



Water analytics- the evidence chain from measurement to model to decision

E Marian Scott

School of Mathematics and Statistics

WWC, May 2015



C. A. Miller, R. A. O'Donnell, K. Gallacher, A. Elayouty

School of Mathematics and Statistics, University of Glasgow, UK

Maria Franco Villoria

Department of Economics and Statistics, University of Turin, Italy

Francesco Finazzi

Department of Statistics, University of Bergamo, Italy

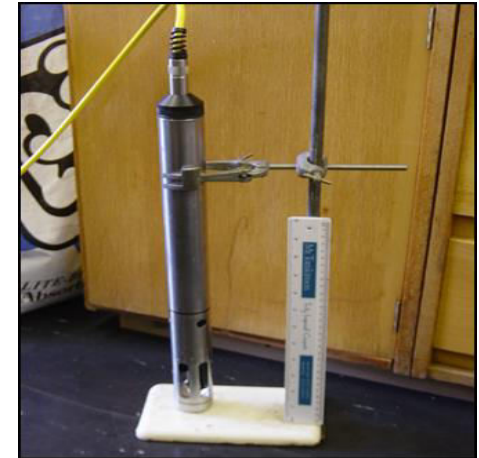
Acknowledgements

SEPA for data and support (MFV), Susan Waldron, University of Glasgow funding (AE and KG), EA (KG), Scottish Water (support for MFV, RO)

Four questions:

- What is changing?
- What are the changes?
- What is driving the changes?
- How certain are we?

Sensor technology delivering
**enhanced dynamic detail of
environmental systems at
unprecedented scale**

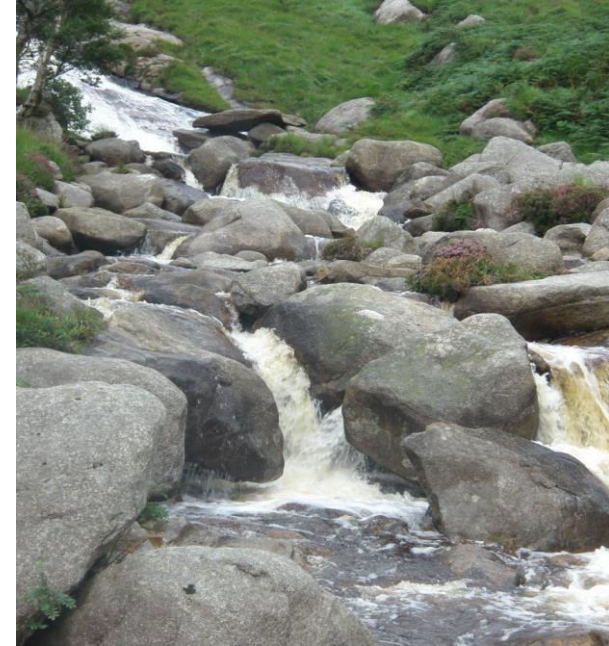


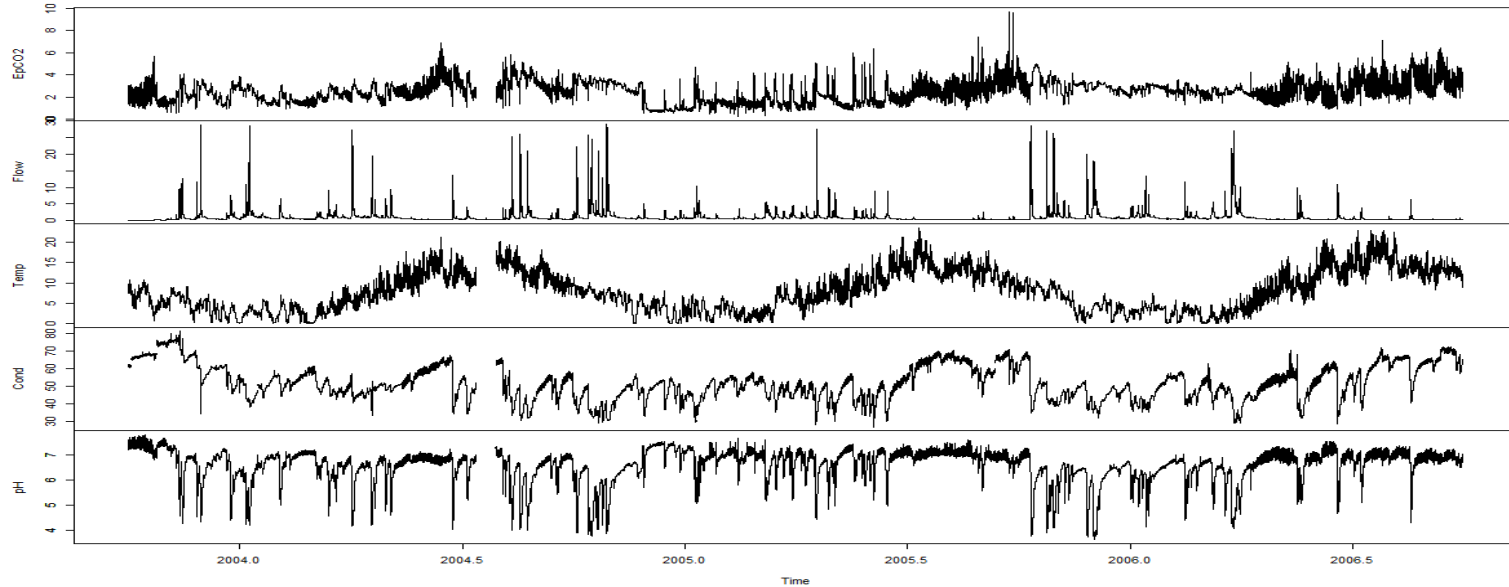
What analytic challenges?

The data needs are driven by policy, regulation and management.

Analytic challenges:

- multi-pollutant data from monitoring networks
- Many covariates: meteorological, land morphology, from different data streams
- Dealing with uncertain climate change (mitigation and adaptation)
- **BIG data-Real time analysis for control and management**



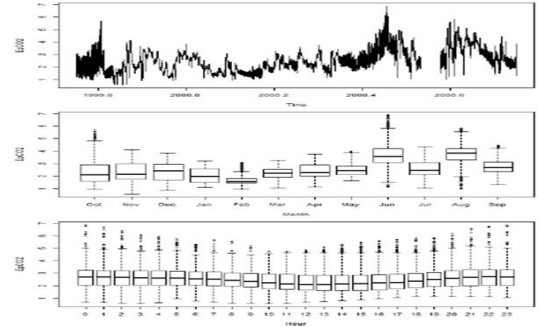
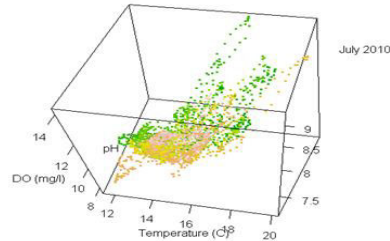
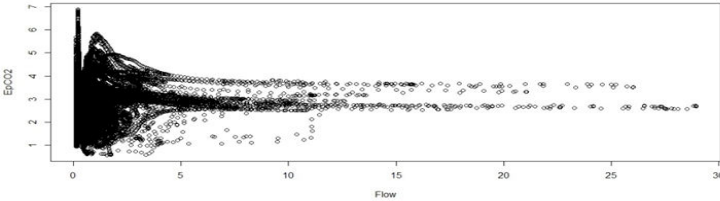
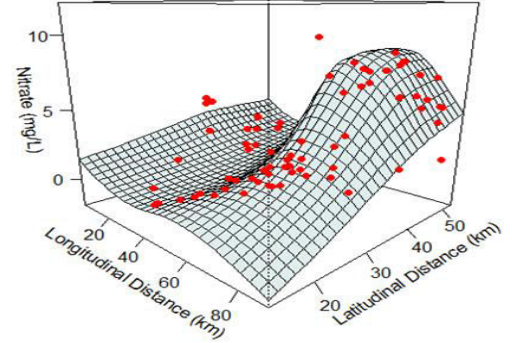
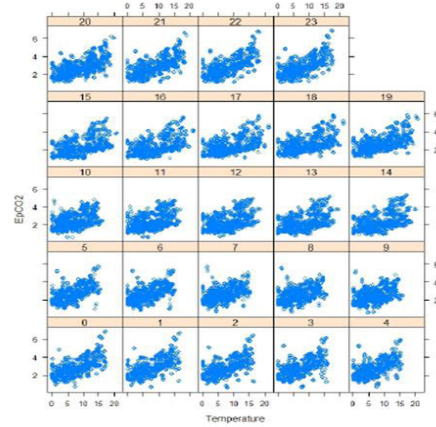


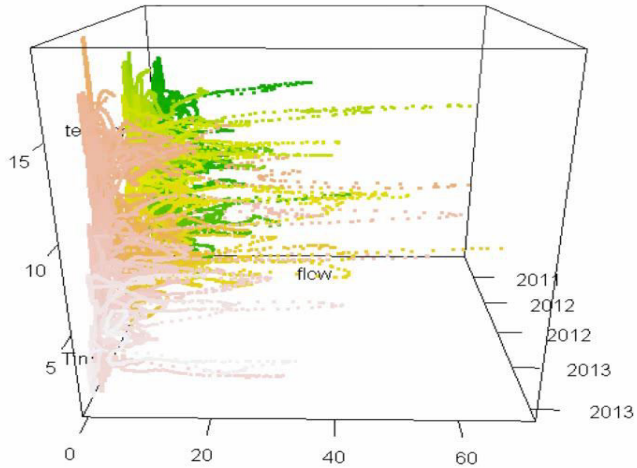
Partial pressure CO_2 , flow, temperature, conductivity and pH over 3 years- **Why? understanding carbon dynamics in small streams**

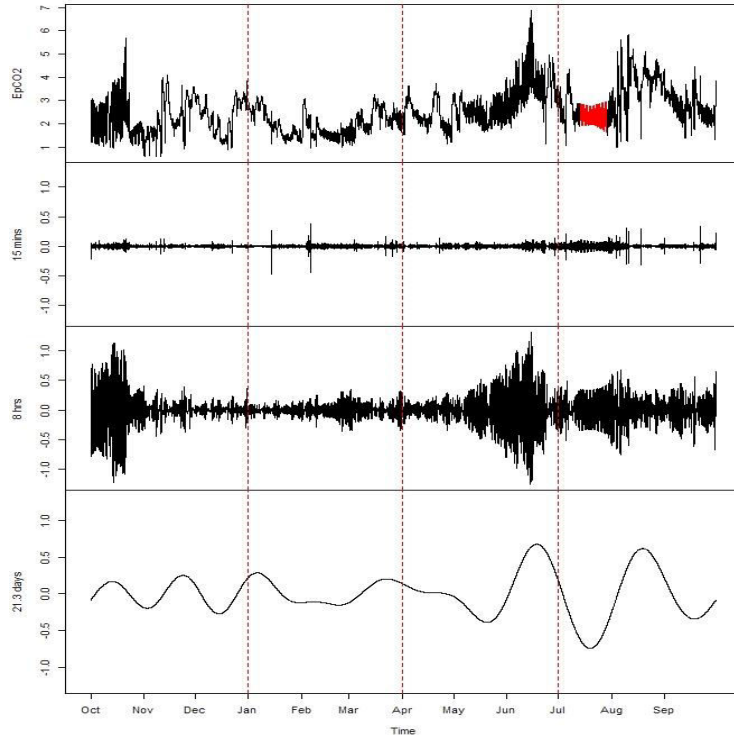


What are the questions?

- How is epCO₂ (or the measurand) changing?
- What are the drivers of change?
- Events, anomalies, unusual conditions







- 15 minute data, 1 hydrological year-a wavelet analysis



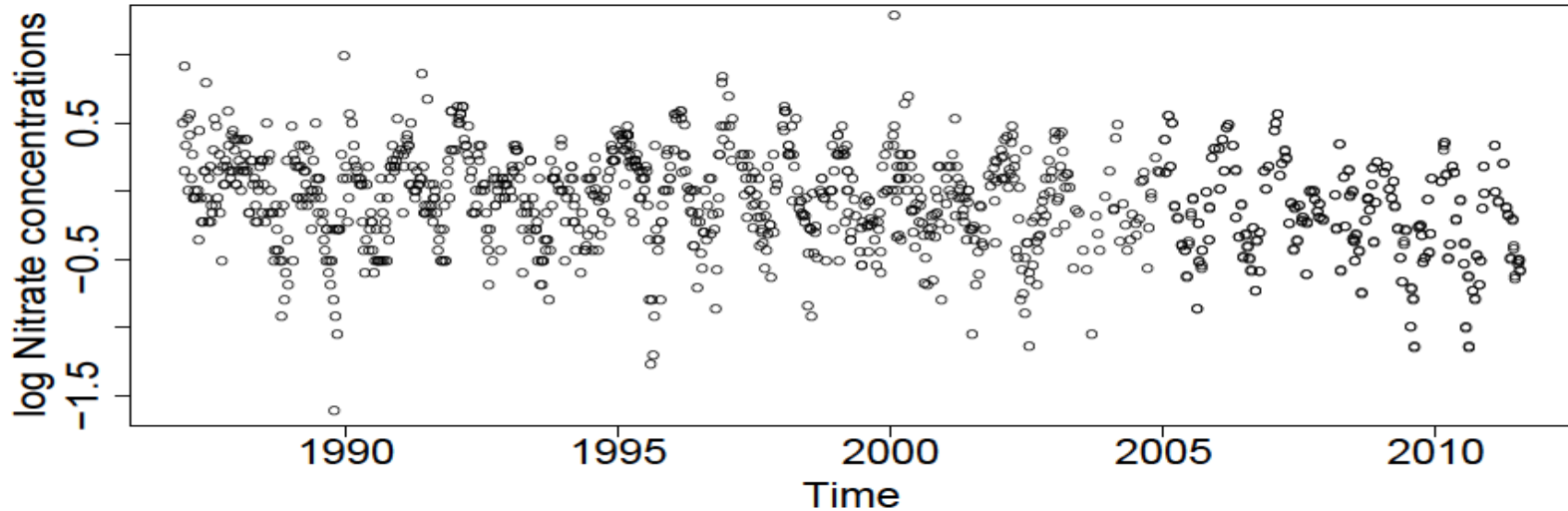
- The highest variability is in the 8 hour signal. This reflects changes in photosynthetic / respiratory dominance, changing seasonally.



- more variation during summer than winter reflecting differences in CO₂ input versus consumption



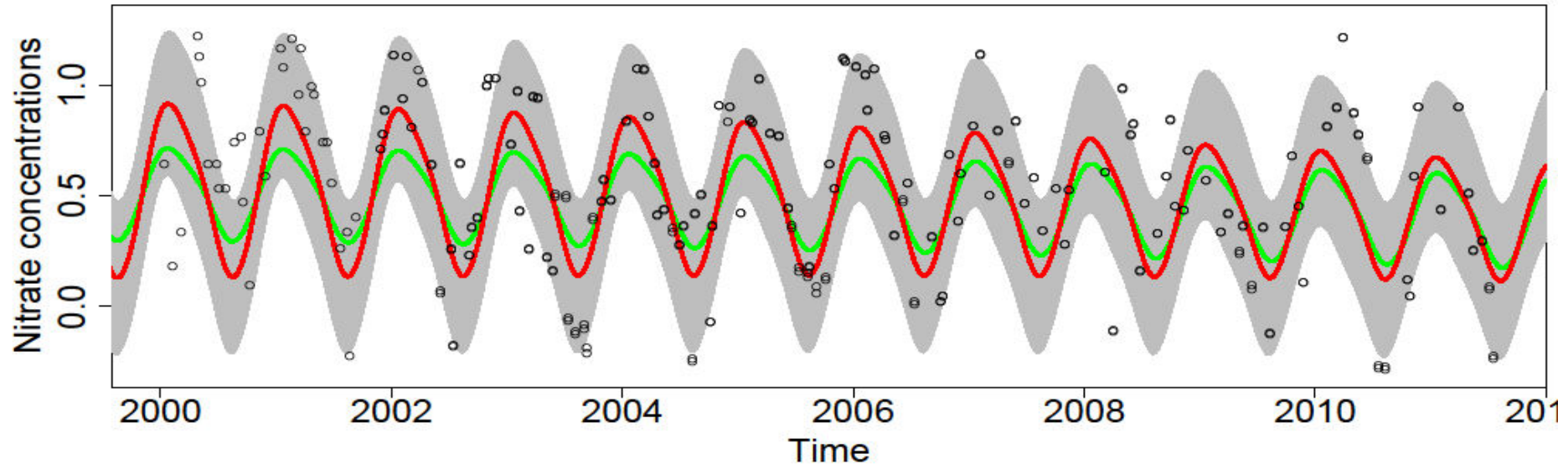
- How should we describe the trend? Linear or smooth? – how smooth?
- **are there seasonal components?** Constant in time? Changing phase or amplitude over time?





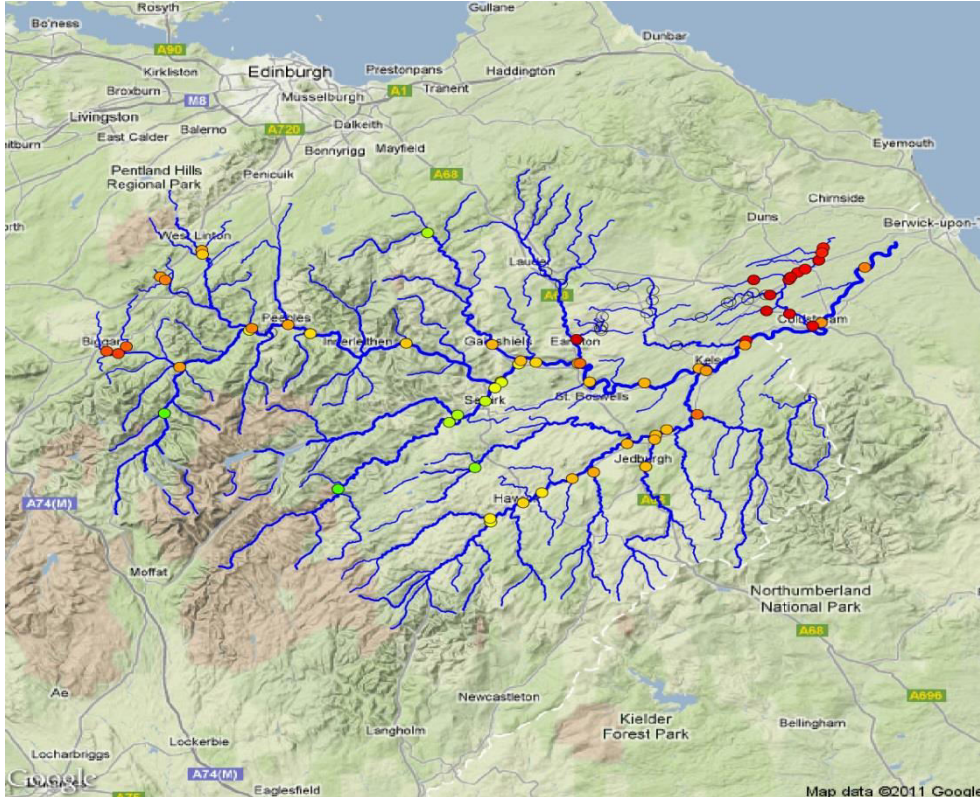
Building models

- How should we describe the trend and the seasonal components? Using smooth functions (based in splines) we can capture complex patterns, test different models and capture our uncertainty.

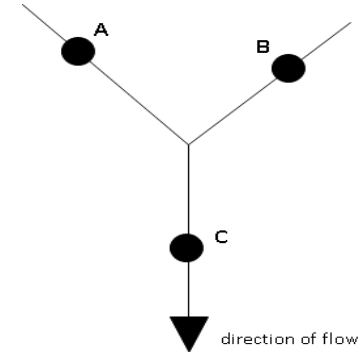




River networks



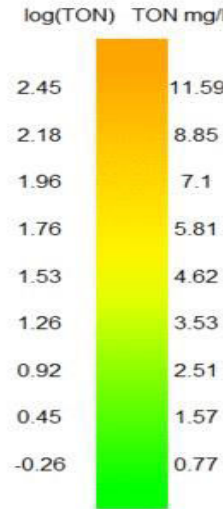
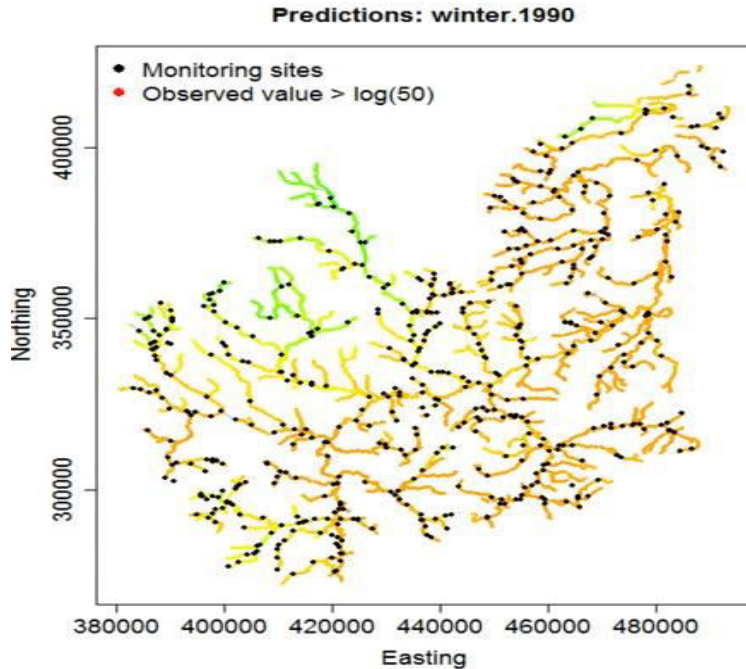
Spatial models for stream networks using stream distance rather than Euclidean distance. The user can specify if monitoring sites are 'flow connected' (A and C or B and C) or 'flow unconnected' (A and B).



Flexible regression models over river networks, O'Donnell, Rushworth, Bowman, Scott and Hallard (2014)



Looking within one Large Hydrological Area



The following models were fitted to the dominant network in Trent:

$$\log(\text{TON}) = \text{Eastings} + \text{Northings} + \varepsilon$$

$$\log(\text{TON}) = \text{Eastings} + \text{Northings} + z_e + \varepsilon$$

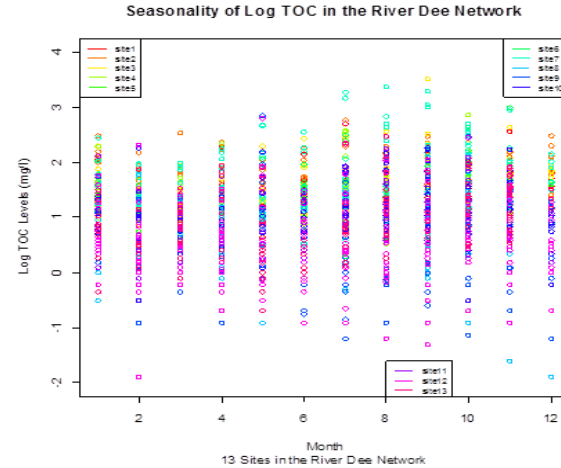
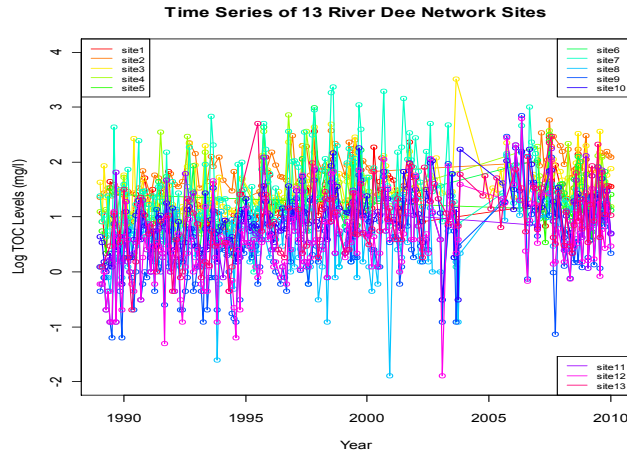
$$\log(\text{TON}) = \text{Eastings} + \text{Northings} + z_u + \varepsilon$$

$$\log(\text{TON}) = \text{Eastings} + \text{Northings} + z_u + z_e + \varepsilon$$



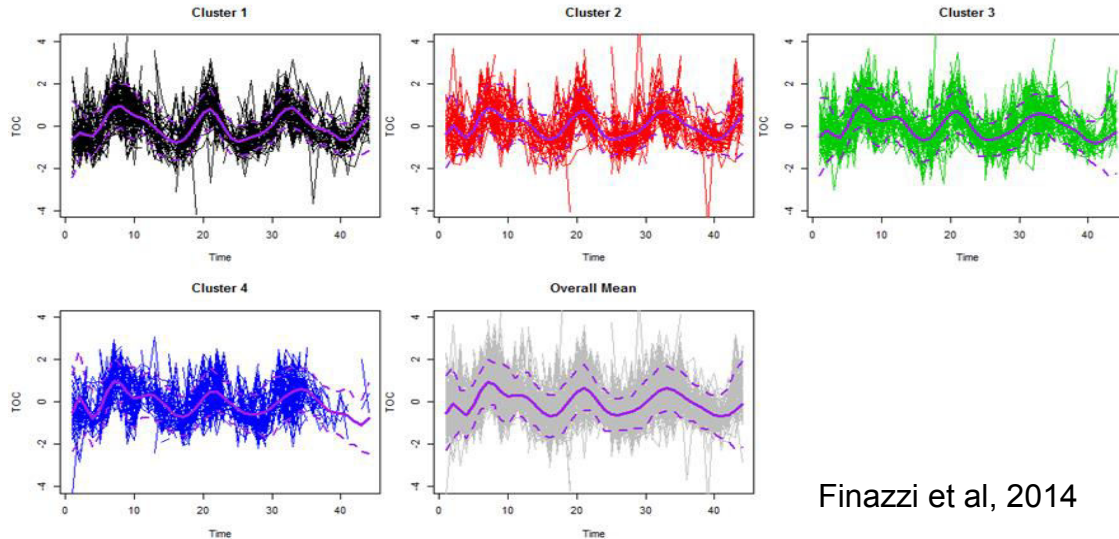
BIG data modelling

With a network of monitoring sites- *a data deluge*, the time series at each site can be regarded as a curve, the curve then becomes the “*data point*”. The statistical model is based on the curves or functions which are assumed to be smooth. This is known as functional data analysis (FDA).

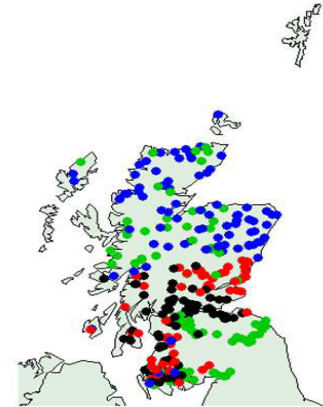




Functional clustering methodology has been applied to Total Organic Carbon (TOC) data from 333 monitoring locations across rivers in Scotland over 44 months, covering the period January 2007 - August 2010. (in collaboration with SEPA).



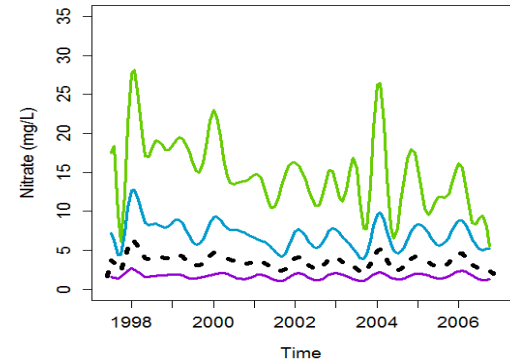
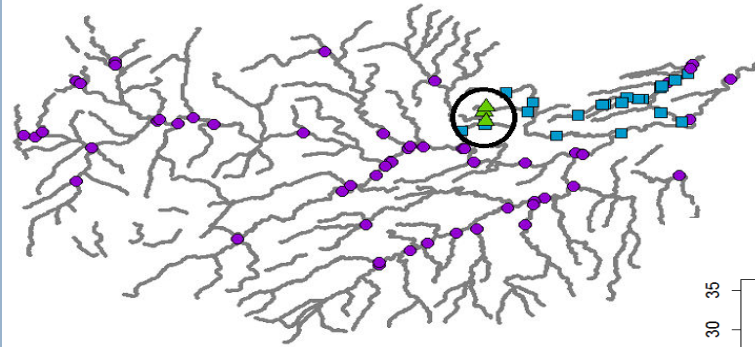
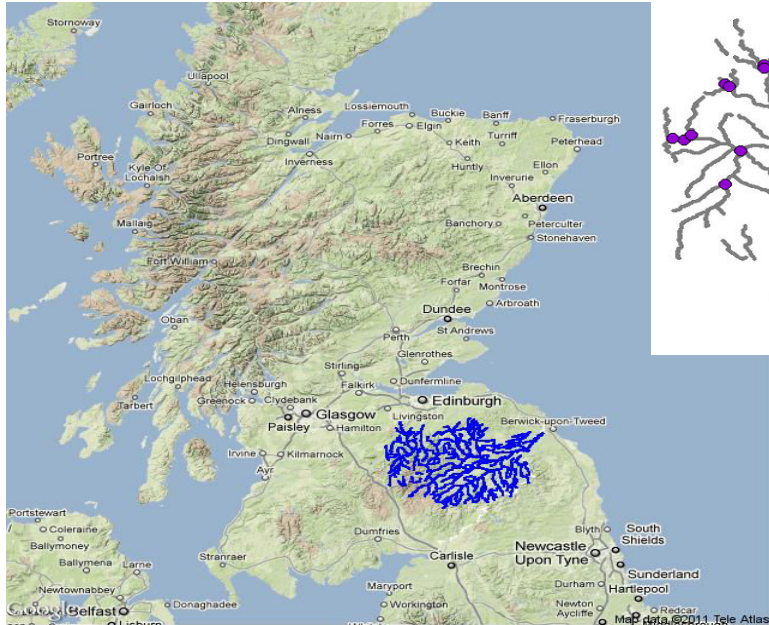
Hierarchical (Euclidean correlation) Rivers





Nitrate in Tweed river basin

Functional clustering methodology has been applied to nitrate data in Tweed river basin- 77 monitoring sites, 10 years of data. **Why? Can we simplify the network?**





- Data characteristics- quantity and quality, missingness, limits of detection.
- Non stationary, complex nature of the relationships
- For networks of sensors- building fast and efficient spatio-temporal models which scale, functional data analysis provides part of that solution
- uncertainty evaluation and visualisation

- 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)*
- Haggarty R, Miller C A, Scott E M, Wylie F, Smith M (2012) Functional clustering of water quality data in Scotland. *Environmetrics*
- Finazzi, F., Haggarty, R., Miller, C., Scott, M., Fasso, A. A comparison of clustering approaches for the study of the temporal coherence of multiple time series, *Stochastic Environment Research and Risk Assessment* (2013) .
- Miller C, Magdalina A, Willows R, Bowman A, Scott E M, Lee D, Burgess C, Pope L, Pannullo F, Haggarty R (2014). Spatiotemporal statistical modelling of long term change in river nutrient concentrations in England and Wales. *Sci Tot Env*, 466-467.
- El-Ayouty A, Scott E M, Miller C, Waldron S, Franco-Villoria M . Challenges in Modelling Detailed and Complex Environmental Data Sets: A Case Study Modelling the Excess Partial Pressure of Fluvial CO₂, submitted J Ecological and Environmental Statistics