditional τ -th quantile of response variable Y as a function $f(X)$ of explanatory variable X :

$$\hat{q}_\tau = \operatorname{argmin}_{f \in \mathcal{F}} \sum_{i=1}^n \rho_\tau(y_i - f(x_i)). \quad (2)$$

- (Newey & Powell, 1987) proposed a computationally attractive alternative, called **expectile regression**, replacing a check loss function with asymmetric L^2 loss. A conditional τ -th expectile $\hat{q}_\tau(X)$ of a response variable Y given X can be estimated by

$$\hat{e}_\tau = \operatorname{argmin}_{f \in \mathcal{F}} \sum_{i=1}^n \gamma_\tau(x)(y_i - f(x_i)), \quad (3)$$

minimizing sum of asymmetric squared residuals, where $\gamma_\tau(x) = |x|^2|\tau - I(x < 0)|$.

Spatio-Temporal Additive Model

- (O'Donnell et al., 2014) proposed a flexible regression model to consider the unique spatial structure of data which arise from a river network implementing **P-spline**. (Eilers & Marx, 1996)
- They assumed that the average response observed at location s on z -th day in year t could be represented as a sum of each component and their interaction terms

$$y = c + m_s(s) + m_z(z) + m_t(t) + m_{s,z}(s, z) + m_{s,t}(s, t) + m_{z,t}(z, t) + \varepsilon. \quad (4)$$

- Each component was estimated by a weighted sum of B-spline basis functions:
 - B-spline of order 0 (and therefore piecewise constant) for spatial component m_s ;
 - cubic B-splines for m_z and $m - t$;
 - and their tensor products for interaction terms.
- Confluence penalty: If streams a and b , at which value of spatial component is β_a and β_b , join with each other and become stream c , with spatial component value β_c , confluence penalty is

$$w_a^2(\beta_a - \beta_c)^2 + w_b^2(\beta_b - \beta_c)^2. \quad (5)$$

- Here $w_a = \frac{f_a}{f_c}$, $w_b = \frac{f_b}{f_c}$ and f_a, f_b, f_c are flows at stream a, b , and c , respectively.
- Penalty at stream flows from segment a to b without confluence is defined as $(\beta_a - \beta_b)^2$.

The Proposed Method

- In this work, we extend the approach to model spatio-temporal data observed along the stream network proposed by (O'Donnell et al., 2014) combining with expectile regression.
- We assume that τ -th expectile of response is represented as a sum of three spatio-temporal components and their interaction terms:

$$e_\tau(y|s, z, t) = c_\tau + m_{s,\tau}(s) + m_{z,\tau}(z) + m_{t,\tau}(t) + m_{sz,\tau}(s, z) + m_{st,\tau}(s, t) + m_{zt,\tau}(z, t). \quad (6)$$

- Denoting the observed response y_i at station s_i on z_i -th day in year t_i for $i = 1, \dots, n$, each component is estimated by minimizing following asymmetrically squared sum of residuals:

$$\sum_{i=1}^n w_i(\tau)(y_i - f_\tau(s_i, z_i, t_i))^2, \quad (7)$$

where

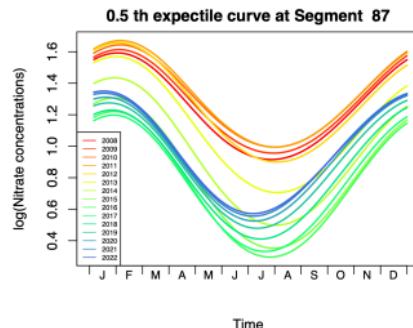
$$w_i(\tau) = \begin{cases} \tau & \text{if } y_i \geq f_\tau(s_i, z_i, t_i) \\ 1 - \tau & \text{if } y_i < f_\tau(s_i, z_i, t_i) \end{cases}. \quad (8)$$

Note that equation (6) can be summarized via design matrix of B-spline basis.

The Proposed Method (cont.)

A forecasting method based on the expectile smoothing

- Regarding $\hat{e}_\tau(y|s, z, t) = f_t(z|s, \tau)$ as a function of z , we concentrate of the problem of curve forecasting at $(T + h)$ -th year based on curves $f_1(z|s, \tau), \dots, f_T(z|s, \tau)$.
- For each station s and expectile level τ , we apply vector autoregressive (VAR) model to FPCs of $f_t(z|s, \tau)$ ($t = 1, \dots, T$) to get h -step-ahead prediction $\hat{f}_{T+h}(z|s, \tau)$. (Aue et al., 2015)
- Repeating the procedure for all values $\tau \in (0, 1)$, we can predict an expectile process $\tau \mapsto \hat{f}_{T+h}(z|s, \tau)$ of the response at location s on z -th day, which can be transformed to the distribution function by (Waltrup et al., 2014) .



The Proposed Method (cont.)

- We summarize a detailed procedure for the forecast:
 1. For a given set \mathcal{T} of expectile levels, fit the model (6) to obtain $\hat{e}_\tau(y|s, z, t) = \hat{f}_t(z|s, \tau)$ for each $\tau \in \mathcal{T}$.
 2. For each observation site s :
 1. For each $\tau \in \mathcal{T}$, conduct functional principal component analysis on $\{\hat{f}_t(z|s, \tau)\}_{t=1}^T$ to obtain FPCs $\gamma_{t,\ell}(s, \tau)$ and eigenfunction $\psi_\ell(z|s, \tau)$ ($\ell = 1, \dots, L$).
 2. Get h -step-ahead forecast of FPC, $(\gamma_{T+h,\ell}(s, \tau))_{\ell=1}^L$, using VAR model.
 3. Reconstruct h -step-ahead forecast of expectile curve $\hat{f}_{T+h}(z|s, \tau)$ as $\hat{f}_{T+h}(z|s, \tau) = \sum_{\ell=1}^L \gamma_{T+h,\ell}(s, \tau) \psi_\ell(z|s, \tau)$.
 3. Predict a distribution function of the response on z -th day of $(T + h)$ -th year at observation site s as follows:
 1. Sort $(\hat{f}_{T+h}(z|s, \tau))_{\tau \in \mathcal{T}}$ in increasing order; denote them as $(\hat{f}_{T+h}^*(z|s, \tau))_{\tau \in \mathcal{T}}$.
 2. Interpolate $(\hat{f}_{T+h}^*(z|s, \tau))_{\tau \in \mathcal{T}}$ linearly and obtain an expectile process $\tau \mapsto \hat{f}_{T+h}^*(z|s, \tau)$, and then estimate a distribution function.

Practical Details

- **Tuning regularization parameters:** We select the tuning parameters λ s by minimizing Schwarz Information Criterion. To lower computational burden, we iteratively update parameter values rather than evaluating all candidates.
- **Choosing a set \mathcal{T} of expectile levels:** We use

$$\mathcal{T} = \{0.01, 0.05, 0.1, 0.2, 0.3, \dots, 0.9, 0.95, 0.99\}. \quad (9)$$

- **Selection of the number of FPCs and VAR order:** We use 3 FPCs, explaining more than 99% of variation, and AIC and Schwarz criterion are used to determine the order of the VAR model.
- **Performance evaluation** We evaluate the continuous ranked probability score (CRPS). Denoting the observation at location and time point of interest as X , and its predicted distribution function as \hat{F} , we calculate the CRPS as follows:

$$\text{CRPS} := \sum_{k=1}^K \left(\hat{F}(x_k) - I(X \leq x_k) \right)^2, \quad (10)$$

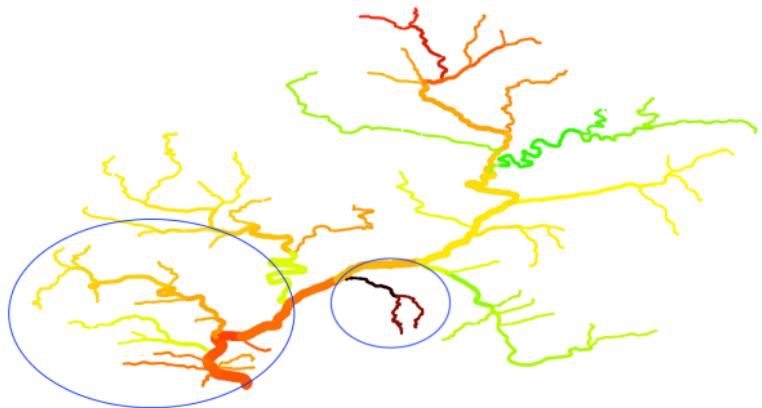
where $(x_1, \dots, x_K)^T$ represents a suitable vector of evaluation points.

Result: Spatial Components

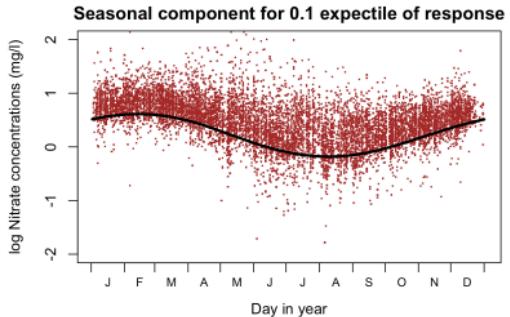
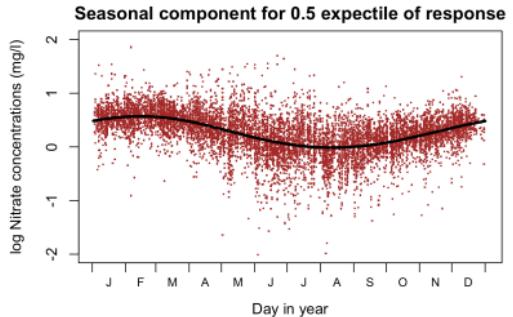
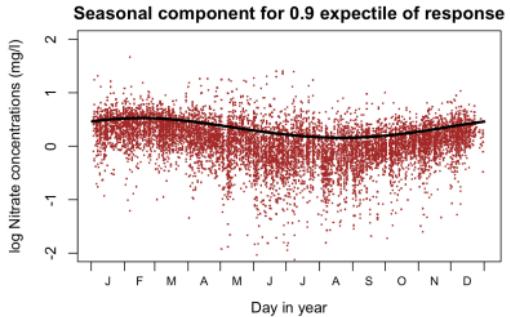
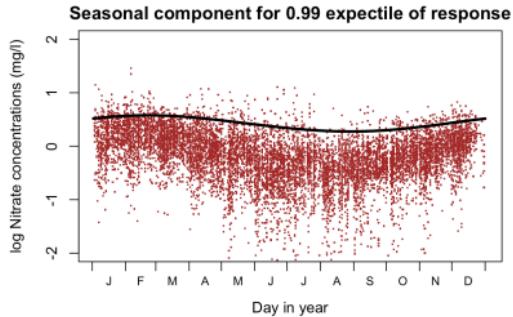
Spatial Component for 50% expectile



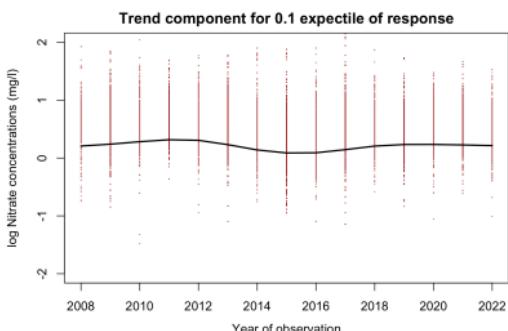
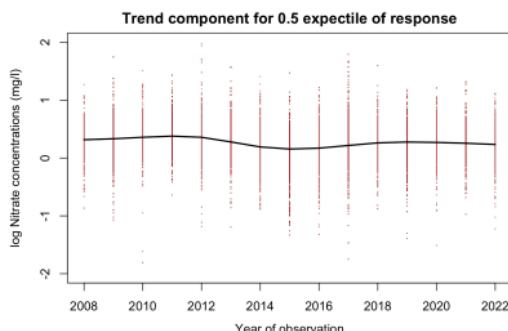
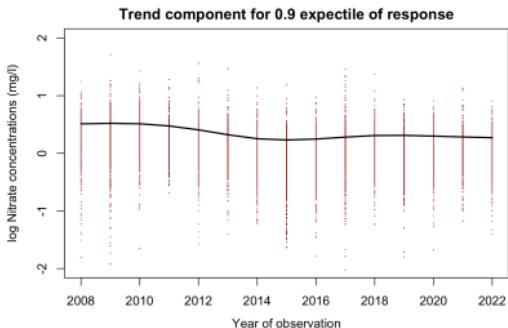
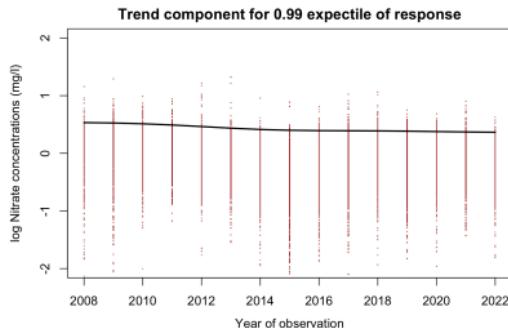
Spatial Component for 90% expectile



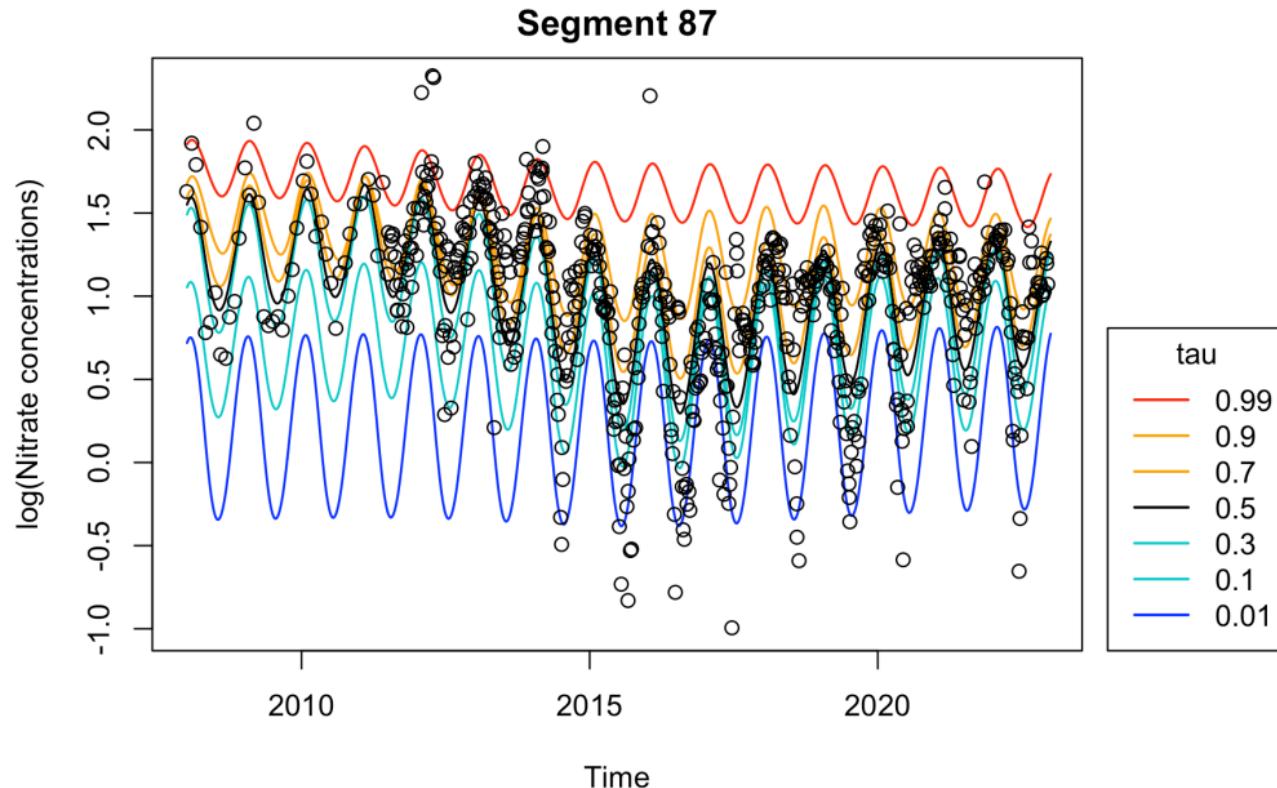
Result: Seasonal Components



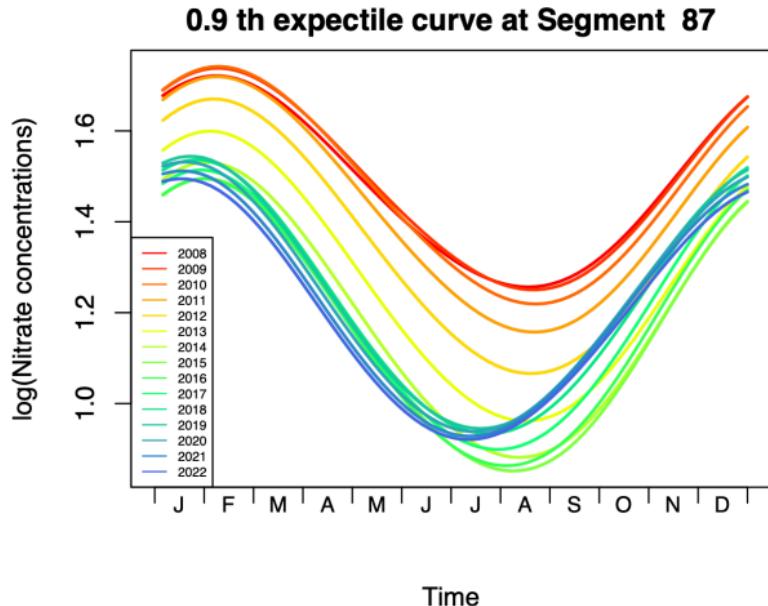
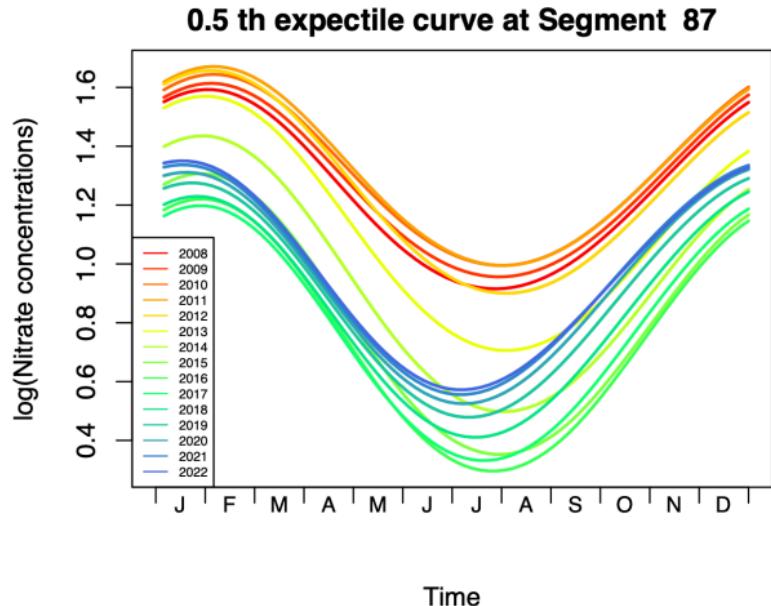
Result: Trend Components



Result: Estimated Curves at Segment 87



Results: Annual Expectile Curves



Performance Evaluation

- We compare CRPS score with benchmarks to validate performance, which are defined as follows:
 - **Benchmark 1:** For each segment, yield point estimate as an average of all training data, and predict as the averaged value regardless of the date. Then estimate the distribution as having point mass on the averaged value.
 - **Benchmark 2:** For each segment, yield probabilistic forecast as a quantile of all training data, regardless of the date. Then distribution function can be estimated as an inverse of quantile process.
 - **Benchmark 3:** Same as the proposed method, but use only the predicted value of mean, and estimate the distribution as having point mass on the mean value.

Performance Evaluation (cont.)

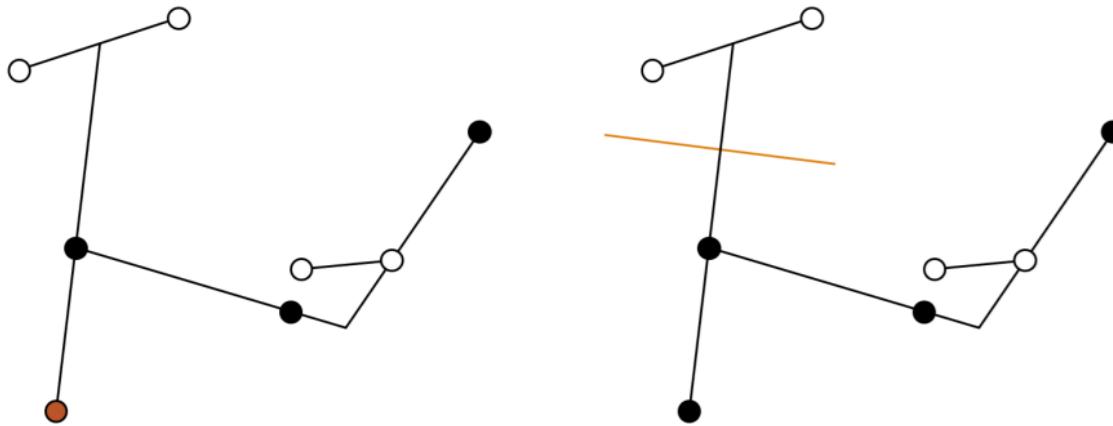
	1 year ahead prediction Average CRPS score	2 year ahead prediction Average CRPS score
Proposed method	0.0422 (0.0371)	0.0385 (0.0309)
Benchmark 1	0.0827 (0.0628)	0.0743 (0.0585)
Benchmark 2	0.0584 (0.0433)	0.0531 (0.0372)
Benchmark 3	0.0565 (0.0485)	0.0496 (0.0404)

Table 1: Average CRPS scores of the proposed method and benchmarks with their standard deviations in parentheses.



Adaptive Boosting on River Networks

Motivation & Linear Networks



- A **Linear network** L is an union of finite collection of line segments in a plane, i.e., $L = \cup_{i=1}^n l_i$, where l_1, \dots, l_n are n line segments in a plane.
- In this study, we deal with connected linear networks with no cycle.
- We assume the data was generated by a point process on linear network.
- We propose adaptive boosting with decision tree, which uses the linear network structure as an explanatory variable.

Related Works

- (Baddeley et al., 2014) modeled spatially varying spine density using the relative distribution and regression trees.
- (Ver Hoef, 2018) found examples where the positive definite covariance assumption does not hold in linear network data.

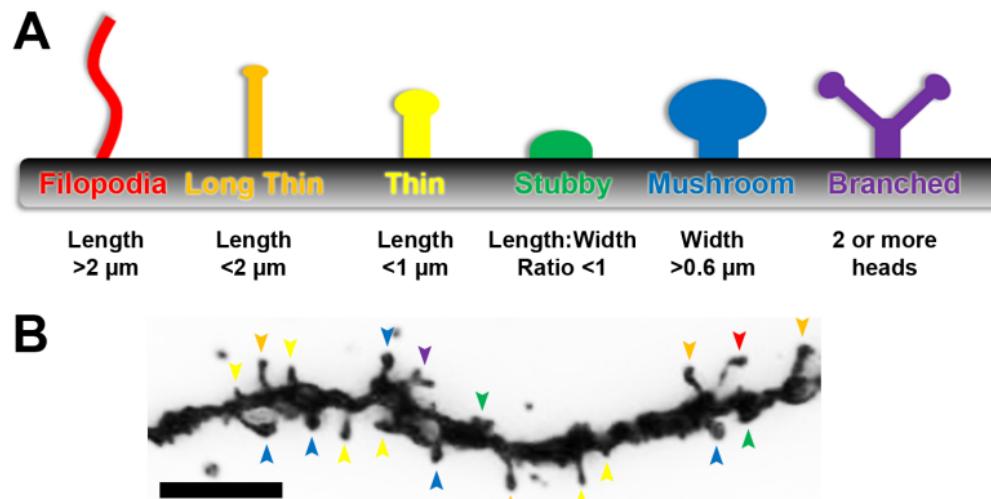
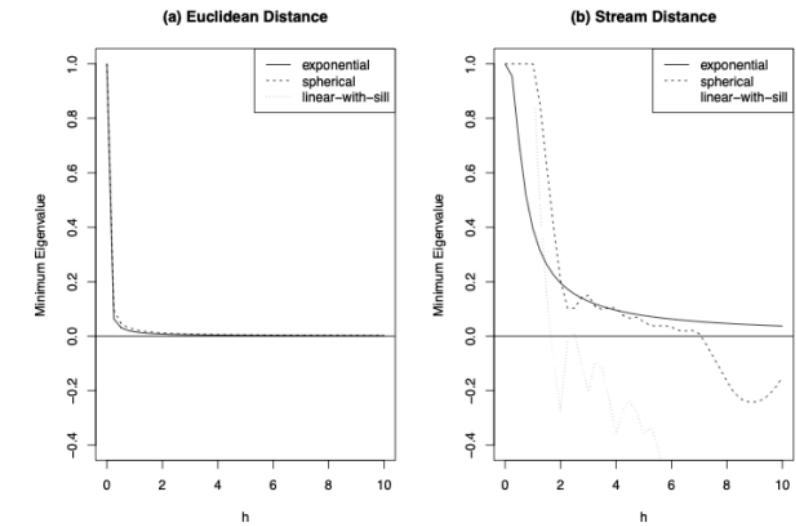
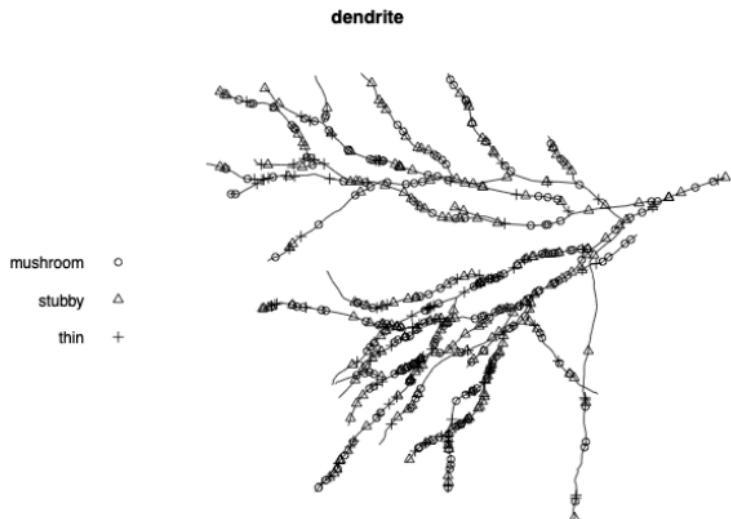


Figure 8: Different dendrite types.

Related Works (cont.)

- Dendrite data: 566 spines observed on one branch of the dendritic tree of a rat neuron (in R spatstat.linnet package)



- Images from Ver Hoef (2018)

Adaptive Boosting Algorithm

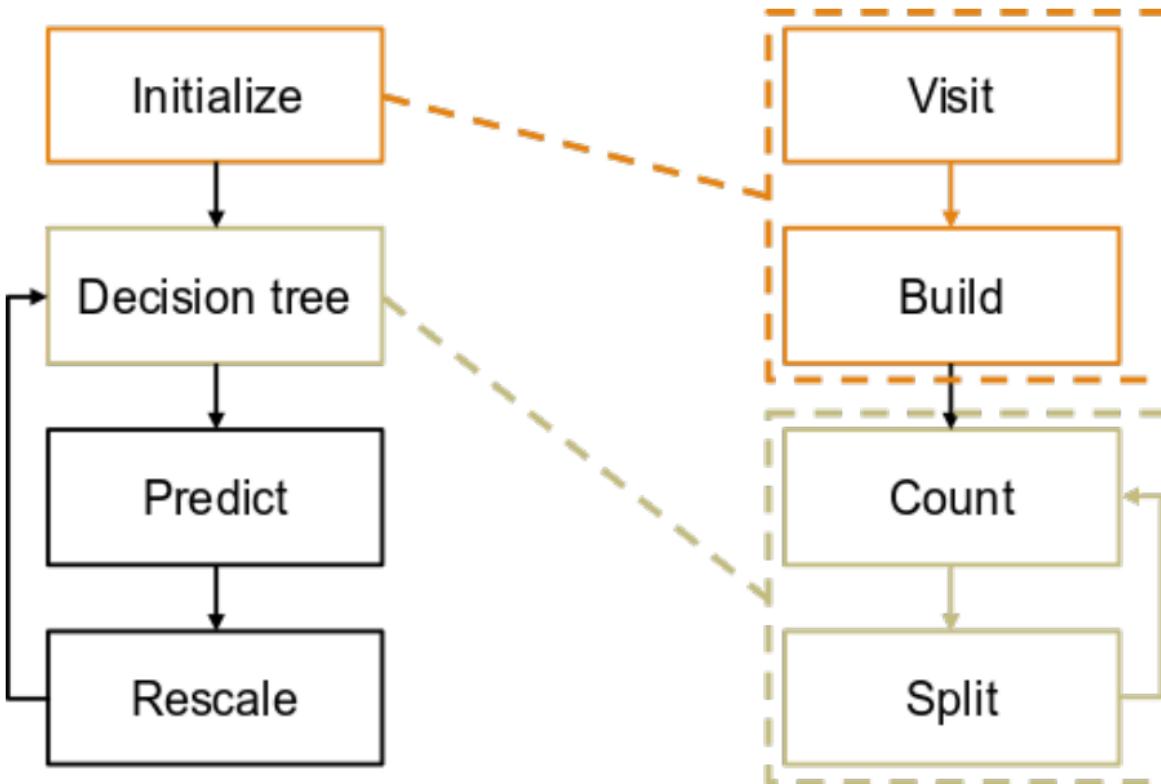


Figure 9: An outline of the splitting algorithm used in the proposed method.

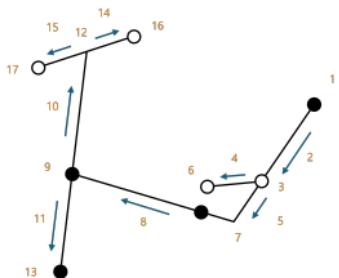
Adaptive Boosting Algorithm (cont.)

- Adaptive boosting with a decision tree has four phases:

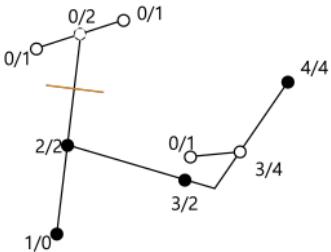
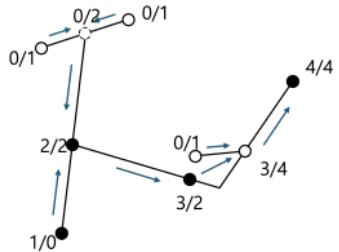
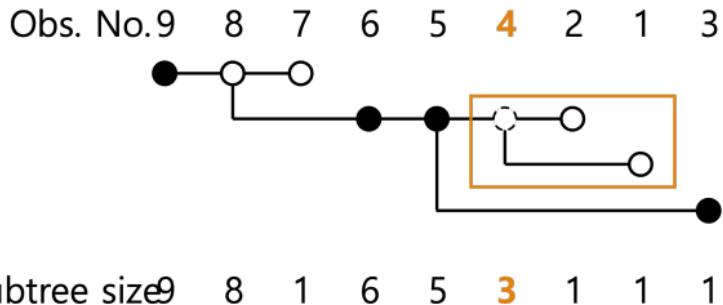
1. **Initialize** phase: We assign equal weights $w_i = \frac{1}{N}$ on each observation.
 1. **Visit** subphase: We select a vertex on the linear network (called the **root**), then traverse the linear network with breadth-first search, including both vertices and edges.
 2. **Build** subphase: We arrange the observations in a tree structure, preserving the adjacency on the linear network.
2. **Decision tree** phase: We construct a **decision tree** using the weights.
 1. **Count** subphase: We need to compute the sum of weights on a subtree of the observation tree.
 2. **Split** subphase: A **split point** is chose from any observation on the observation tree, then the split of the observation tree is obtained by deleting an edge between the split point and its parental observation.
3. **Predict** phase: We obtain the fitted values of the data with the decision tree.
3. **Rescale** phase: We assign less weight on correctly classified observations and more weight on misclssified ones.

Adaptive Boosting Algorithm (cont.)

- Visit, build, count, and split subphase

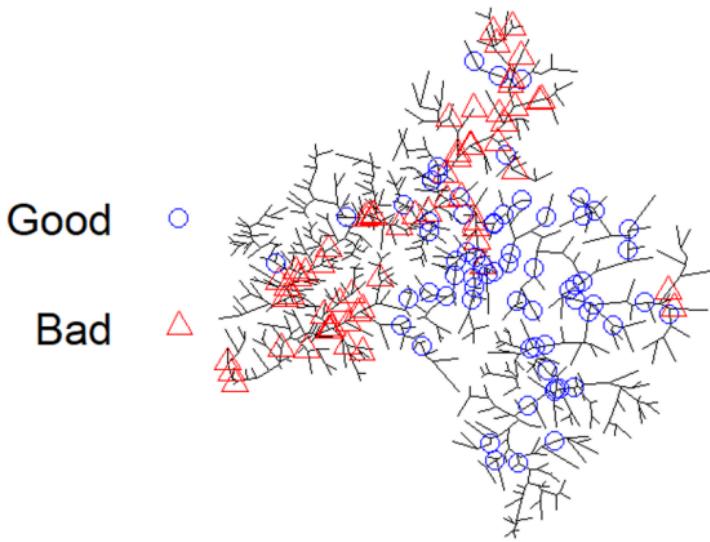
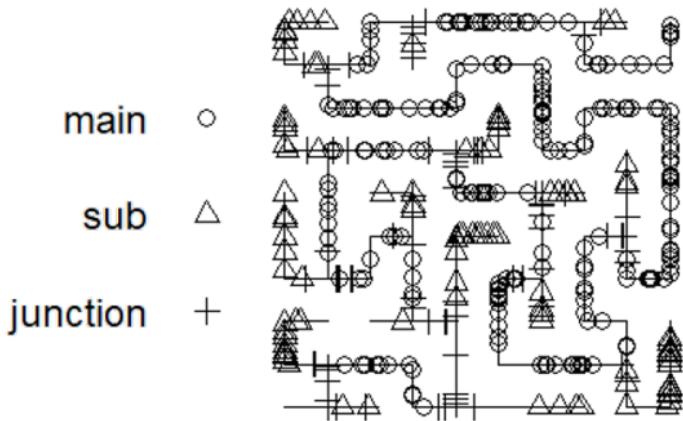


- Predict subphase



Data Analysis

- (Left) A realization of “maze” dataset and (Right) Geum River dataset.



Data Analysis (cont.)

Maze

- Setting: 3-class classification
- Dataset: Linear network with 99 line segments, 396 observations with label “main”, “sub”, “junction”, 100 datasets generated independently with the same rate
- Evaluation: paired differences of leave-one-out cross validation (LOOCV)

Number of Iterations	5	10	20	30	50	75	100
Mean accuracy: proposed	67.6%	72.1%	75.6%	78.1%	79.6%	80.6%	81.5%
Mean accuracy: contrast	64.0%	68.0%	73.4%	75.5%	77.4%	78.5%	78.9%

Table 2: Performances of two methods by the number of iterations.

Data Analysis (cont.)

Geum-River

- Dataset: River network representing Geum-River, 129 observations with label “Good” and “Bad”, classified by TOC
- Water quality “Good” if $\text{TOC} \leq 4(\text{mg/L})$, “Bad” if $\text{TOC} > 4$
- TOCs are measured in June 2022

Number of Iterations	5	10	15	20
Accuracy: proposed	80.6%	84.5%	83.7%	82.2%
Accuracy: contrast	79.8%	78.2%	79.8%	79.8%

Table 3: Accuracies of two methods by the number of iterations.

- The contrast method has weakness in predicting the main stream quality.
- The proposed method has weakness in predicting the branch quality.

Data Analysis (cont.)

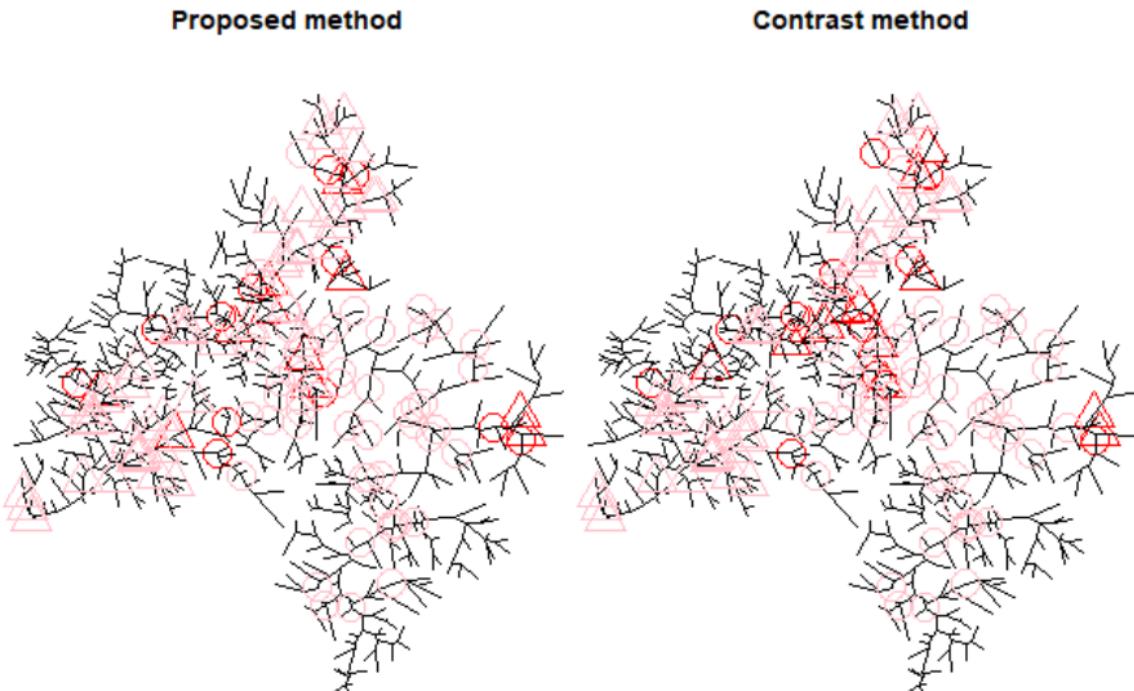


Figure 10: Correctly classified (pink) and misclassified (red) observations by proposed (Left) and contrast (Right) methods.

Conclusion

Summary

Expectile-based Probabilistic Forecasting

- The proposed method combines a nonparametric spatio-temporal additive model and expectile regression, allowing it to account for the unique structure of the river network.
- By adopting a roughness penalty, we obtain smooth curves in both spatial and temporal domains, with low computational burden.
- By integrating forecasting methods developed on the area of functional data analysis, the proposed method can generate expectile curves several years ahead at each segment.

Adaptive Boosting on Linear Networks

- We propose an adaptive boost algorithm using this decision tree on a linear network.
- We showed that our proposed method has better accuracy than the well-known method in datasets “maze” and “Geum River”.

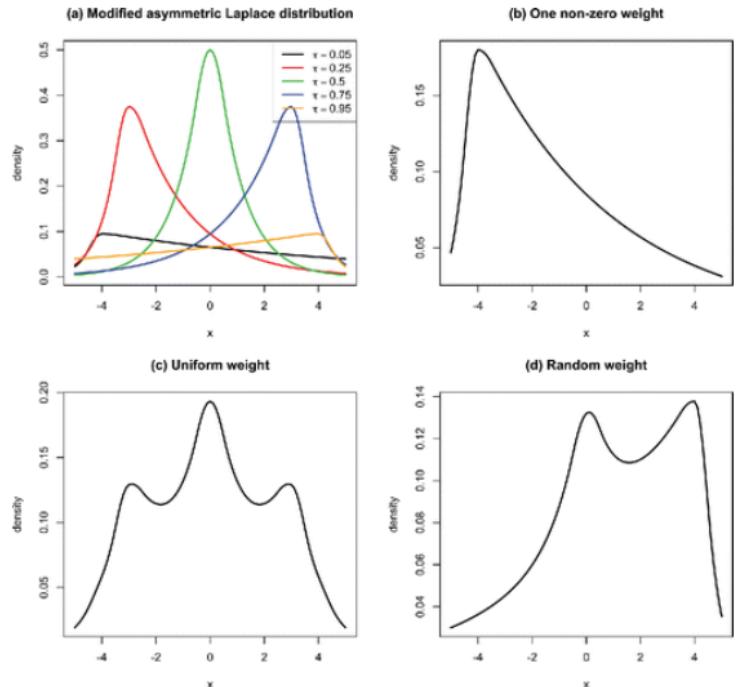
Further Works

Quantile Regression

- Spatio-temporal regression for the river network data can be combined with **quantile regression**.
- As suggested by (Franco-Villoria et al., 2018), we can replace asymmetric L^2 loss with a Huber loss

$$\rho_{\tau,c}(u) = \begin{cases} (\tau - 1)(u + 0.5c) & u < -c \\ \frac{0.5(1-\tau)u^2}{c} & -c \leq u < 0 \\ \frac{0.5\tau u^2}{c} & 0 \leq u < c \\ \tau(u - 0.5c) & c \leq u \end{cases} \quad (11)$$

- The function $\rho_{\tau,c}(u)$ with $c = 1.345$ and $\tau = 0.5$ is equivalent to the Huber loss function. (Huber, 1964; Lim & Oh, 2016)



- Images from Lim & Oh (2016)

Further Works (cont.)

Multivariate Analysis of Water Quality Data

- Conditional multivariate extremes (Heffernan & Tawn, 2004)

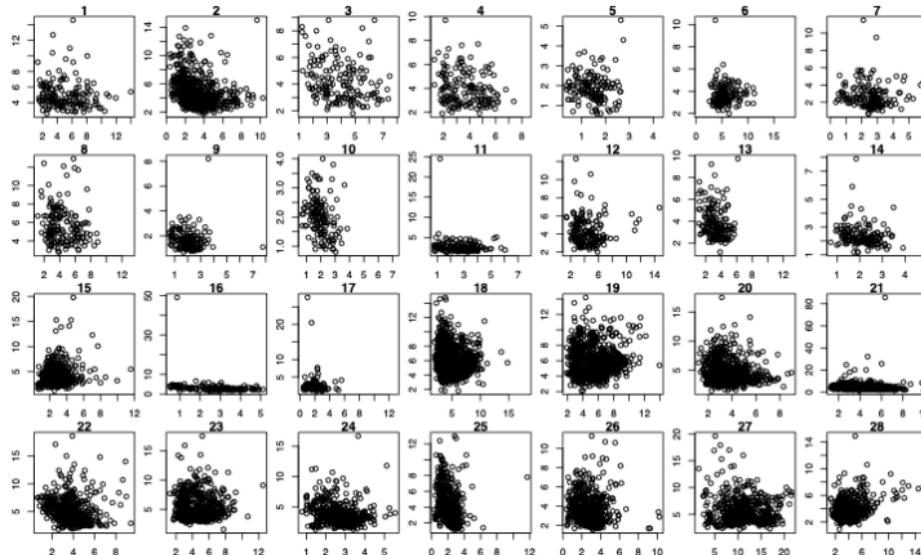


Figure 11: Scatterplot of TN (x-axis) vs. TOC (y-axis) at the Miho River.



Thank You!



Bibliography

- Aue, A., Norinho, D. D., & Hörmann, S. (2015). On the Prediction of Stationary Functional Time Series. *Journal of the American Statistical Association*, 110(509), 378–392. <https://doi.org/10.1080/01621459.2014.909317>
- Baddeley, A., Jammalamadaka, A., & Nair, G. (2014). Multitype Point Process Analysis of Spines on the Dendrite Network of a Neuron. *Journal of the Royal Statistical Society Series C: Applied Statistics*, 63(5), 673–694. <https://doi.org/10.1111/rssc.12054>
- Eilers, P. H. C., & Marx, B. D. (1996). Flexible smoothing with B-splines and penalties. *Statistical Science*, 11(2), 89–121. <https://doi.org/10.1214/ss/1038425655>
- Franco-Villoria, M., Scott, M., & Hoey, T. (2018). Spatiotemporal modeling of hydrological return levels: A quantile regression approach. *Environmetrics*, 30(2). <https://doi.org/10.1002/env.2522>
- Gallacher, K., Miller, C., Scott, E. M., Willows, R., Pope, L., & Douglass, J. (2017). Flow-directed PCA for monitoring networks. *Environmetrics*, 28(2), e2434. <https://doi.org/10.1002/env.2434>

Bibliography (cont.)

- Heffernan, J. E., & Tawn, J. A. (2004). A conditional approach for multivariate extreme values (with discussion). *Journal of the Royal Statistical Society: Series B (Statistical Methodology)*, 66(3), 497–546. <https://doi.org/10.1111/j.1467-9868.2004.02050.x>
- Huber, P. J. (1964). Robust Estimation of a Location Parameter. *The Annals of Mathematical Statistics*, 35(1), 73–101. <https://doi.org/10.1214/aoms/1177703732>
- Kim, K., Oh, H., & Park, M. (2022). Principal component analysis for river network data: Use of spatiotemporal correlation and heterogeneous covariance structure. *Environmetrics*, 34(4). <https://doi.org/10.1002/env.2753>
- Koenker, R., & Bassett, G. (1978). Regression Quantiles. *Econometrica*, 46(1), 33–50. <https://www.jstor.org/stable/1913643>
- Lim, Y., & Oh, H.-S. (2016). A Data-Adaptive Principal Component Analysis: Use of Composite Asymmetric Huber Function. *Journal of Computational and Graphical Statistics*, 25(4), 1230–1247. <https://doi.org/10.1080/10618600.2015.1067621>

Bibliography (cont.)

- Newey, W. K., & Powell, J. L. (1987). Asymmetric Least Squares Estimation and Testing. *Econometrica*, 55(4), 819. <https://doi.org/10.2307/1911031>
- O'Donnell, D., Rushworth, A., Bowman, A. W., Scott, E. M., & Hallard, M. (2014). Flexible regression models over river networks. *Journal of the Royal Statistical Society: Series C (Applied Statistics)*, 63(1), 47–63. <https://doi.org/10.1111/rssc.12024>
- Park, S., & Oh, H.-S. (2022). Lifting scheme for streamflow data in~river networks. *Journal of the Royal Statistical Society: Series C (Applied Statistics)*, 71(2), 467–490. <https://doi.org/10.1111/rssc.12542>
- Santos-Fernandez, E., Ver Hoef, J. M., Peterson, E. E., McGree, J., Isaak, D. J., & Mengersen, K. (2022). Bayesian spatio-temporal models for stream networks. *Computational Statistics & Data Analysis*, 170, 107446. <https://doi.org/10.1016/j.csda.2022.107446>
- Ver Hoef, J. M. (2018). Kriging models for linear networks and non-Euclidean distances: Cautions and solutions. *Methods in Ecology and Evolution*, 9(6), 1600–1613. <https://doi.org/10.1111/2041-210x.12979>

Bibliography (cont.)

Waltrup, L. S., Sobotka, F., Kneib, T., & Kauermann, G. (2014). Expectile and quantile regression—David and Goliath?. *Statistical Modelling*, 15(5), 433–456. <https://doi.org/10.1177/1471082x14561155>