



# Mastering the Data Explosion in the Earth and Environmental Sciences

Australian Academy of Science Elizabeth and Frederick White Conference

## Dealing with unknown discontinuities in data and models

**Kerry Gallagher**

**John Stephenson**

**Chris Holmes**

**Imperial College  
London**



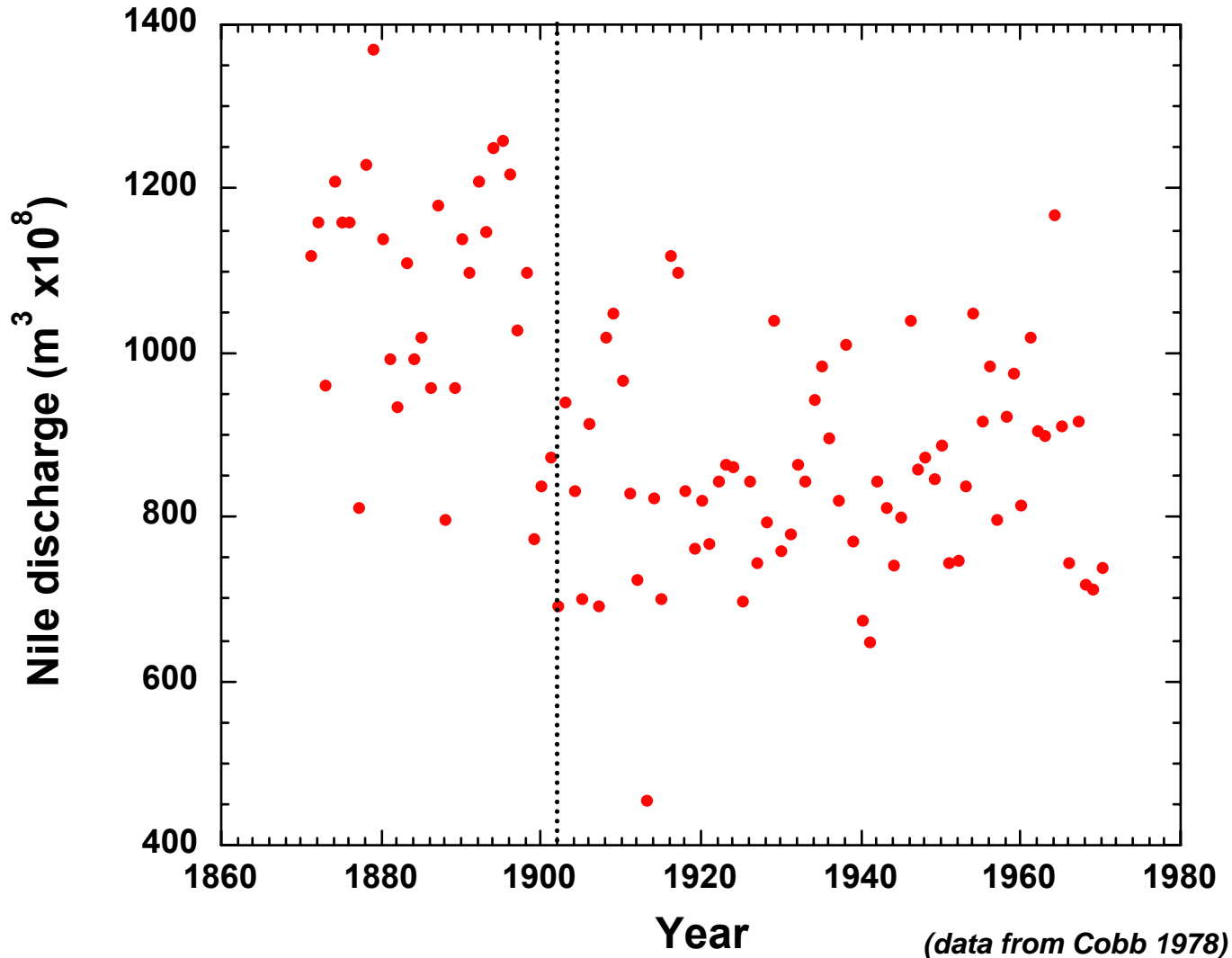
# **Discontinuities occur in both data and processes in the Earth and Environmental Sciences**

*Spatial : faults, topography, lithology, phase,  
composition,...*

*Temporal : climate, seismicity, tectonics,...*

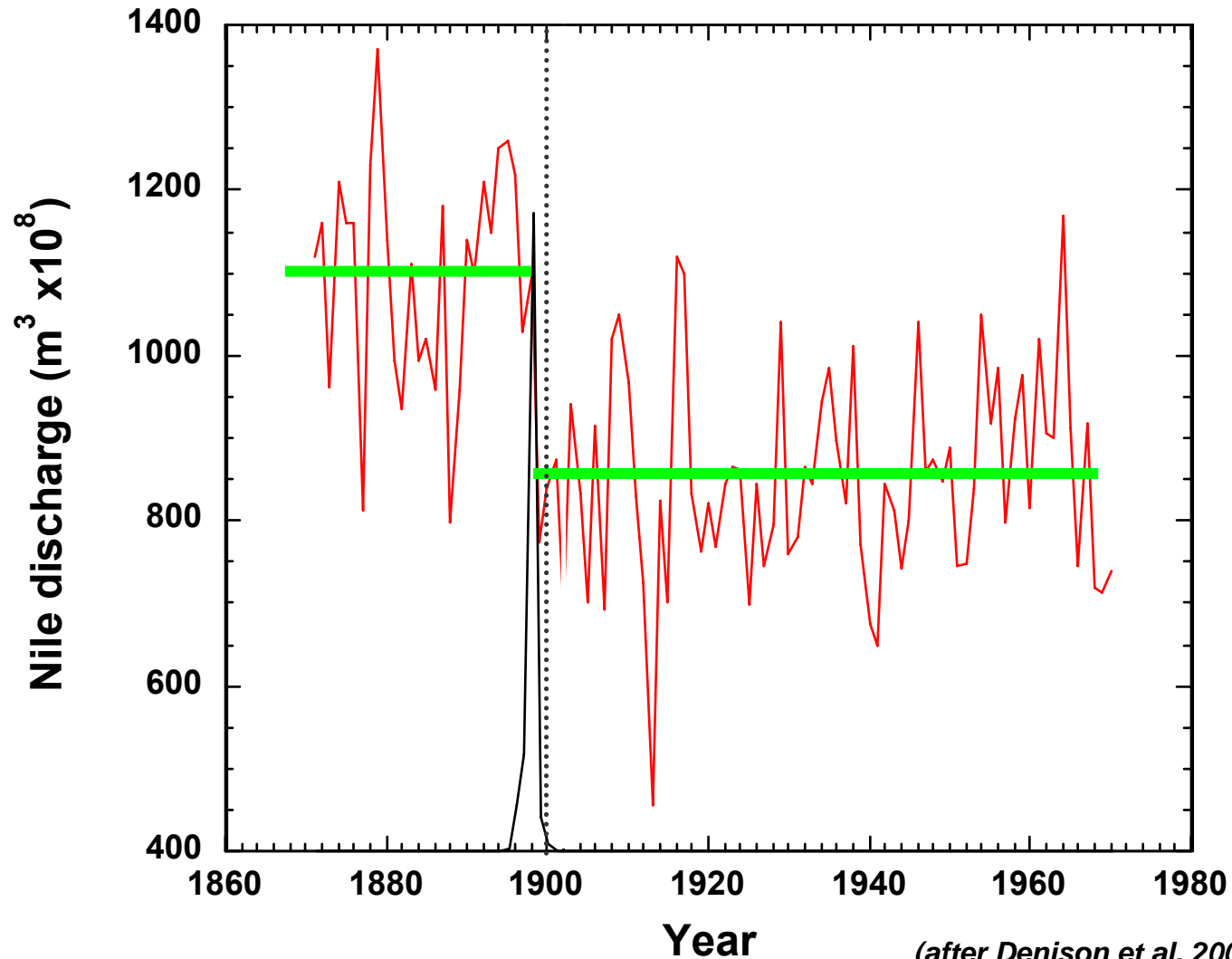
# What is the appropriate question ?

*What was the significance of the opening of the Aswan Dam ?*



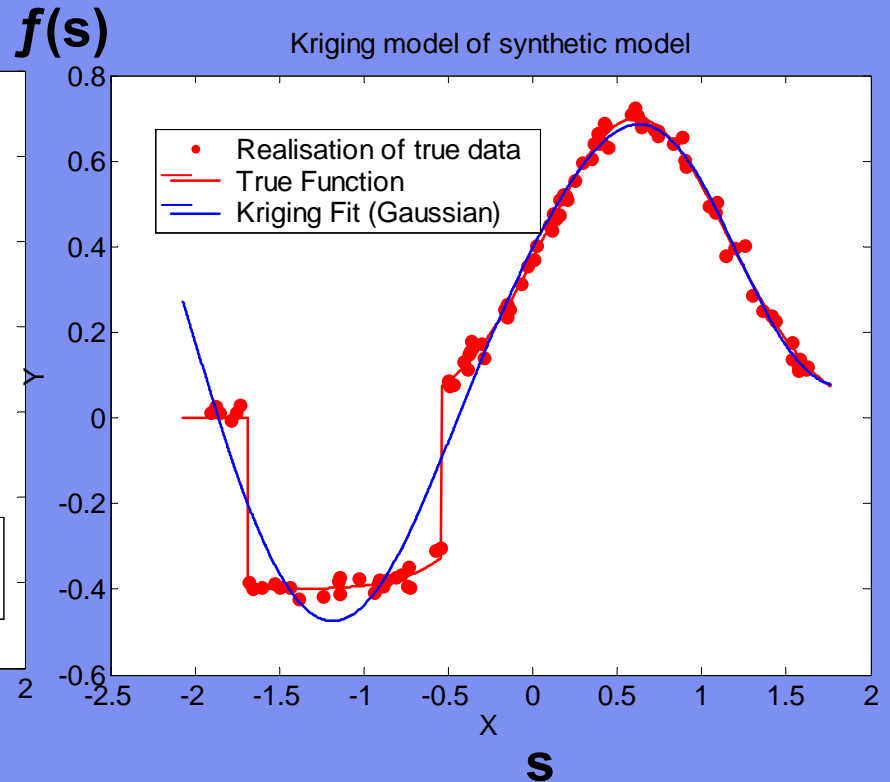
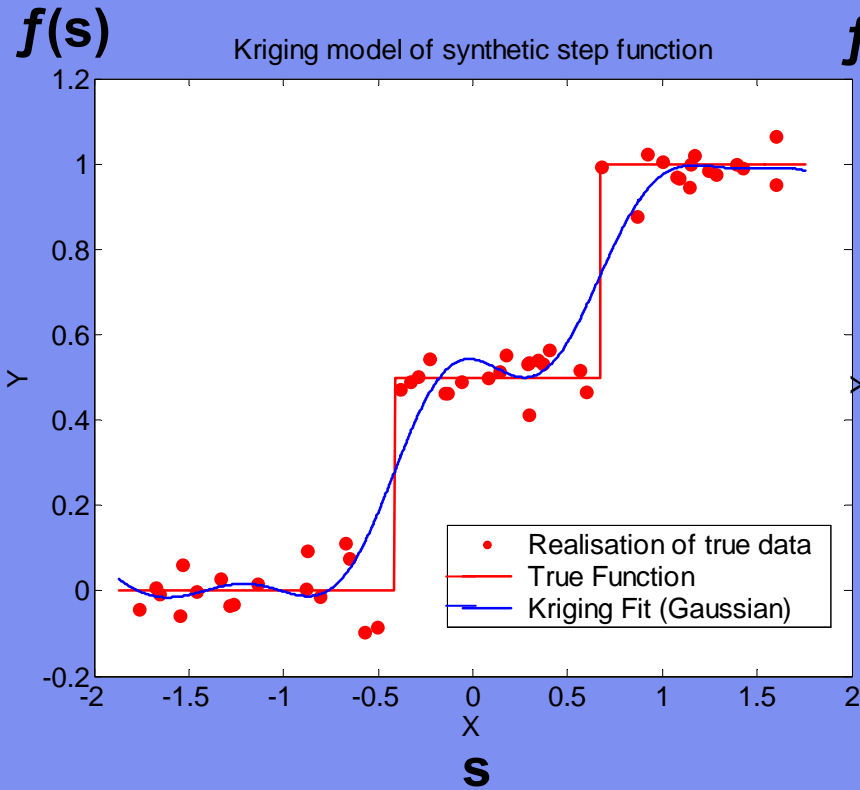
# When was the change ?

$$f(t) = \mu_1 I(t \leq t_c) + \mu_2 I(t > t_c)$$



# Data interpolation and prediction with discontinuities

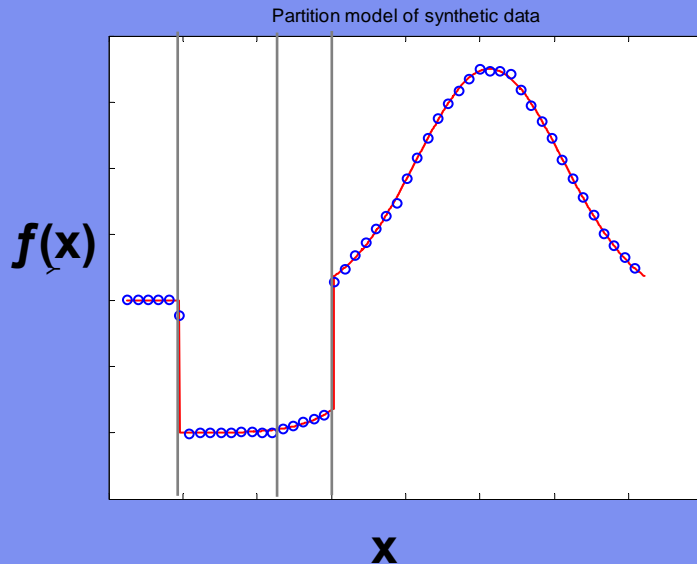
## Standard methods may be too smooth



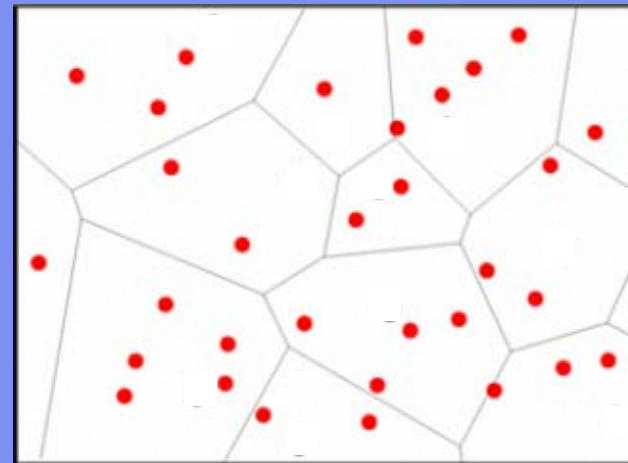
Need a method that can deal with an unknown number of discontinuities in unknown locations

# Partition Modelling

1D

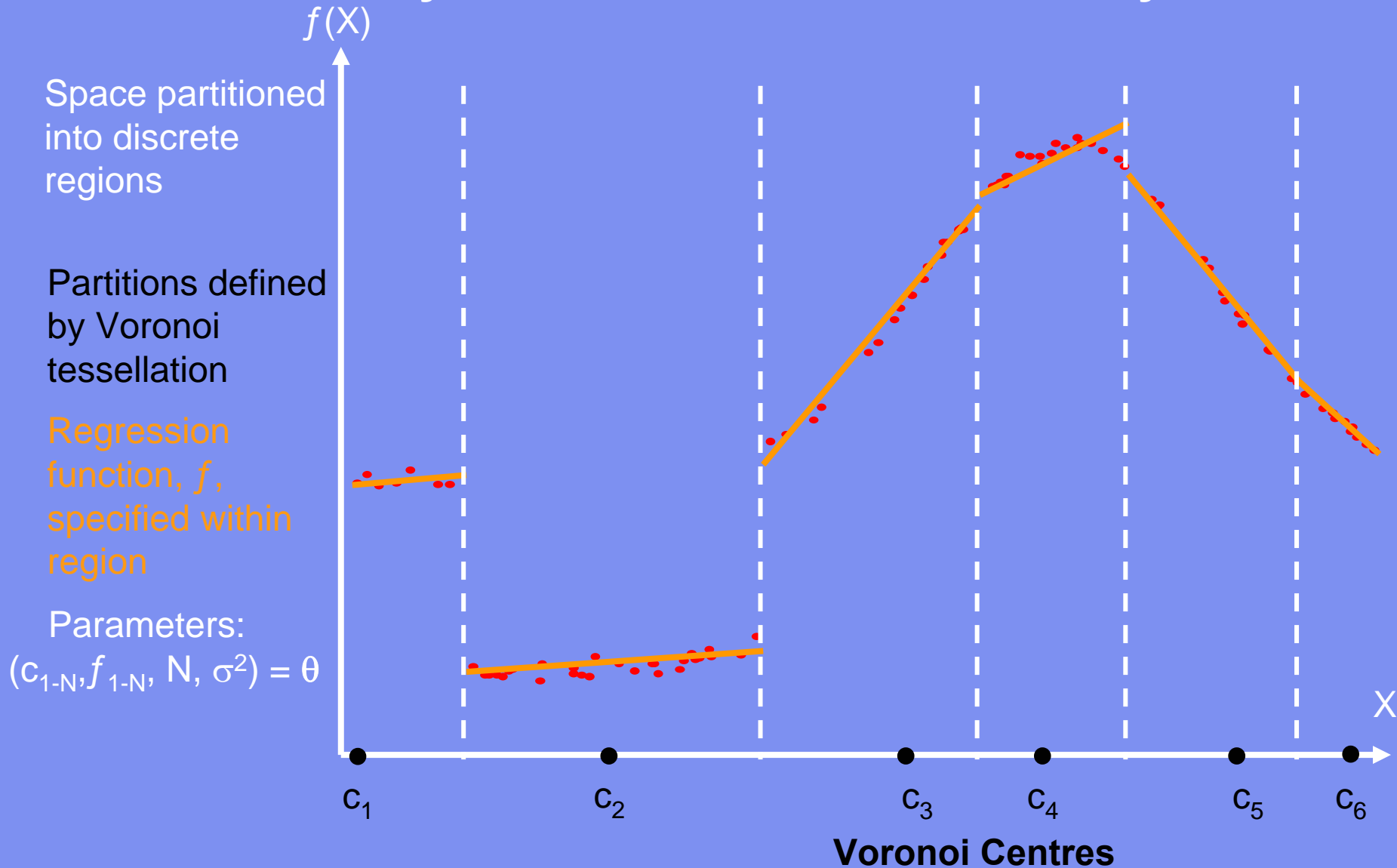


2D



# Formulating a Partition Model

How many discontinuities, where are they ?



# Generating Partition Models

Prediction

$$p(y | D) = \int_{\Theta} p(y | \theta, D) p(\theta | D) d\theta$$

Posterior distribution

Monte Carlo integration

$y$  = value to be predicted

$D$  = observed data

$\theta$  = model parameters

$$p(y | D) \approx \frac{1}{N} \sum_{i=1}^N p(y | \theta_i, D) p(\theta_i | D)$$

Bayes' Theorem

$$p(\theta | D) \propto p(D | \theta) p(\theta)$$

Posterior

Likelihood

Prior

Use Markov chain Monte Carlo (MCMC) to sample the posterior distribution,  $p(\theta|D)$



# Sampling with (transdimensional) MCMC

Initialise  $\theta$

Iterate

- Propose new  $\theta'$
- Calculate likelihood with new  $\theta'$
- Accept new  $\theta'$  or retain current  $\theta$

Acceptance  
criterion

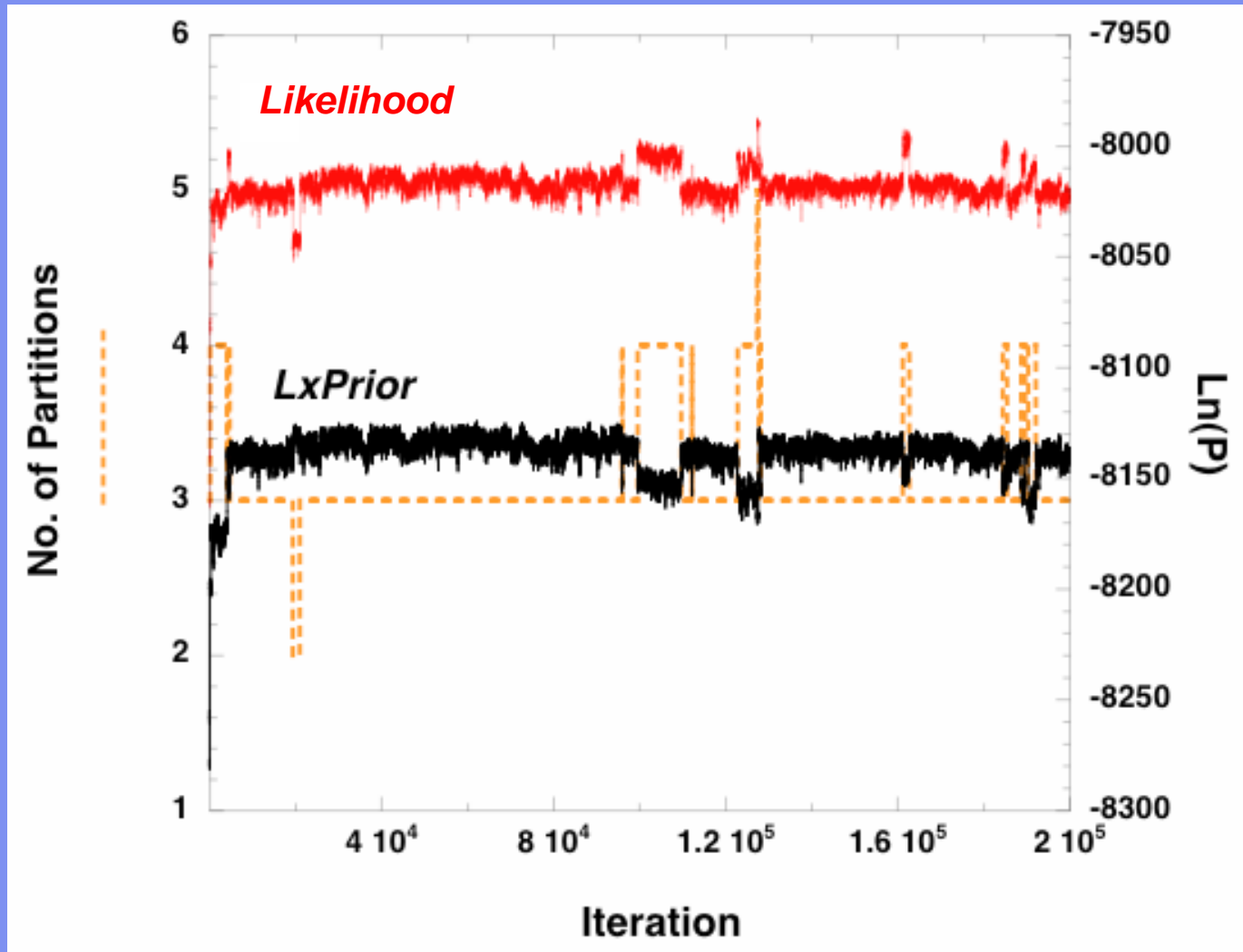
$$\alpha(\theta, \theta') = \min \left\{ 1, \frac{\overbrace{p(\theta')}^{\text{Prior}} \overbrace{p(D|\theta')}^{\text{Likelihood}} \overbrace{p(\theta|\theta')}^{\text{Model Proposal}}}{\overbrace{p(\theta)}^{\text{Prior}} \overbrace{p(D|\theta)}^{\text{Likelihood}} \overbrace{p(\theta'/\theta)}^{\text{Model Proposal}}} \underbrace{R}_{\text{Jump}} \underbrace{|J|}_{\text{Jacobian proposal}} \right\}$$

Jump Jacobian  
proposal

Distribution of accepted models  $\theta \sim p(\theta/D)$

# Sampling Partition Models

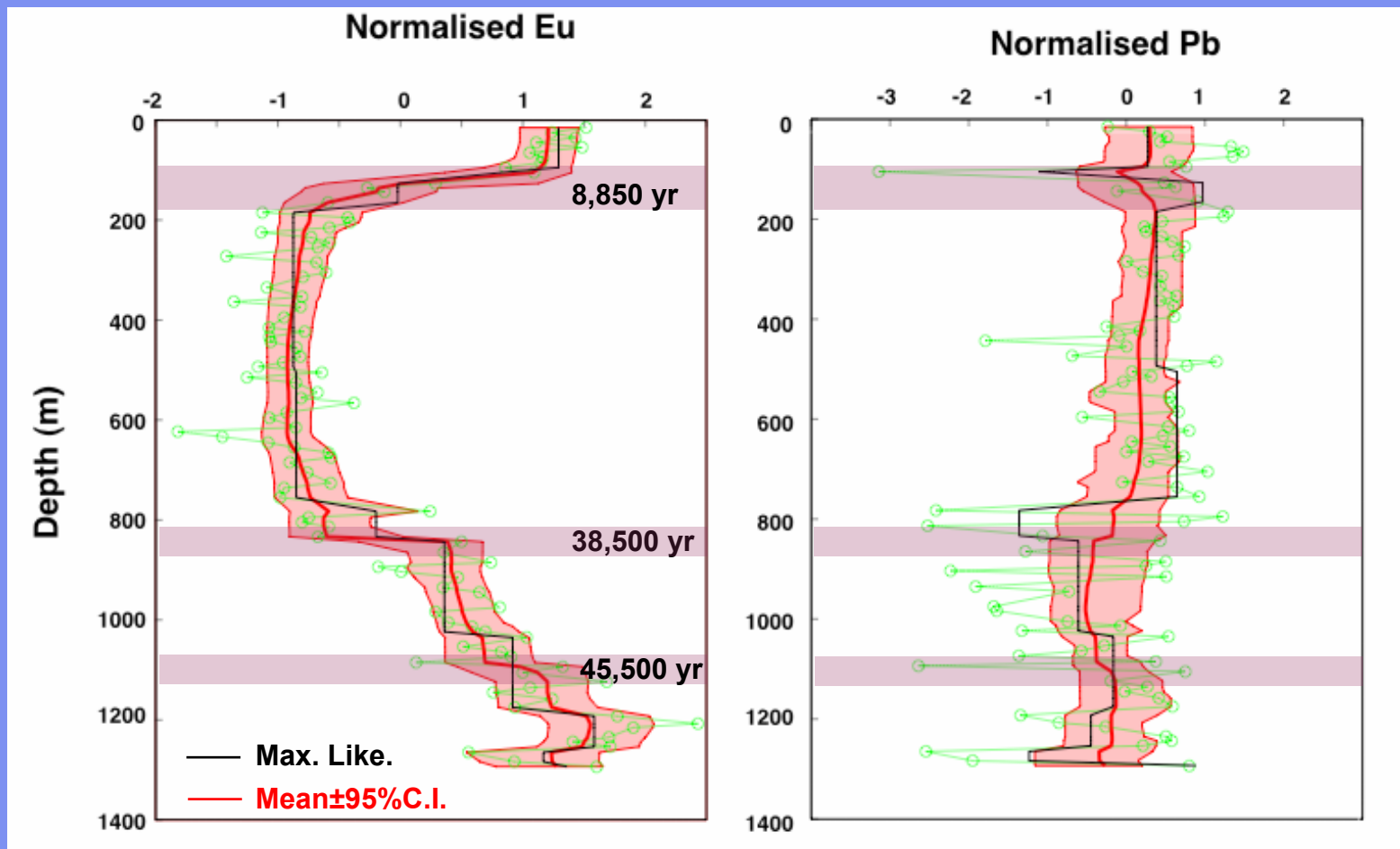
## natural parsimony



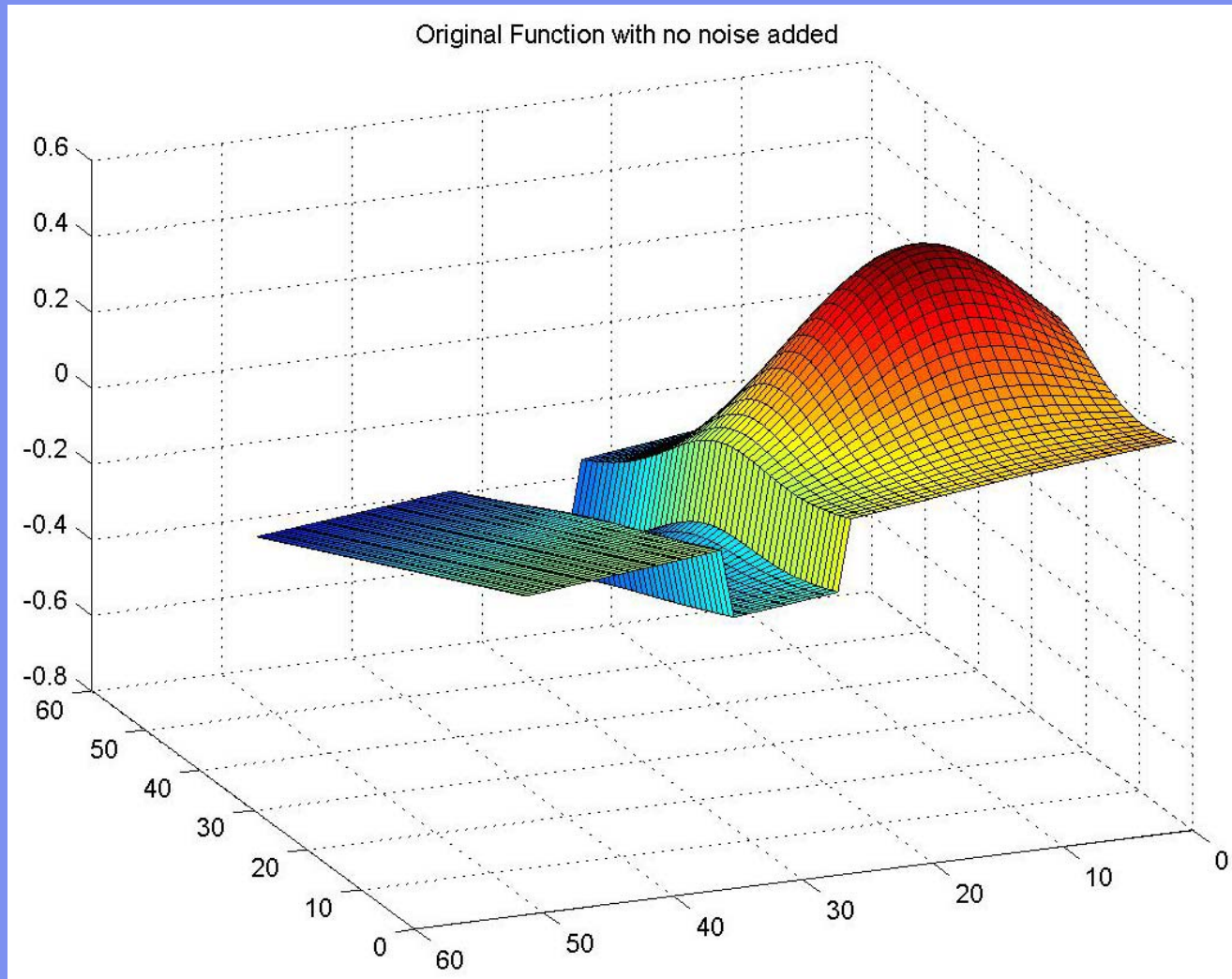
# 1D partition models for data interpolation

## Atmospheric dust input to peat bogs

Looking for common signature in multiple systems



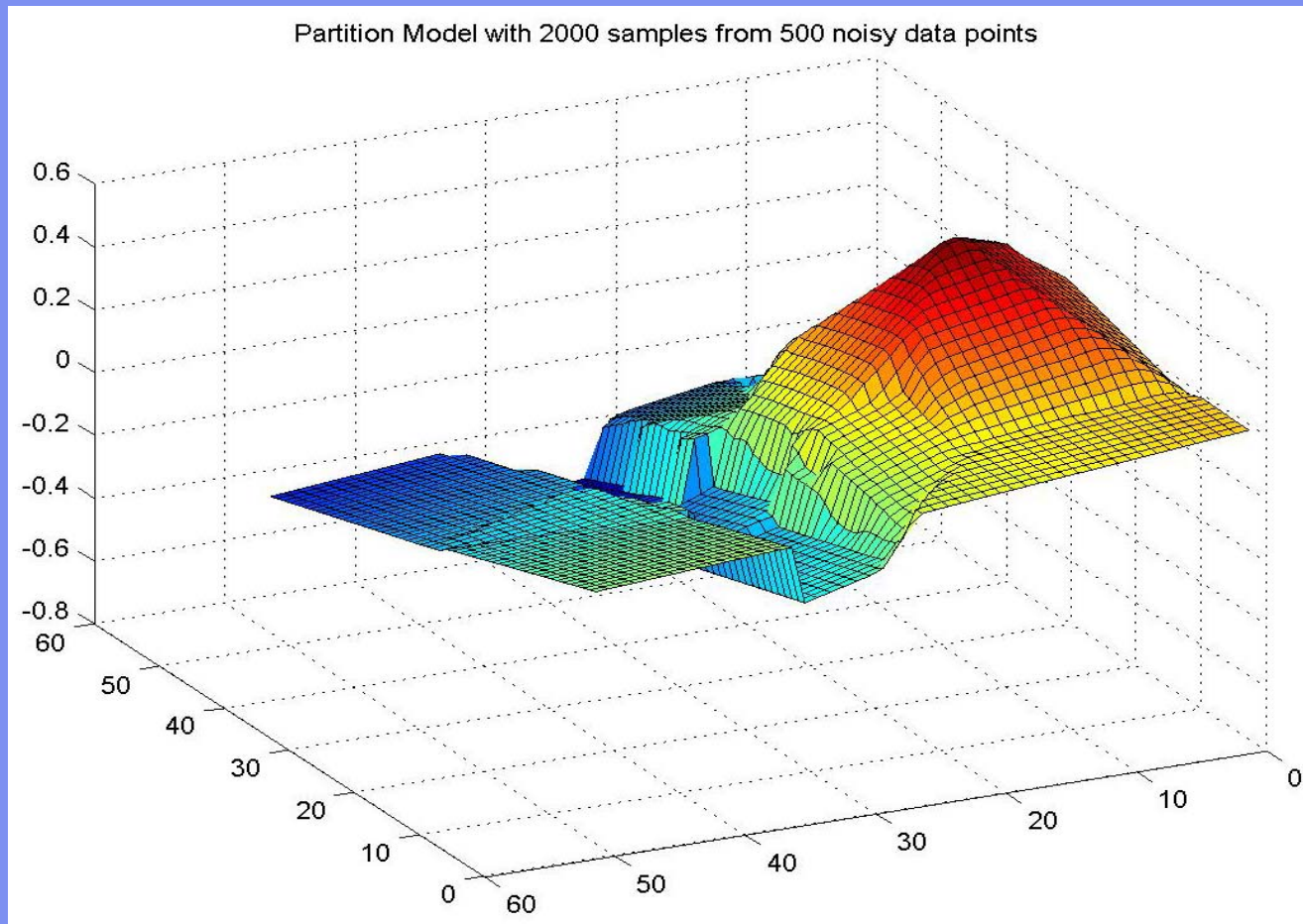
# Partition Models – 2D example function



# Partition Sampling – 2D single realisation

Multiple realisations ...

ensemble average (smooth, but maintain discontinuities)

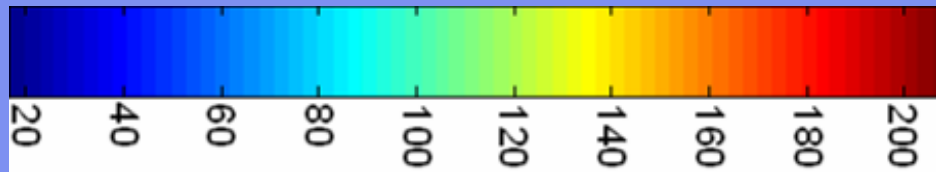
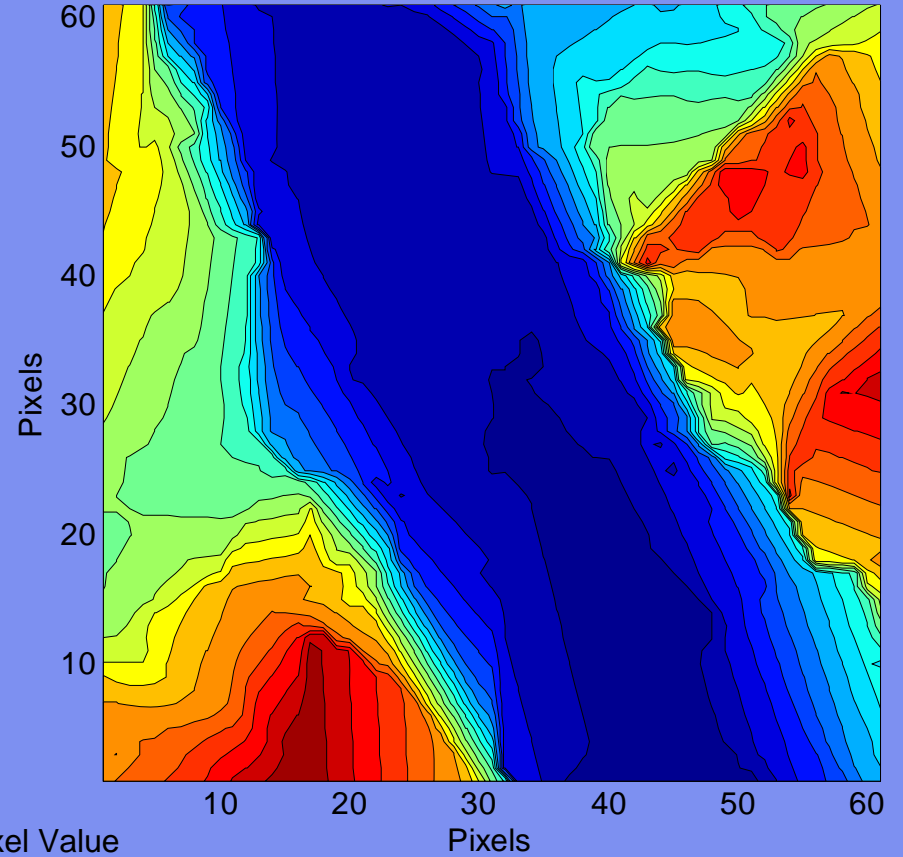
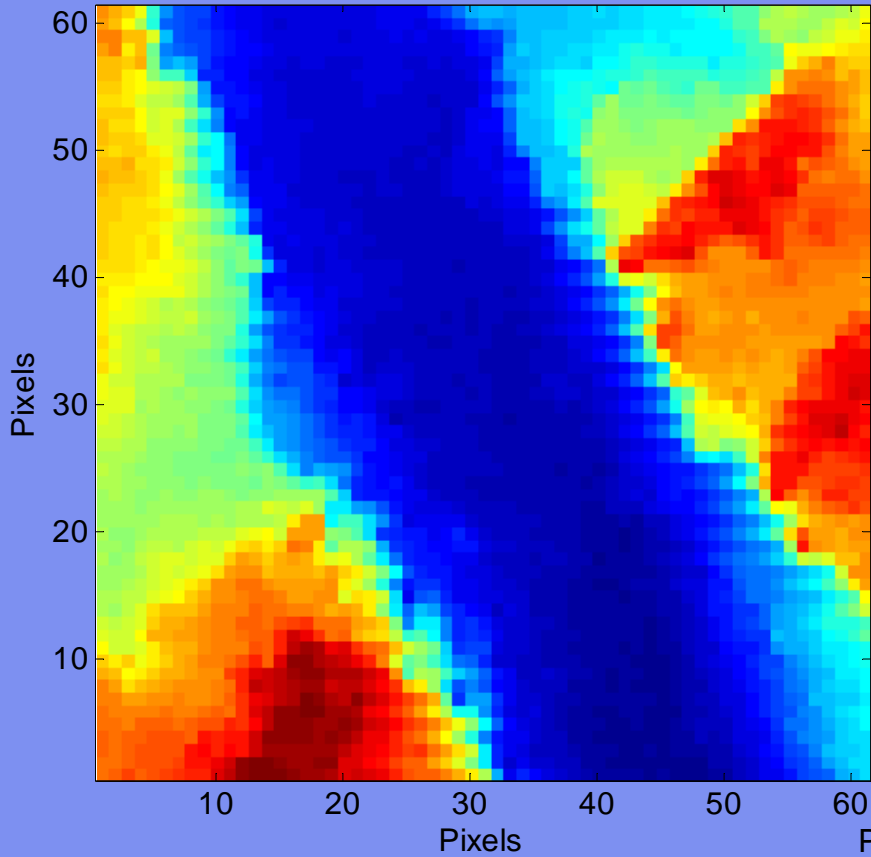


# Partition Model

## Digital Elevation Model (DEM) example

Raw ERS Sample Image

Contour Plot of Partition Model



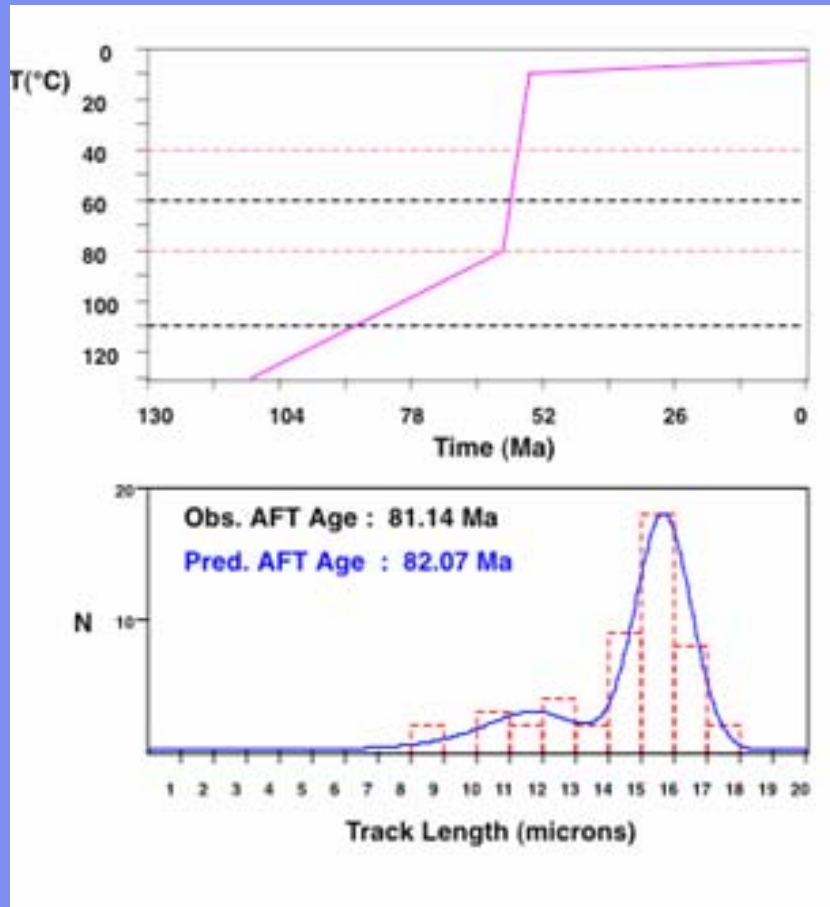
# **Partition Models**

**Application to spatially variable  
physical processes and parameters**

**Example from thermochronology**

# Thermochronology : data are sensitived to temperature history experience by host rock

e.g. apatite fission track analysis

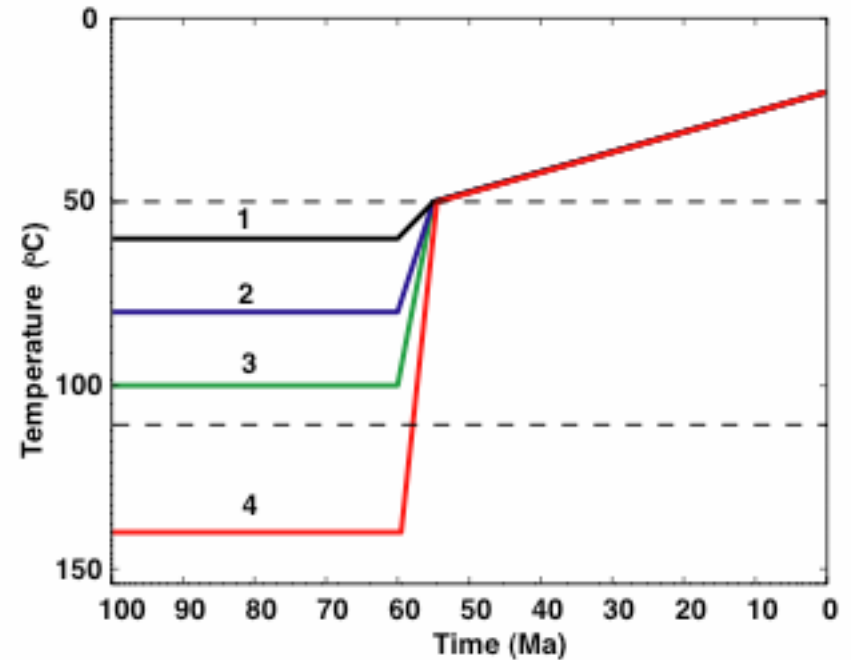
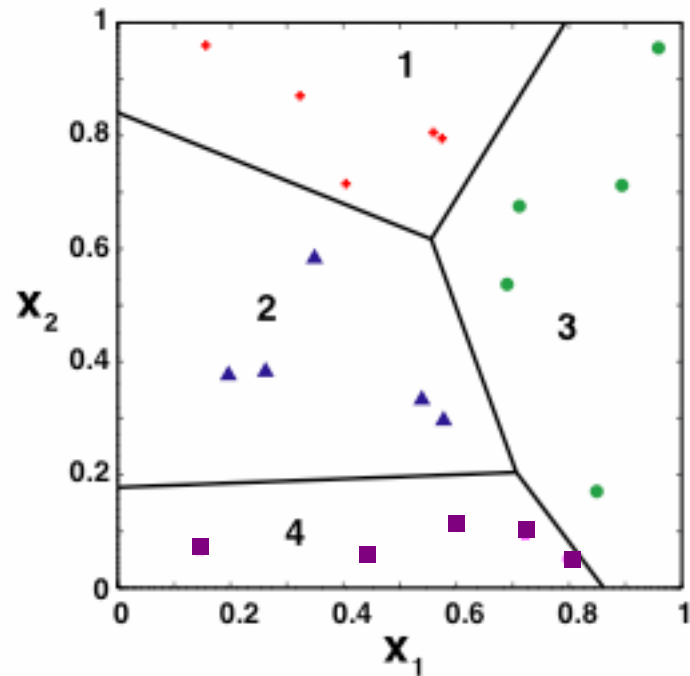


$$p(D|\theta) = f(T(t), \phi)$$

Likelihood is a non-linear function of unknown parameters at each location within each partition



# Model partition distribution and thermal histories



**The problem is to find**

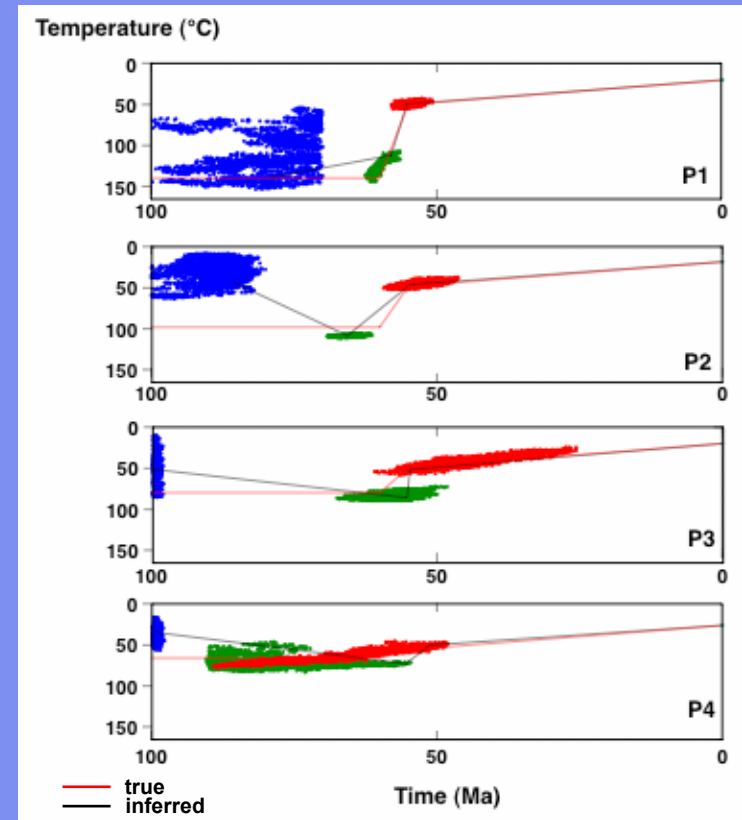
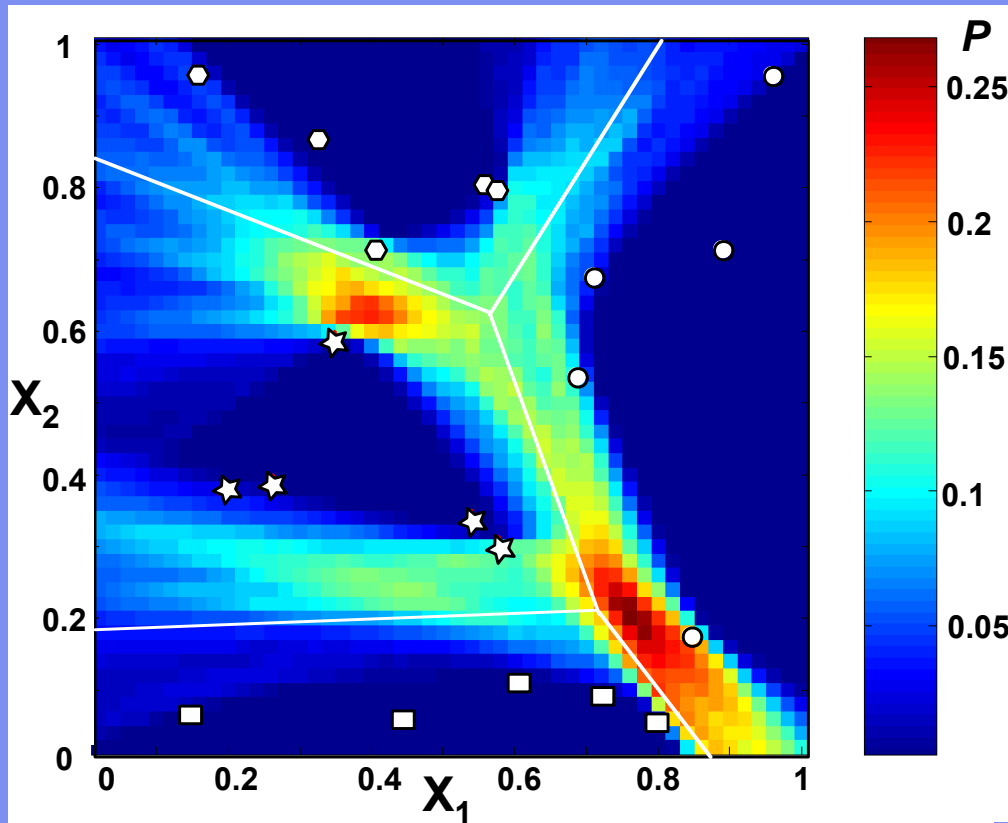
**(a) how to partition the samples in 2D**

**(i) number of partitions**

**(ii) location of the partitions**

**(b) the distribution of thermal histories  
in each partition**

# Inferred partition distribution and thermal histories



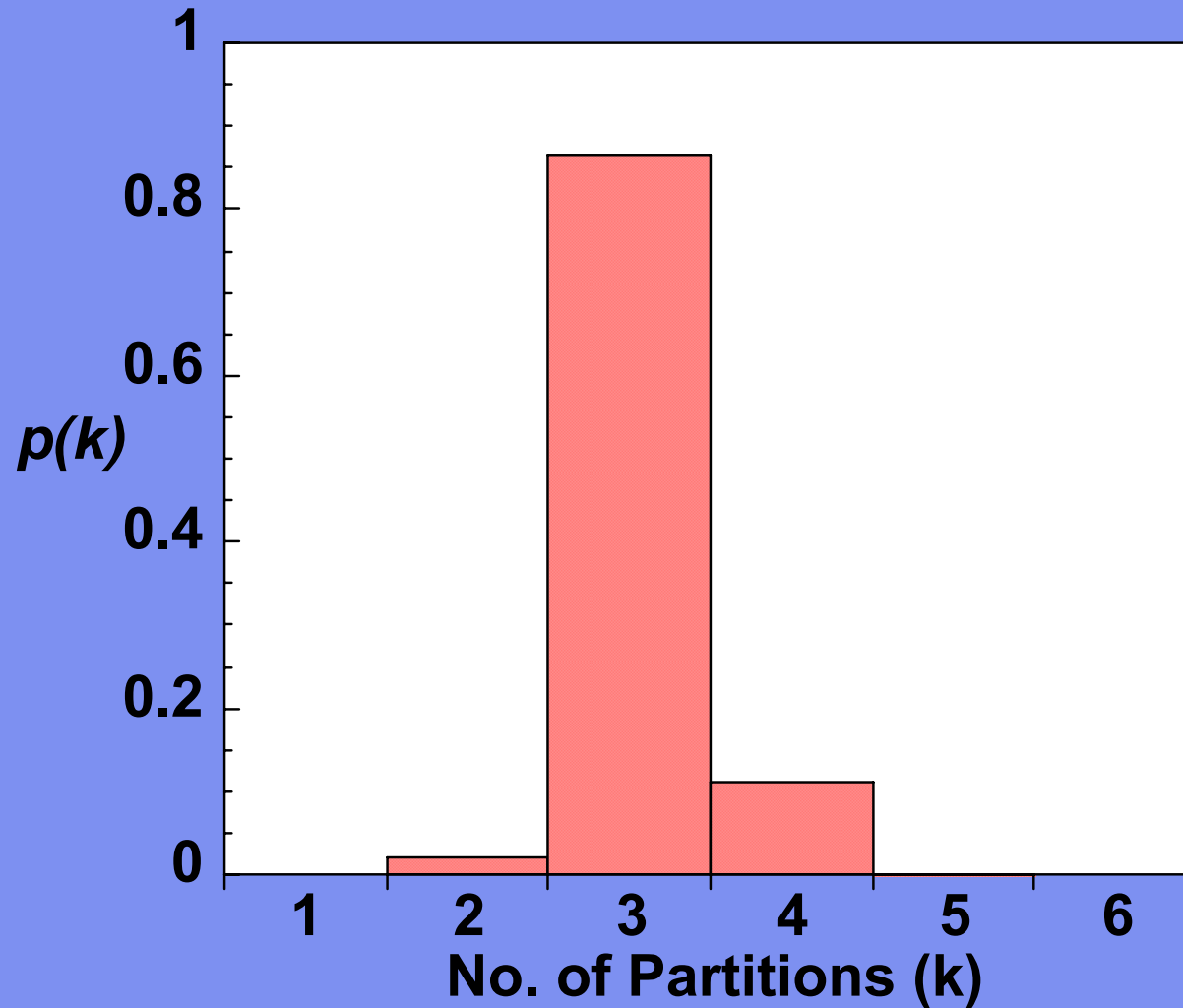
(Stephenson, Gallagher and Holmes 2006)

# Summary

- **Partition models allow for unknown number of discontinuities with unknown geometry in variable dimensions**
- **Bayesian approach deals with the problem in terms of probabilities...intuitive for model choice**
- **Obtain probability distributions (partitions, model parameters, posterior predictions)**
- **Bayesian approach is naturally parsimonious**
- **Potential for self-adaptive/self regularising model parameterisation**

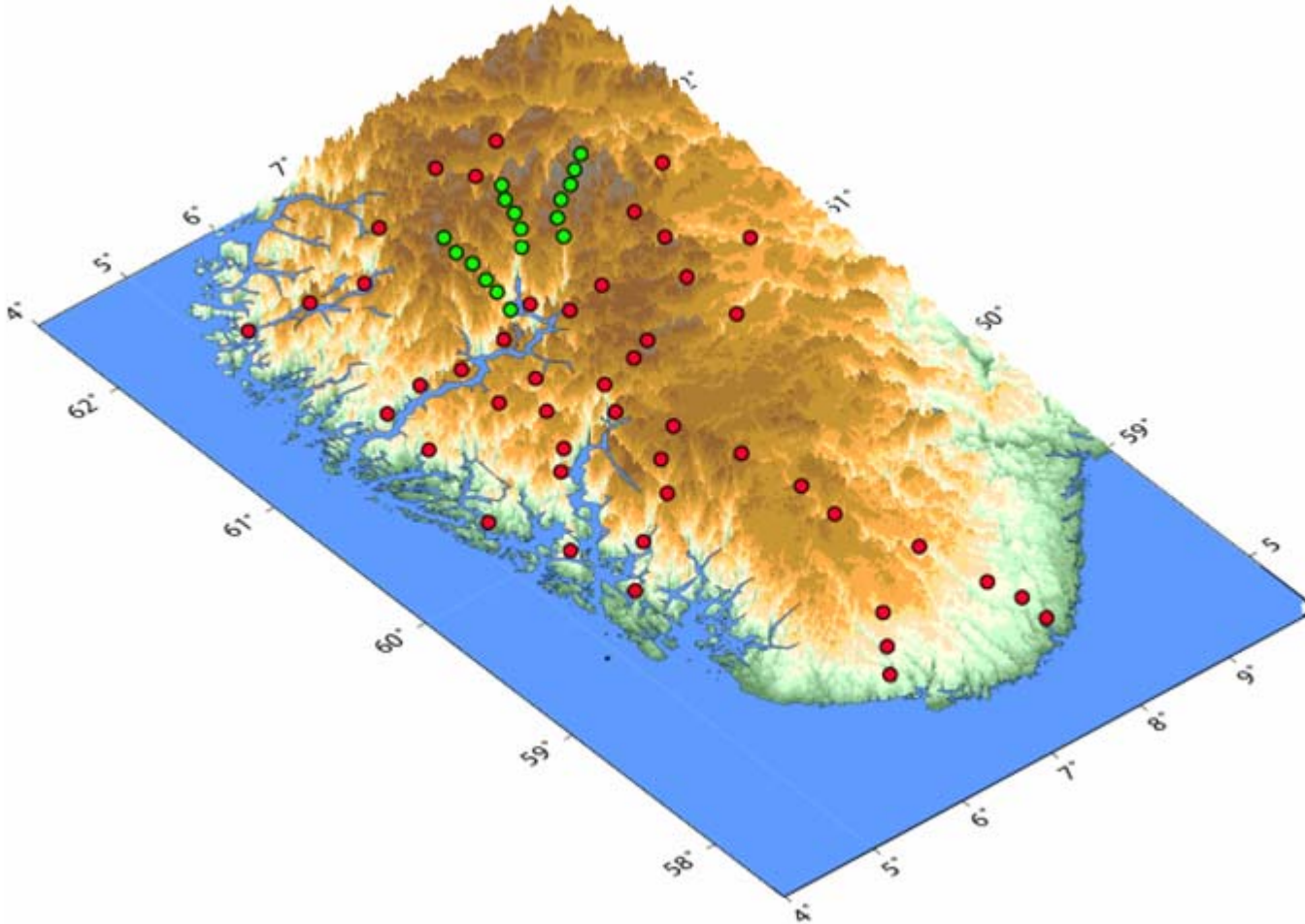
# Sampling Partition Models

## distribution on number of partitions



**Traditionally, each sample is modelled independently..**

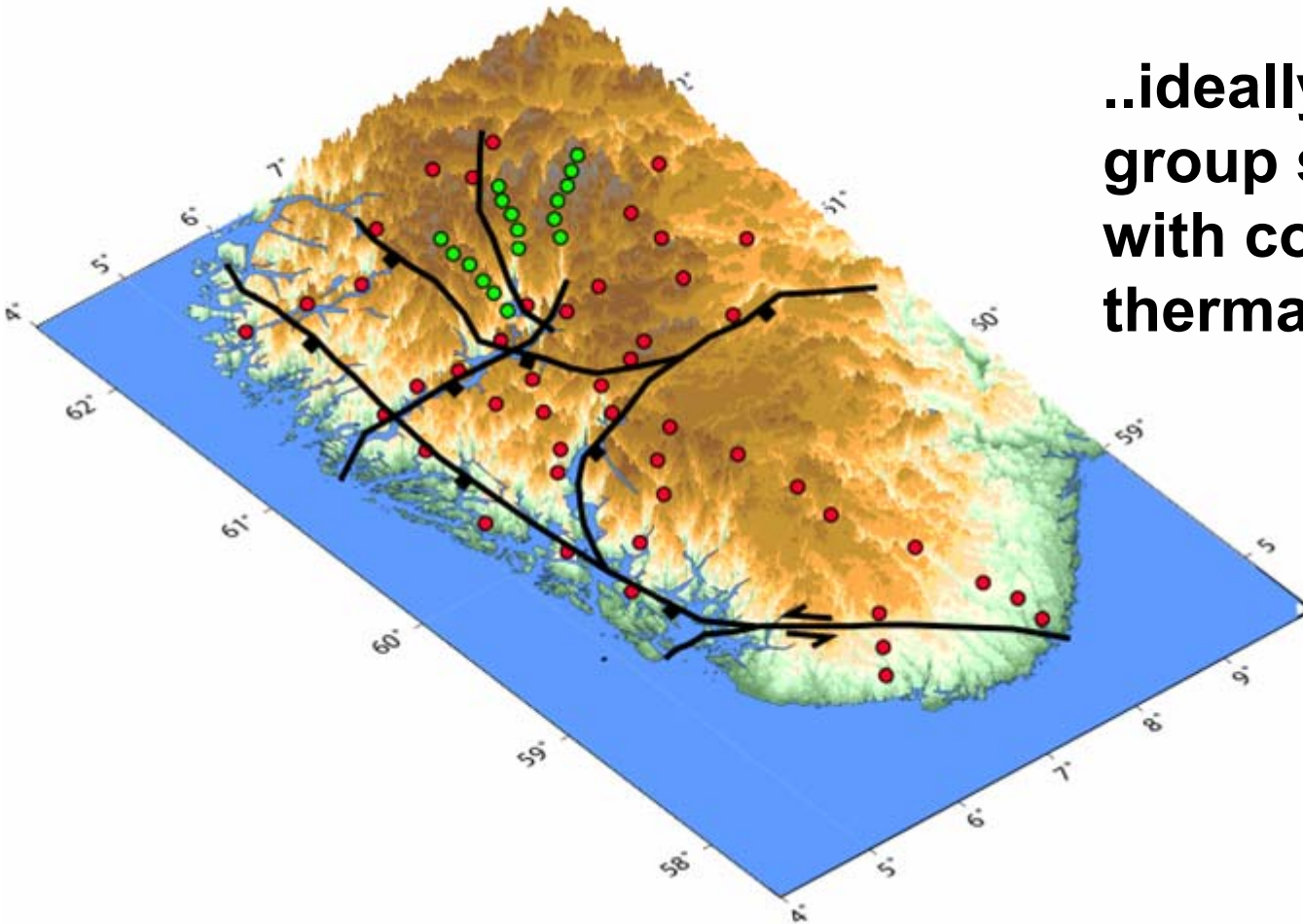
**ignores spatial relationships....**



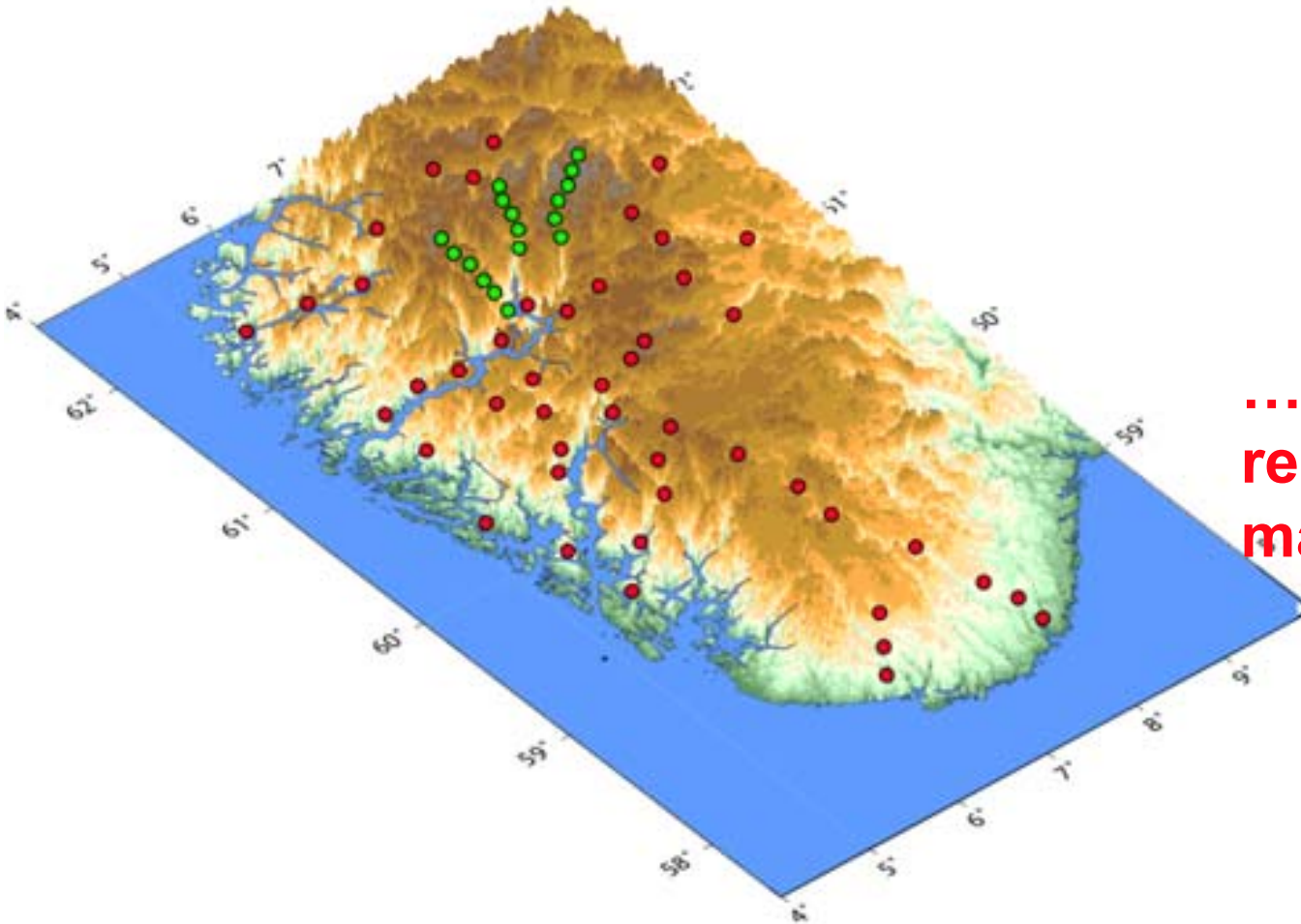
**Traditionally, each sample is modelled independently..**

**ignores spatial relationships....**

**..ideally want to group samples with common thermal history**



**Traditionally, each sample is modelled independently..**



**...but the spatial relationships may be unknown...**