Motivation
Time series segmentation
Outlier Visualization for Barium
Other project ideas
Discussion

# Statistical Analysis of Surface Water Data

## Time series segmentation and Outlier visualization

PENNSTATE

1855

Motivation
Time series segmentation
Outlier Visualization for Barium
Other project ideas
Discussion

# Outline

1. **Motivation**

2. **Time series segmentation**
   - Barium Dataset
   - Sulphate Dataset

3. **Outlier Visualization for Barium**

4. **Other project ideas**

5. **Discussion**

Motivation
Time series segmentation
Outlier Visualization for Barium
Other project ideas
Discussion

# Outline

Motivation
Time series segmentation
Outlier Visualization for Barium
Other project ideas
Discussion

## Surface Water Data

- It includes *irregularly* spaced time series of Barium and Sulphate concentrations at approx. 80 PA counties from 1921-2015

- Deciphering spatial and temporal correlations may provide important insights about water quality deterioration due to energy extraction processes

- **Challenges**: spatial and temporal "sparsity"

- **Our current work**:
  - Naive spatial clustering using k-means
  - Time series segmentation

**4/23**

Motivation
Time series segmentation
Outlier Visualization for Barium
Other project ideas
Discussion

## Surface Water Data

- It includes *irregularly* spaced time series of Barium and Sulphate concentrations at approx. 80 PA counties from 1921-2015

- Deciphering spatial and temporal correlations may provide important insights about water quality deterioration due to energy extraction processes

- **Challenges**: spatial and temporal "sparsity"

- **Our current work**:
  - Naive spatial clustering using k-means
  - Time series segmentation

Statistical Analysis of Surface Water Data

Motivation
Time series segmentation
Outlier Visualization for Barium
Other project ideas
Discussion

Barium Dataset
Sulphate Dataset

# Outline

Motivation
Time series segmentation
Outlier Visualization for Barium
Other project ideas
Discussion

Barium Dataset
Sulphate Dataset

# Model and Goals
## Time series segmentation

- Assumption: There are *m* breakpoints.

- Model[1]:

$$Ba_i(Su_i) = Time_i\beta_1 + u_i \quad (i = 1, ..., T_1)$$
$$\vdots$$
$$Ba_i(Su_i) = Time_i\beta_m + u_i \quad (i = T_{m+1}, ..., T)$$

- Goals:

    - Point estimate of breakpoints.
    - Interval estimate of breakpoints.

[1]Bai, Jushan, and Pierre Perron. "Computation and analysis of multiple structural change models." Journal of applied econometrics 18.1 (2003): 1-22.

Motivation
Time series segmentation
Outlier Visualization for Barium
Other project ideas
Discussion

Barium Dataset
Sulphate Dataset

# Methodology
Time series segmentation

- Notation:
  - $h$: Minimum segment length.
  - $\text{SSR}^r_{t_1:t_2}$: Sum of squared residuals for the time segment $t_1$-$t_2$ with $r$ breakpoints.

- Recursive problem:

$$\text{SSR}_{1:T} = \min_{mh \leq j \leq T-h} \left[ \text{SSR}^{m-1}_{1:j} + \text{SSR}^0_{(j+1):T} \right]$$

- Construction of triangular matrix of sums of squared residuals.
- Time Complexity: $O(T^2)$

Statistical Analysis of Surface Water Data

Motivation
Time series segmentation
Outlier Visualization for Barium
Other project ideas
Discussion

Barium Dataset
Sulphate Dataset

## Confidence intervals for breakpoints

- Asymptotic distribution of breakpoint estimate:

$$A_T(\hat{T}_i - T_i) \xrightarrow{d} \text{some Wiener process}$$

where $A_T$ is a normalization constant.

- Using above asymptotic distribution function of the breakpoint.[2], 95% confidence intervals for point estimates $\hat{T}_1,..., \hat{T}_m$ can be created.

---

[2]Bai, Jushan. "Estimation of a change point in multiple regression models." Review of Economics and Statistics 79.4 (1997): 551-563.

Motivation
Time series segmentation
Outlier Visualization for Barium
Other project ideas
Discussion

Barium Dataset
Sulphate Dataset

# Choosing the number of segments

- Bayesian Information Criterion (BIC):

$$BIC = -2\log(\hat{L}) + k\log(n)$$

  - Limitation in our case: $n >> k$ condition not satisfied by Barium data for some clusters

- Residual sum of squares:

  - Always decreases and so can't be used for choosing segmentation.

Motivation
Time series segmentation
Outlier Visualization for Barium
Other project ideas
Discussion

Barium Dataset
Sulphate Dataset

# Outline

Motivation
Time series segmentation
Outlier Visualization for Barium
Other project ideas
Discussion

Barium Dataset
Sulphate Dataset

# Full time series



Figure: Full time series for Barium dataset

Motivation
Time series segmentation
Outlier Visualization for Barium
Other project ideas
Discussion

Barium Dataset
Sulphate Dataset

## Dense time series



Figure: Dense time series for Barium dataset

Motivation
Time series segmentation
Outlier Visualization for Barium
Other project ideas
Discussion

Barium Dataset
Sulphate Dataset

# BIC and RSS



Figure: BIC and RSS plot for dense part of Barium dataset

Motivation
Time series segmentation
Outlier Visualization for Barium
Other project ideas
Discussion

Barium Dataset
Sulphate Dataset

## Time segmentation for dense part

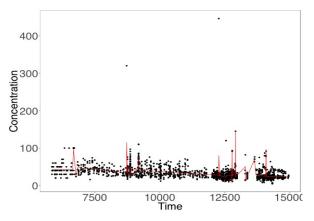

Figure: Time segmentation for dense part of Barium dataset with 18 breakpoints

Motivation
Time series segmentation
Outlier Visualization for Barium
Other project ideas
Discussion

Barium Dataset
Sulphate Dataset

# Outline

Motivation
Time series segmentation
Outlier Visualization for Barium
Other project ideas
Discussion

Barium Dataset
Sulphate Dataset

## Full time series



Figure: Full time series for Sulphate dataset

Motivation
Time series segmentation
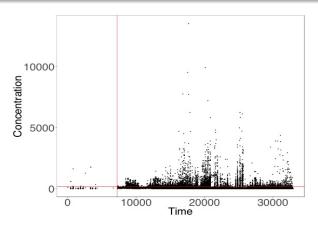Outlier Visualization for Barium
Other project ideas
Discussion

Barium Dataset
Sulphate Dataset

## Dense time series



Figure: Dense time series for Sulphate dataset

Motivation
Time series segmentation
Outlier Visualization for Barium
Other project ideas
Discussion

# Outline

Motivation
Time series segmentation
Outlier Visualization for Barium
Other project ideas
Discussion

# Outlier map for Barium dataset



Figure: Outlier map for Barium dataset

Motivation
Time series segmentation
Outlier Visualization for Barium
Other project ideas
Discussion

# Outline

Motivation
Time series segmentation
Outlier Visualization for Barium
Other project ideas
Discussion

## Other project ideas

- Bedrock zone spatial clustering and time series segmentation.
- Analyzing river system networks:
  - Yearly correlation networks.
  - Identifying the community of outliers.
- High dimensional tests
  - Simultaneous testing for methane concentrations in 1000 wells before and after drilling.
  - False discovery rate to control the type I error.

Motivation
Time series segmentation
Outlier Visualization for Barium
Other project ideas
Discussion

# Outline

1. **Motivation**

2. Time series segmentation
   - Barium Dataset
   - Sulphate Dataset

3. Outlier Visualization for Barium

4. Other project ideas

5. **Discussion**

Motivation
Time series segmentation
Outlier Visualization for Barium
Other project ideas
Discussion

## Discussion

Summary of current works:

- Time series segmentation for dense regions
- Interactive Shiny app

Future works:

- Robust estimation of breakpoints.
  - $SAR^r_{t_1:t_2}$: Sum of absolute residuals for the time segment $t_1$-$t_2$ with $r$ breakpoints.

$$SAR_{1:T} = \min_{mh \leq j \leq T-h} \left[ SAR^{m-1}_{1:j} + SAR^0_{(j+1):T} \right]$$

- Agglomerative clustering and Spatial segmentation based on trend

Motivation
Time series segmentation
Outlier Visualization for Barium
Other project ideas
Discussion

## Discussion

Summary of current works:

- Time series segmentation for dense regions
- Interactive Shiny app

Future works:

- Robust estimation of breakpoints.
  - $SAR^r_{t_1:t_2}$: Sum of absolute residuals for the time segment $t_1$-$t_2$ with $r$ breakpoints.

$$SAR_{1:T} = \min_{mh \leq j \leq T-h} \left[ SAR^{m-1}_{1:j} + SAR^0_{(j+1):T} \right]$$

- Agglomerative clustering and Spatial segmentation based on trend