# Detecting anomalies in sensor network data

**Richard Jarrett**

CSIRO Mathematical & Information Sciences, Melbourne, Australia

Cherry Bud Workshop, 26 March 2008

# Summary

- Background on "sensor networks"
- Project of "anomaly detection" in water distribution systems
- Methods used for identifying anomalous events for a single sensor
- Issues about multiple sensors
- Estimation of "travel time" between sensors

# Sensors and sensor networks

CSIRO is a Government funded research organisation in Australia, with 6500 employees, focussing on major national issues. 40% of our budget comes from work with business and industry.

'Sensors and sensor networks' is a major focus area for CSIRO, whose aim is:

*To create technologies to radically reduce the cost and improve the quality of data gathering to*

- *enhance the understanding of our natural environments and*
- *provide the ability to manage & exploit Australia's resources.*

# Sensors and sensor networks

## CSIRO work is focussed on

- Development of new sensors
- Data transmission protocols
- Distributed processing/autonomous decisions

## For our Division of CSIRO, the interest lies in

- How reliable is the data we are collecting?
- What do we do with the data that is collected?
- How many sensors, how frequently we measure?
- Optimal placement of the sensors

# Sensors and sensor networks

- ## WRON (Water Resources Observation Network):
  - Water accounting – using flow sensors and other information to find out how much water there is and where it's going
  - Water forecasting – predictive models based on matching sensor outputs to runoff and flow models, checking calibration

- ## CMIS/CLW projects (AwwaRF, Sydney Water, Water Corp):
  - Gauges measuring depth and flow in sewer systems
  - Calibration issues for both gauges and models
  - Now looking at measurements in water distribution systems
  - Aim to detect anomalous events and take action
  - Also used for detecting calibration problems
  - Now targetting "travel time" between sensor locations
  - Will ultimately be able to follow "events" through the system

# "Anomalous events" in sensor networks

- **Study funded by CSIRO and the American Water Works Research Foundation (AwwaRF)**
  - Literature review of current methodologies for analysis and evaluation of on-line water quality data
  - Application of the most promising methods to data sets obtained from a number of Australian and US water utilities

- **The methods considered will eventually lead to**
  - Better understanding of water distribution systems
  - Techniques which enable identification of anomalous events
  - In particular, events which might be linked to security issues

# Data available

- Australia
  - City West Water (Melbourne, Australia)
  - Hunter Water (Newcastle, Australia)
  - South East Water (Melbourne, Australia)

    (All used the same instrumentation with pH, ORP (oxidation-reduction potential), TEMP every 10 min.)
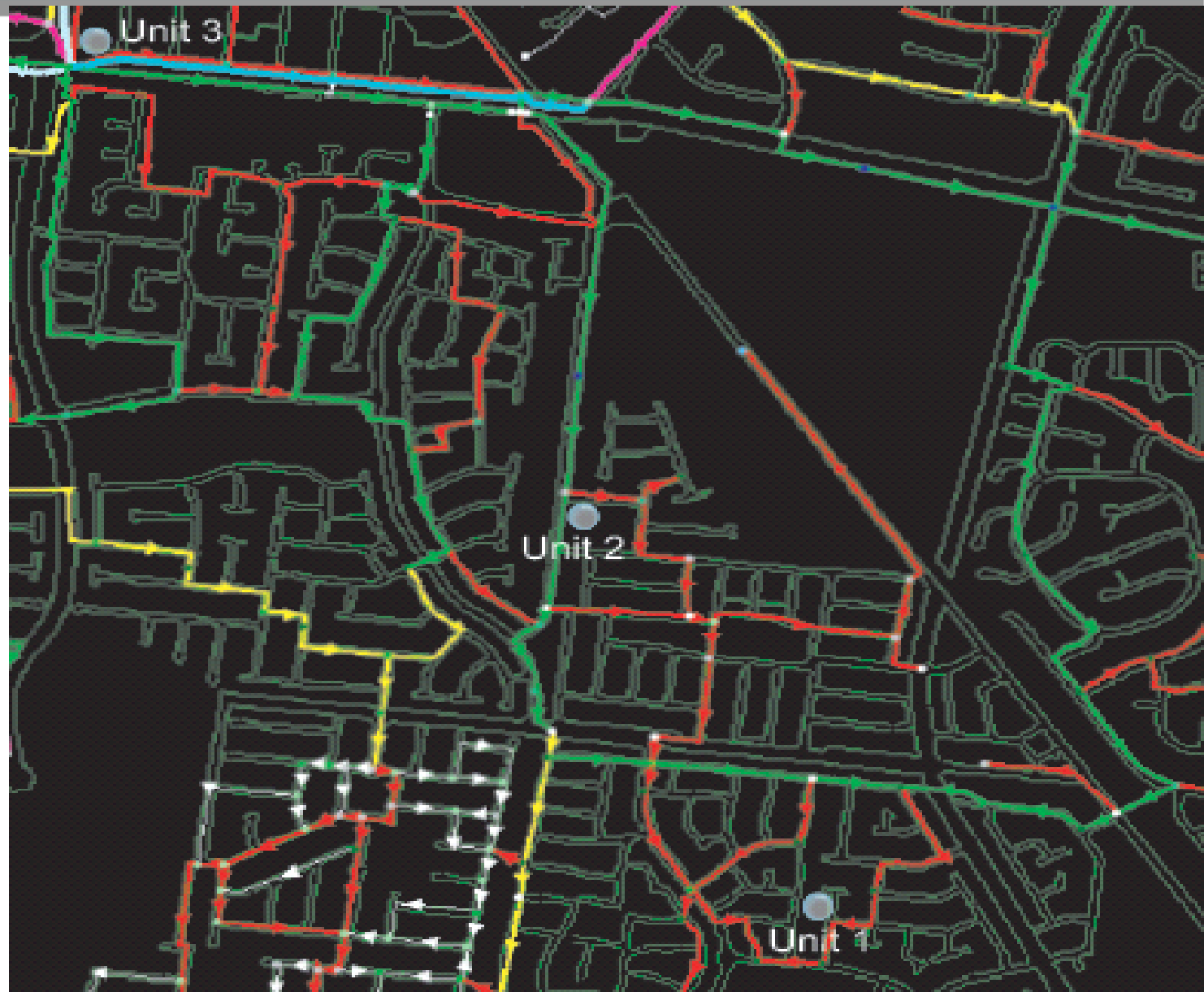
- United States
  - Philadelphia (every 1 min)
  - Oklahoma City (every 15 min)

    (Mainly pH, Electrical Conductivity (EC), Turbidity, Residual Chlorine;  Flow and Pressure are sometimes available but these refer to the manifold where the sensor is, not the pipe)

- Typically, about a year of data in each case, for a number of sites

# "Anomalous events" in sensor networks

- This shows a typical network
- Arrows show 'usual' direction of flow
- Water comes in at Unit 3 and flows down and to the right



Cherry Bud Workshop

# Metadata

- "Metadata" is vital for an understanding of the system and identification of possible reasons for anomalies

- System data
    - Details of variables/equipment/units
    - Codes/values used when data is missing or equipment is off-line
    - Method and timing of data retrieval from equipment to computer
    - Time standards, eg daylight saving

- Event data
    - Time and duration of maintenance/calibration events
    - Time and duration of major system problems, eg pump failures, mains breakages, treatment failures
    - School holidays, public holidays
    - Major weather events
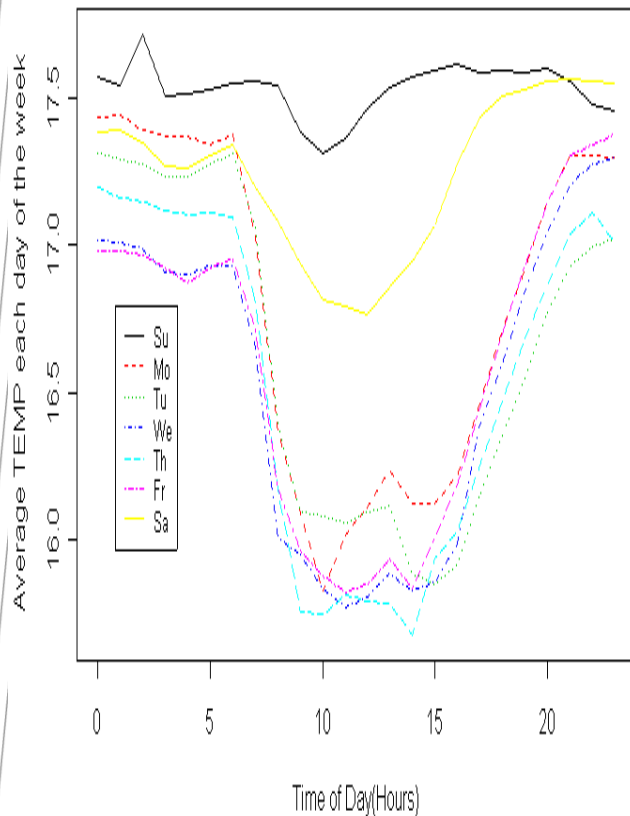
# "Anomalous events" in sensor networks

| Techniques available | Comments |
|---|---|
| Statistical models<br><br>Data mining | Flow rates vary on a daily and weekly basis. This creates daily and weekly patterns in measured variables due to "time spent in pipes" (eg cement-lined pipes change the pH of water)<br><br>These methods do not cope with "slowly varying changes" |
| Time series models<br><br>Control charting techniques | Data too erratic with daily, weekly, seasonal changes Hence not suitable for original data but work well on "differenced" data |
| Adaptive models, eg Kalman filter | Adapt well to slow changes but still allow detection of rapid changes |
| Multivariate versions for multiple variables per sensor and/or multiple sites | Lack of correlation so multivariate results similar to univariate<br><br>Multiple sites requires knowledge of "travel time" |

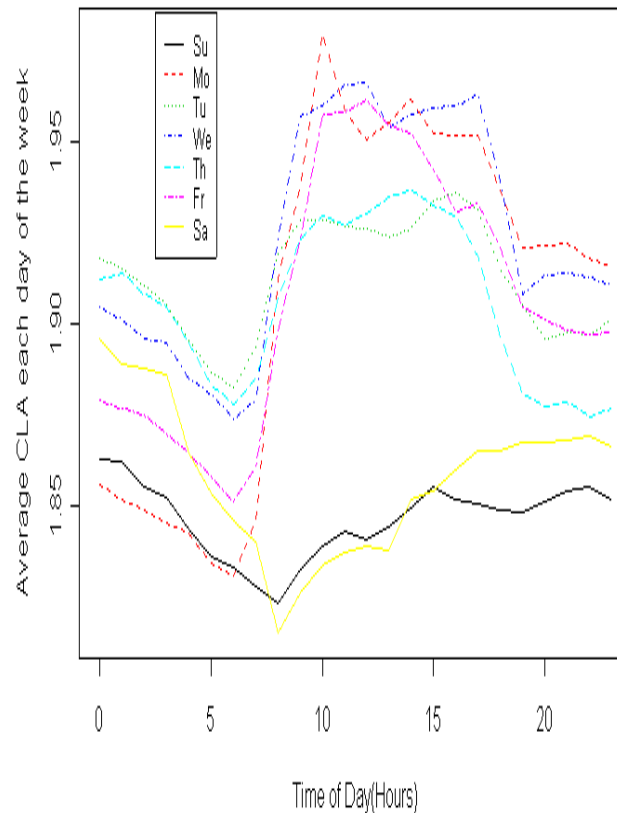# "Anomalous events" in sensor networks

- **Particular problems with on-line sensors**
  - Data is generally not as reliable as laboratory-based analyses
  - Tendency for instrumental drift, so there is a need for maintenance and re-calibration at (typically monthly) intervals
  - Local disturbances can occur
  - Volume of data, often once a minute from >20 sites, is large
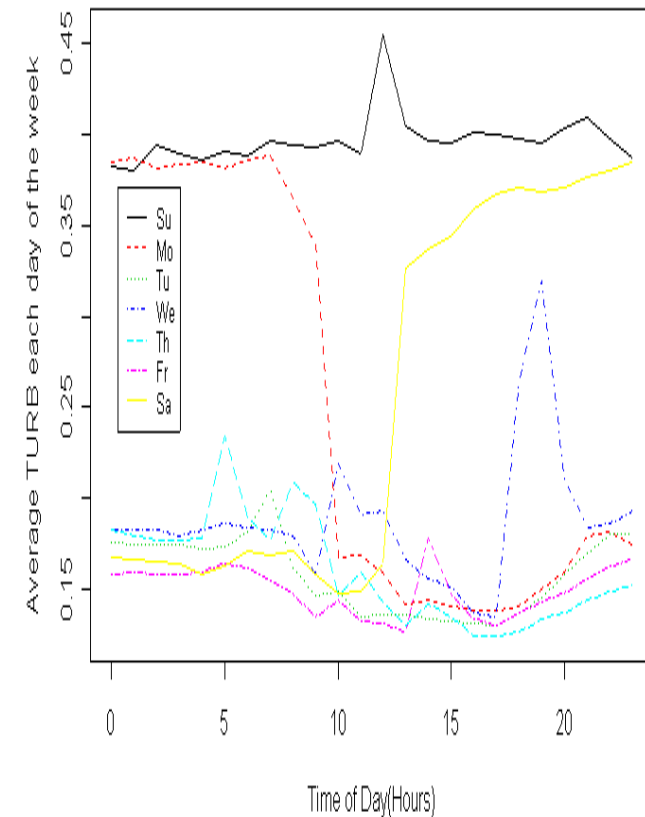
# Example 1: Day-of-week patterns

- Most water utilities show patterns based on time of week
- Here, CLA (Chlorine) remains low Sat-Sun but has a clear cycle Mon-Fri
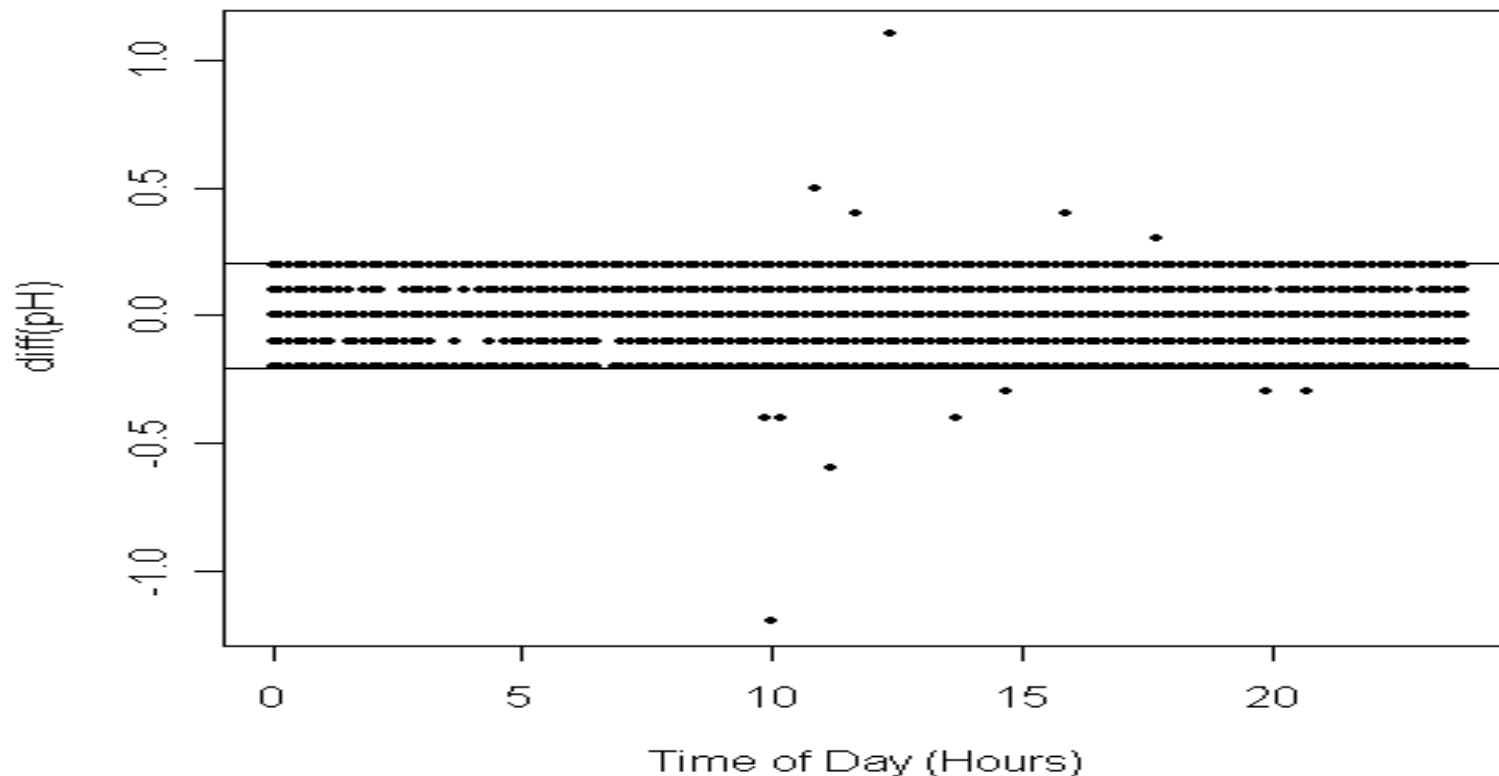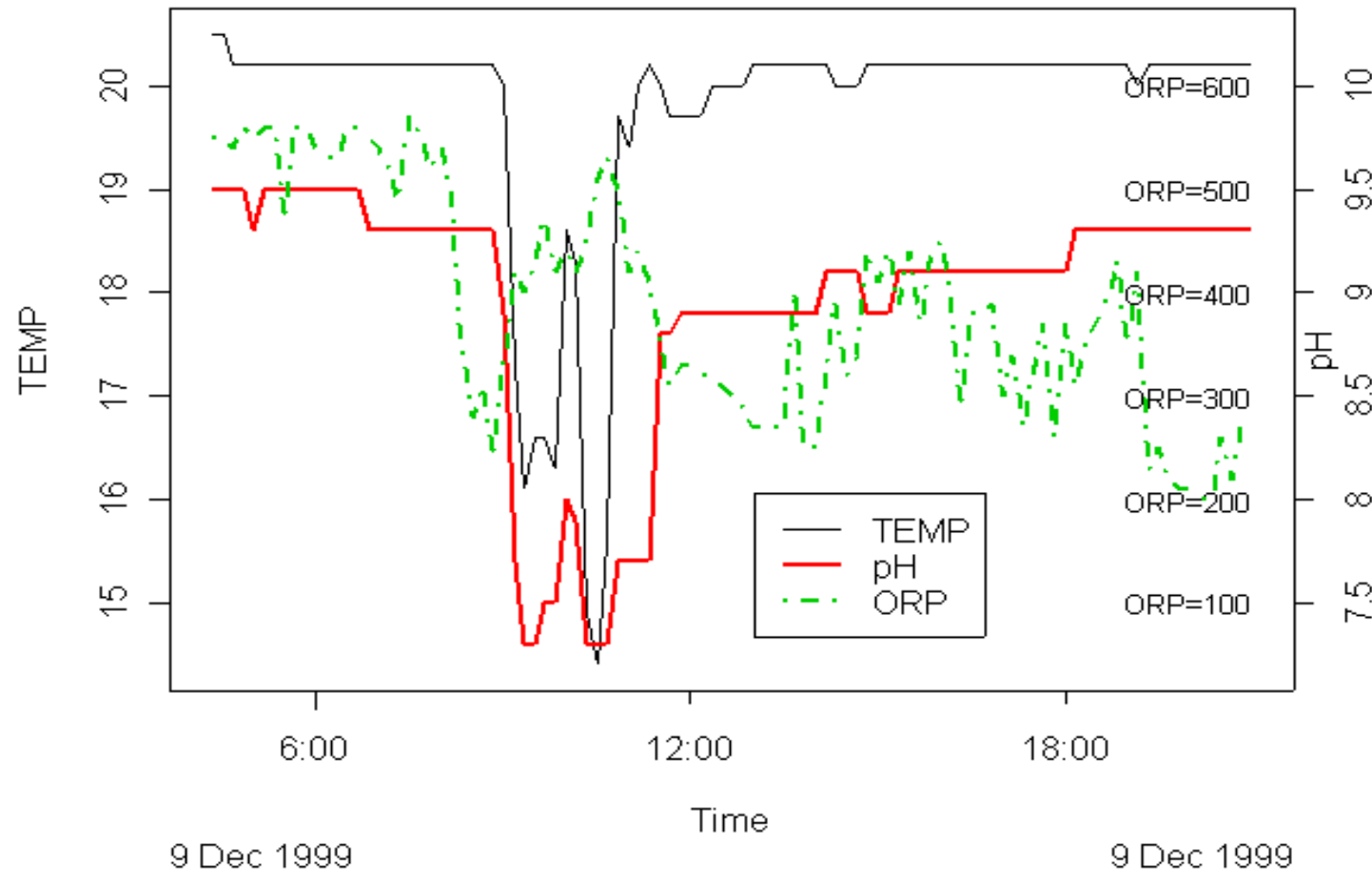
# Example 2: Control charts using first differences

- First differences remove the slow changes/cycles
  - Generally stationary, but can still have minor time-of-day effects, so we plot by Time of Day (for the whole year)
  - Usual Control Limits too narrow – generally we take about $5\sigma$
  - Here, 13 points (out of 55,000) are outside the limits:

# Example 3: A major event

- 11/13 of these outliers occur on one day
- Either a maintenance event or a change in flow direction at the sensor
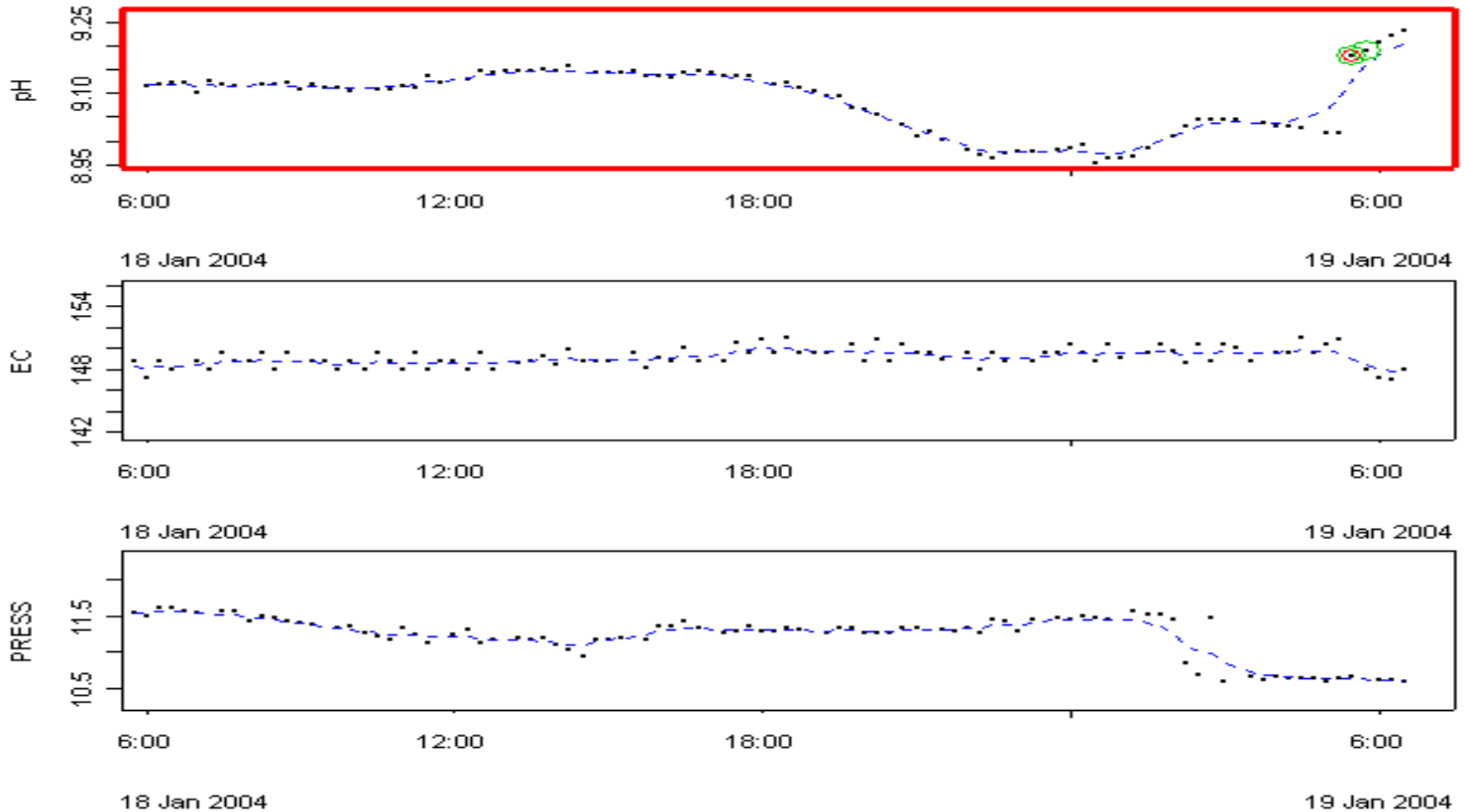
# How should we move forward?

- Control charts of differences
  - Works well because first differences are close to independent (e.g. look at variograms)
  - If data points every 10-15 minutes, ideal for picking up sharp changes over this sort of time period

- Models based on "white noise + Brownian motion" do quite well here

- For these models, EWMA charts would work well; effectively, Kalman filters

- Kalman filters with different levels of filtering enable us to "tune" our detection methods to pick up "events" of different shapes and sizes

# Adaptive schemes: Kalman filters

- We identify 4 alarm codes:
  - Difference between data point and 1-min Kalman filter
  - Departure of slope of 1-min Kalman filter from average over last 31 days
  - Departure of slope of 10-min Kalman filter from average over last 31 days
  - Departure of slope of 60-min Kalman filter from average over last 31 days

- In each case, an alarm is triggered if the quantity is greater than (default) $5 \times$ average abs diff over last 31 days.

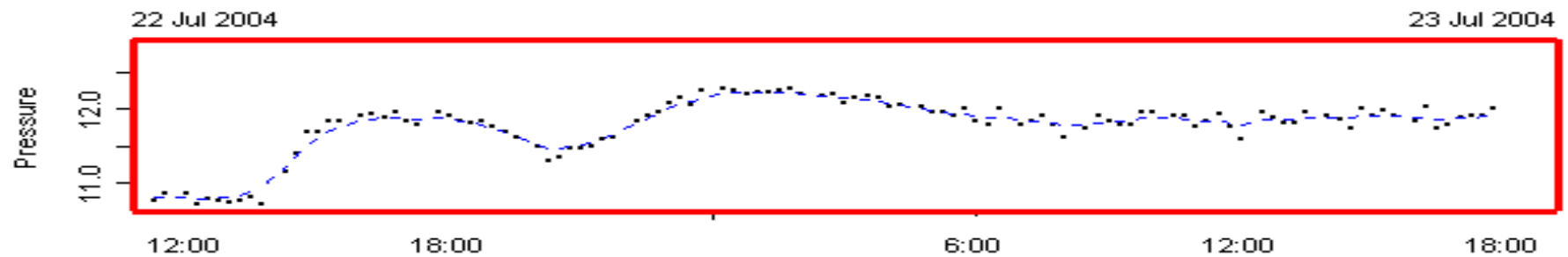- These alarms are colour coded with (×), (+), a small (o), and a larger (O).
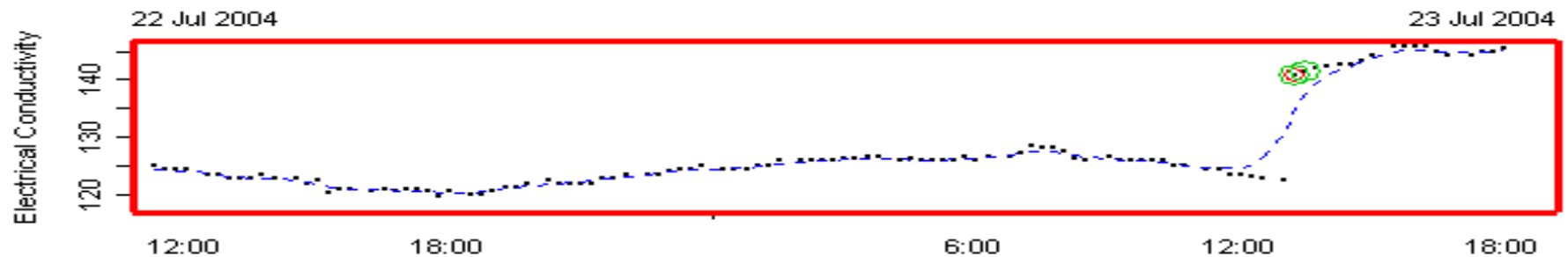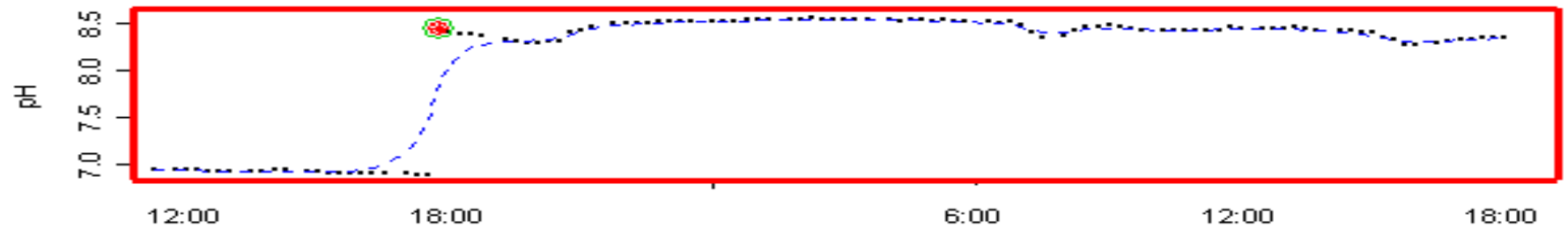
# Example 4: US Utility 2, 15 min data

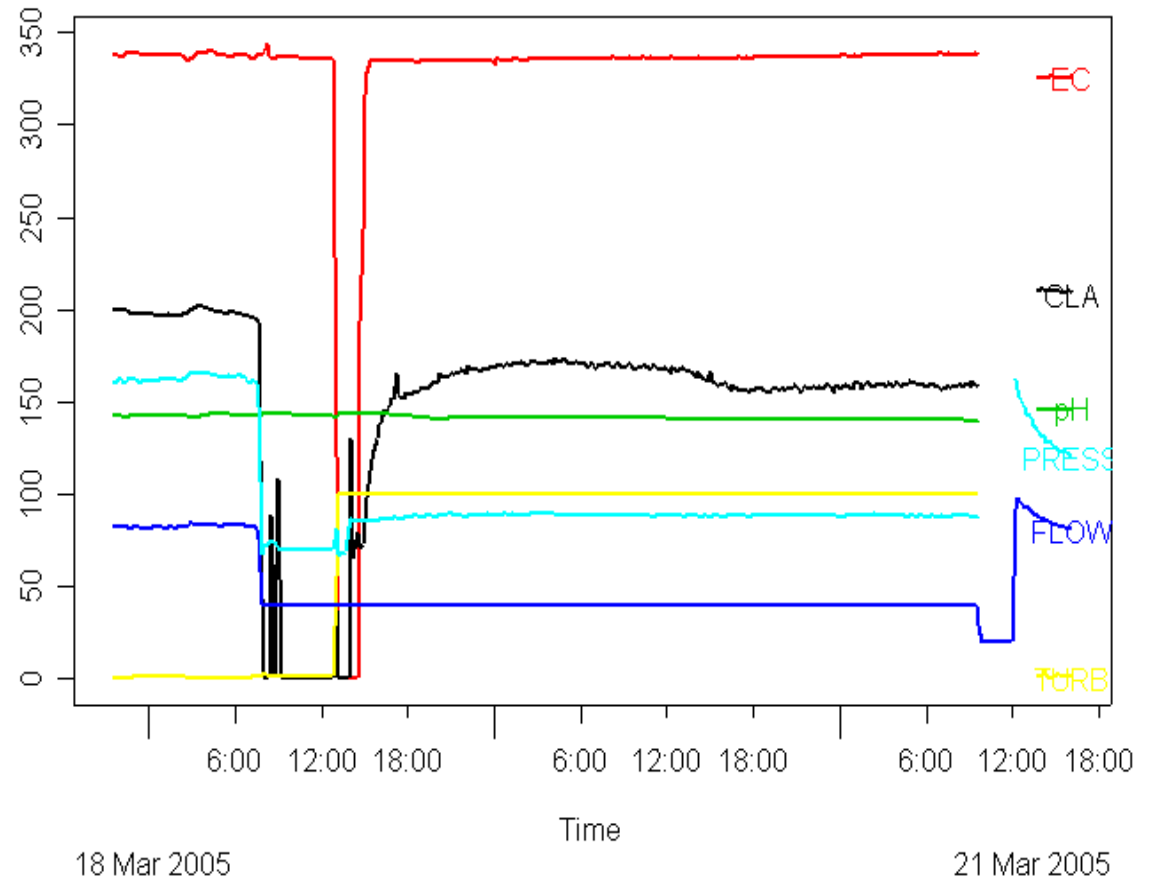- Here is the first alarm we see (••• = data,--- = 60-min filter)

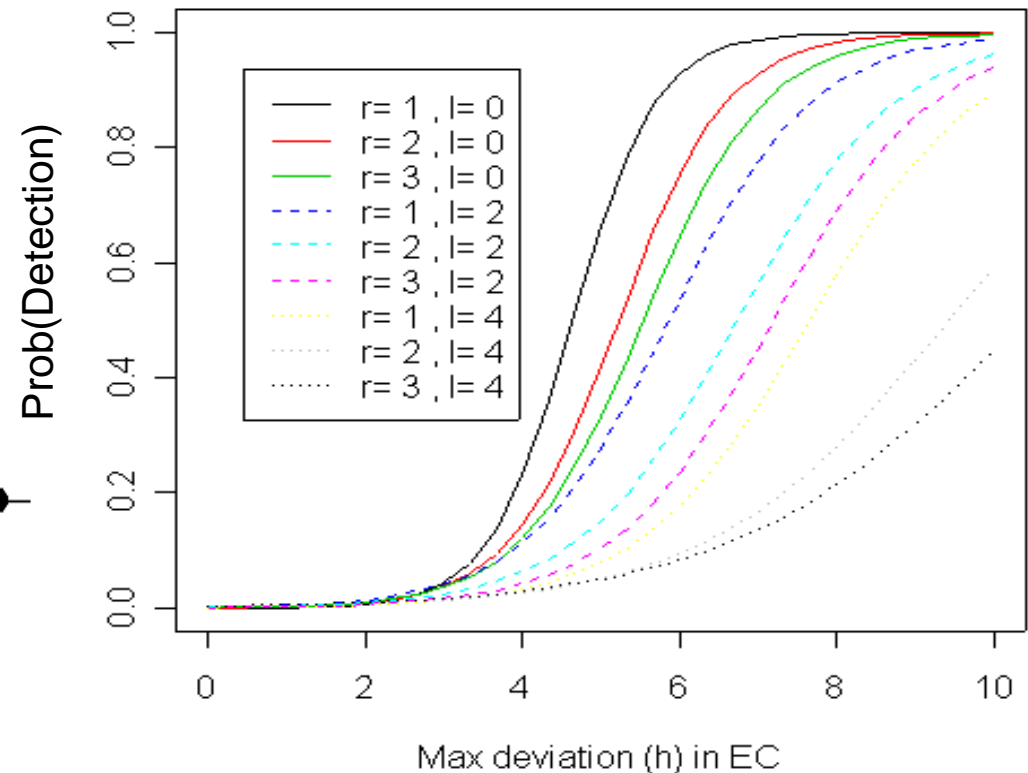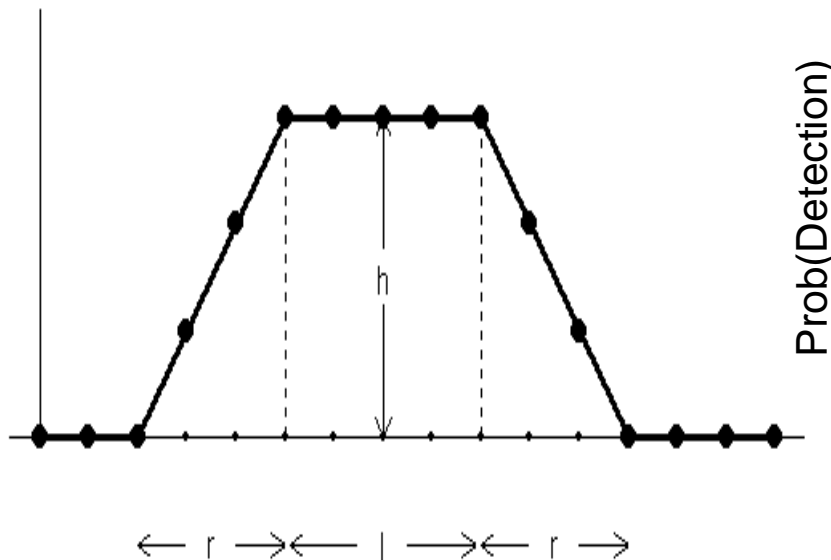# Example 4: US Utility 2, 15 min data

- Another example:

# Example 5: US Utility 1, 1 min intervals

- Looking at multiple variables is useful
- Here, two major events:
  - First is a 4h disruption, possibly due to fouling (blockage) of the sensor
  - Second is 1h, likely to be maintenance event
- Note gradual return to stable levels, often different from previous stable levels

# How well does the Kalman filter work?

- Can look at false positive/negative rates by adding artificial 'events' and looking at the detection rates
- Events chosen have different size (h), slope (r) and duration ($\ell$)
- Example shows ROC (power curves) for 10-min EC data from US



Prob(Detection) vs Max deviation (h) in EC

Legend:
- r= 1 , l= 0
- r= 2 , l= 0
- r= 3 , l= 0
- r= 1 , l= 2
- r= 2 , l= 2
- r= 3 , l= 2
- r= 1 , l= 4
- r= 2 , l= 4
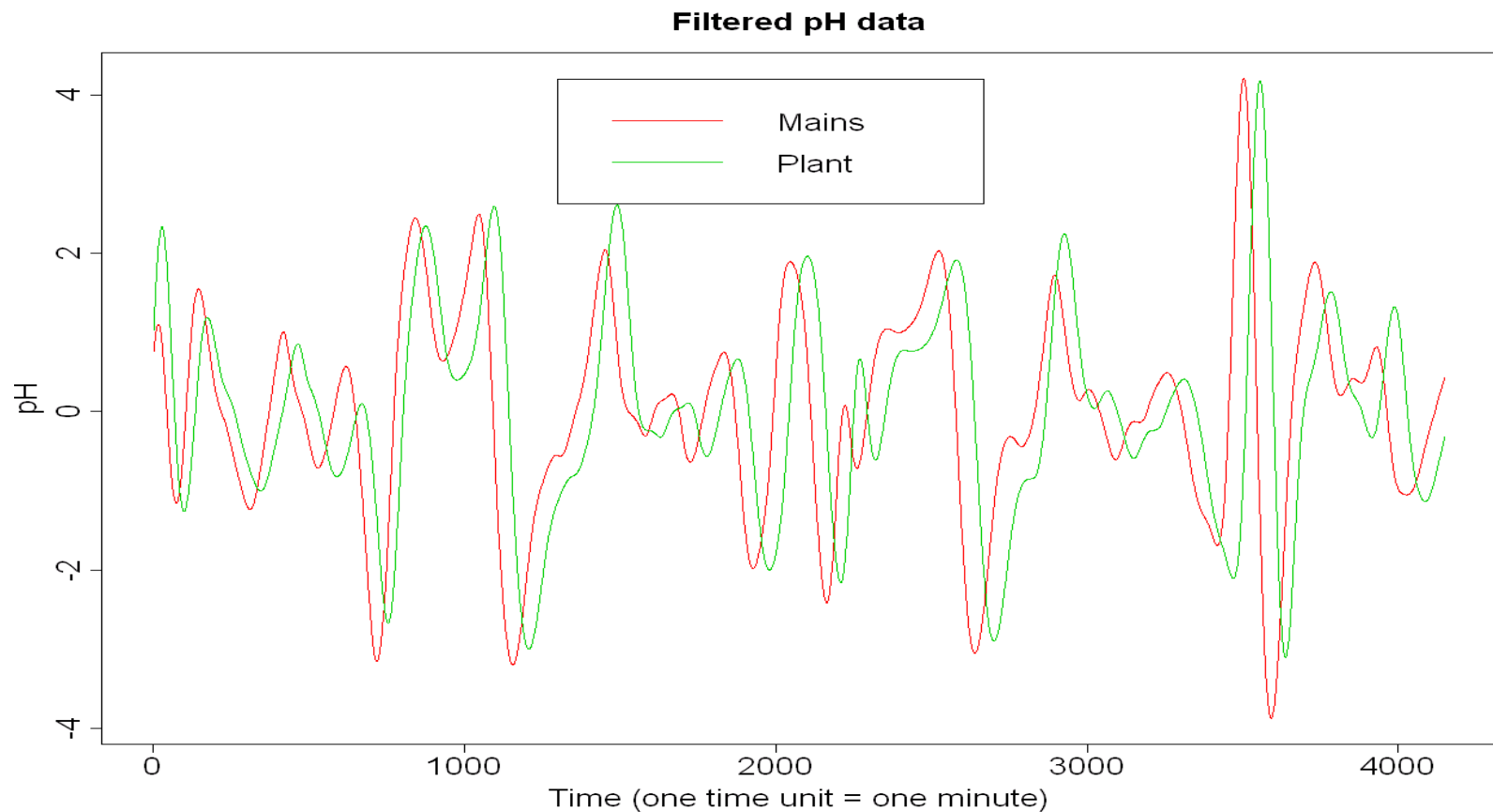- r= 3 , l= 4

# Learnings from the data

- **Control charting techniques**
  - Simplistic, but do not work well with data which shows significant drift over time. Some success when applied to "differences".
- **Time series analysis**
  - Insufficient regular structure for this to be effective.
- **Kalman filter techniques ("State space models")**
  - Well suited to slowly varying processes - the strongest contender
  - Uses past data to assess how well new data conforms to past patterns
  - Can identify anomalies across a range of scales
- **Need to extend this to multiple sensors, but this can only be done if we know 'travel time' between sensors**

# Estimation of "Travel time"

- For many systems (water distribution, sewers, river networks), "travel time" is important
- In the past, "tracer" studies or complex hydraulic models were used, but they are expensive and only give a "snapshot"
- We use the natural perturbations in the sensor data to estimate travel time in real-time
- It would then be possible to
  - confirm anomalous events and track them through the system,
  - have a better idea of their origin, and
  - measure water age and use this to determine the rate of change of contamination indicators such as free chlorine.
- Viterbi algorithms were tried but perform poorly – they don't use the fact that travel time changes smoothly
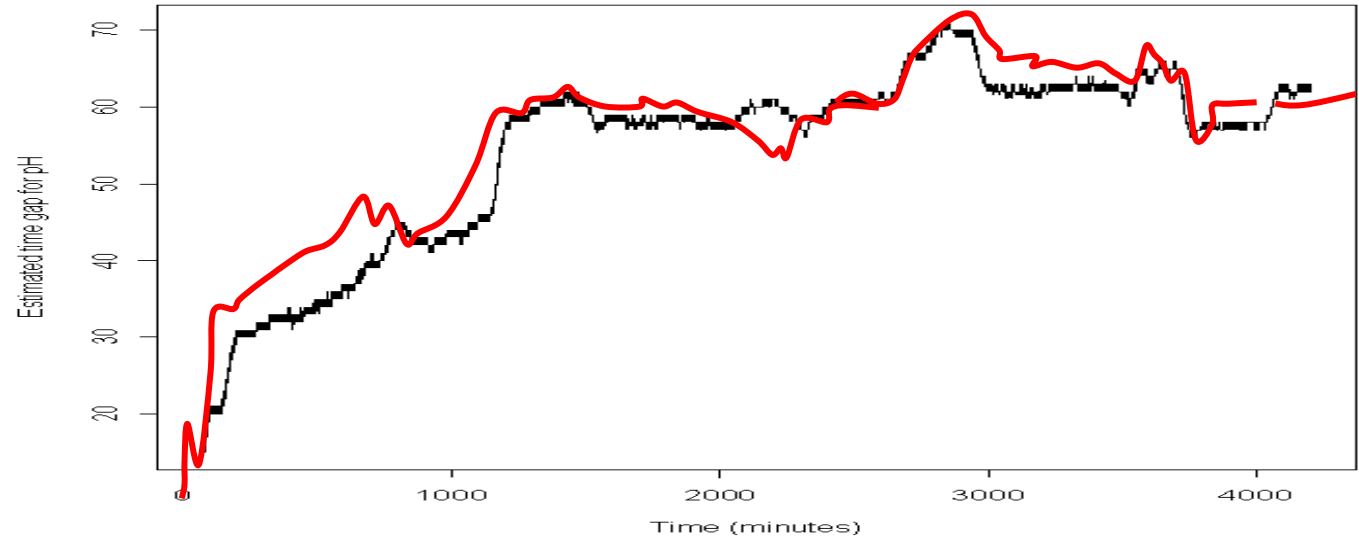- Hidden Markov models for travel time work better

# Example 6: Two sensors 150m apart

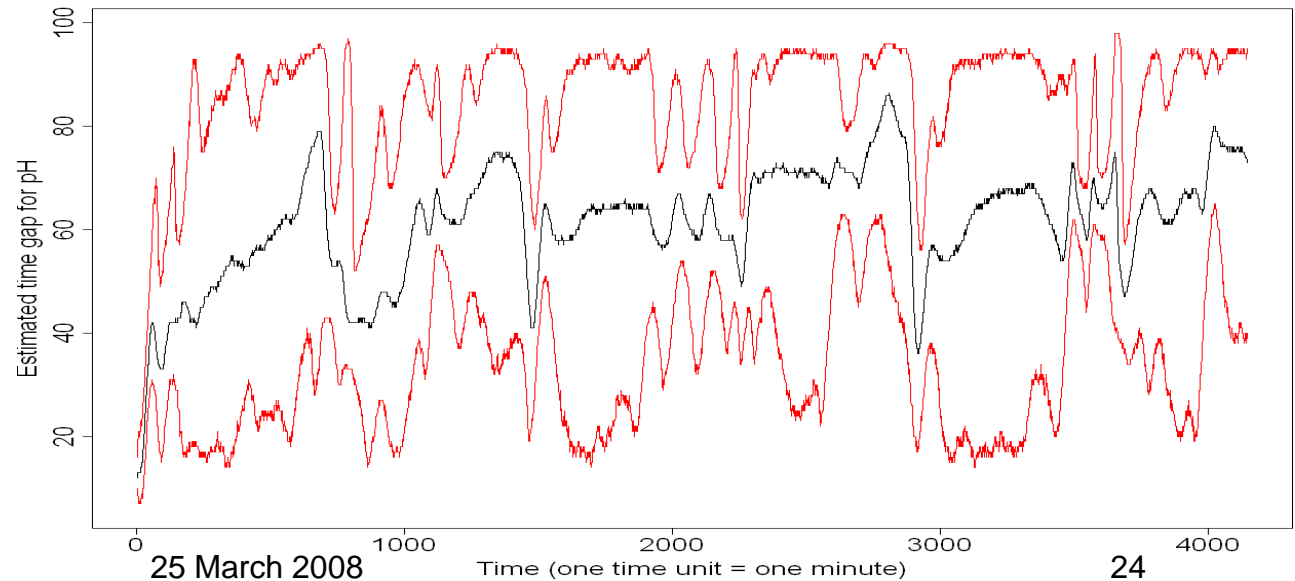- Both sets of data here are filtered to remove high frequency



**Filtered pH data**

# Example 6: Two sensors 150m apart



- Above: Lined up "by eye" using pH (–), EC (–)

- Below: MCMC applied to hidden Markov model (with 95% limits)

# Hidden Markov models for "Travel time"

- Two monitoring sensors at X and Y, water travels from X to Y, readings every $z$ minutes.
- Observe $x_i$ and $y_i$ at time $i$ at locations X and Y, respectively.  Suppose that::

$$y_i = \alpha + \beta x_{i- t_i -a} + e_i$$

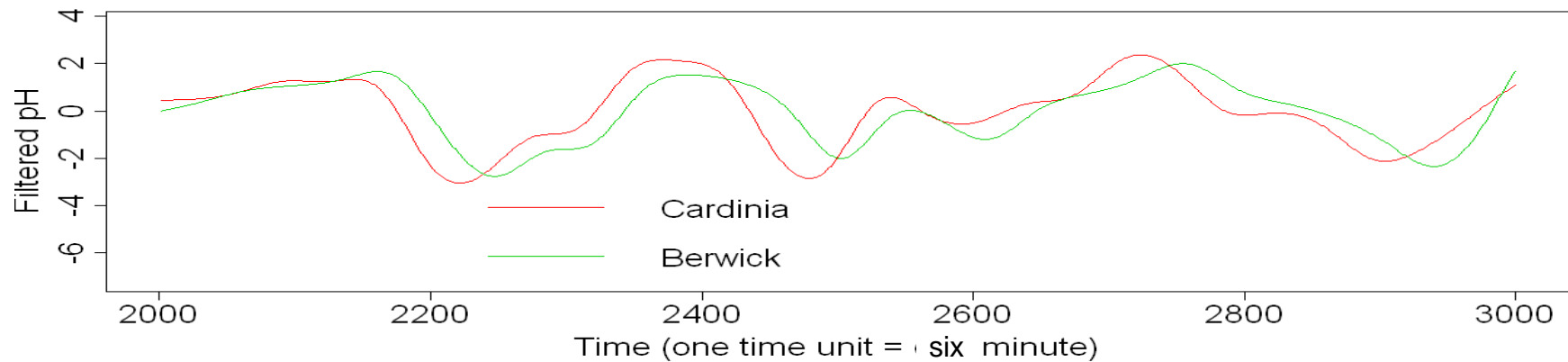where $e_i$ is Normally, zero mean, variance $\sigma^2$,  and $t_i = t_{i-1}+s_i$, with $s_0 = 0$ and

$$
s_i = \begin{cases}
-dz & \text{with probability } p_1 \\
-(d-1)z & \text{with probability } p_2 \\
\ldots \\
0 & \text{with probability } p_{d+1} \\
\ldots \\
+dz & \text{with probability } p_{2d+1}
\end{cases}
$$

- Here, $a$ is the *known* travelling time gap between locations X and Y at time 0 and the *unobserved* $s_i$ are the hidden states.
- Parameter estimation can be undertaken using the EM  algorithm, or using Bayesian methods such as MCMC. The MCMC approach was used here.
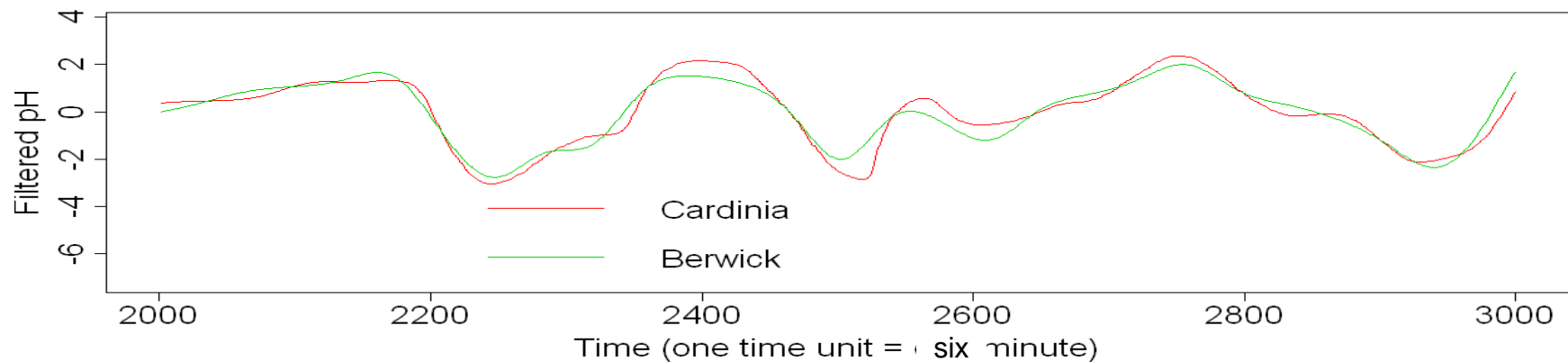
# Example 7: Two sensors 20km apart

- Cardinia (reservoir), only 3% reaches Berwick, 20km away
- Data is filtered, Berwick shifted 6.5h left. Plot covers 4 days.
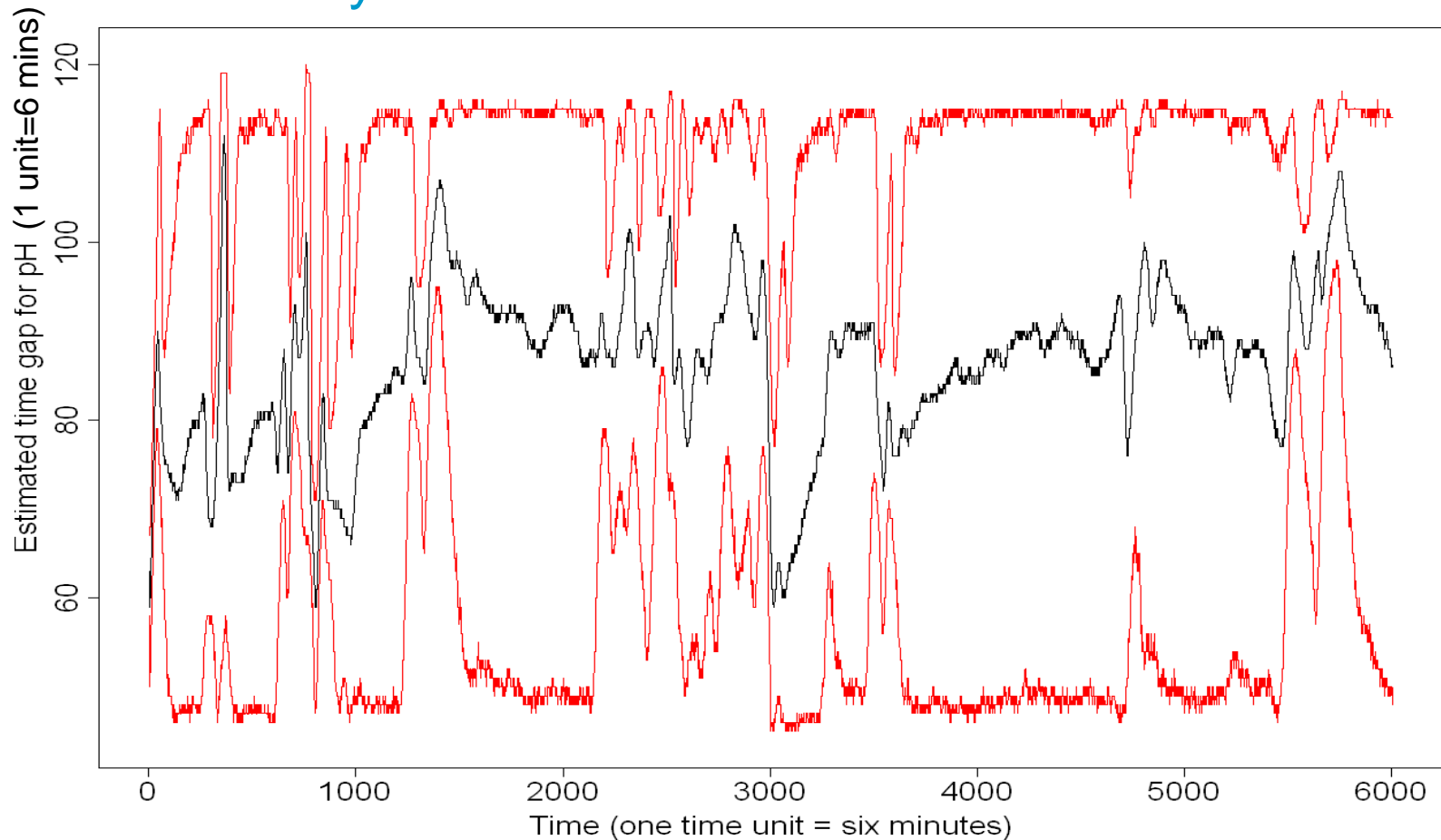


**Filtered pH data**



**Bayesian hidden Markov model output for filtered pH data**

# Example 7: Two sensors 20km apart

- Time delay estimated accurately at some times, poorly at others.
- Plot covers 25 days

# Further work on "Travel time"

- Understanding how water moves through the system is vital for planning purposes (eg chlorine addition)
- Combination of water sources can cause problems with water quality
  - Where we have 'mixed supplies', we hope to identify travel time and the percentage of each water type at each location, in real time
- Research questions:
  - MCMC approach takes ~2h, can we do this in real time?
  - Can we improve accuracy using more variables?
  - How far apart do sensors need to be?
  - How frequently do we need to measure and how much 'natural variation' is needed for this to work?
  - Can we process locally between pairs or triples of sensors, to enable local decision-making?
  - Can we use this to 'track' anomalous events?

# In summary

- Sensor networks offer us huge opportunities/challenges
- These are important problems for water utilities
- We need to:
  - Deal with issues of calibration/maintenance of sensors
  - Build and calibrate models of the system
  - Know how much water we have and where it is and how long it takes to get from A to B
- Potentially massive data sets
- Potential for distributed processing/autonomous decision making

**Richard Jarrett**
**Computational and Mathematical Modelling Program**
**CSIRO Mathematical and Information Sciences**

**Phone:** +61 3 9545 8039
          +61 419 239 452
**Email:**  Richard.Jarrett@csiro.au

# Thank you

CSIRO