

# Structure and Uncertainty

Graphical modelling  
and complex stochastic systems

Peter Green (University of Bristol)  
2-13 February 2009

## What has statistics to say about science and technology?

2

## Statistics and science

Ernest Rutherford (1871-1937)

3

## What has statistics to say about the complexity of modern science?

4

## Gene networks

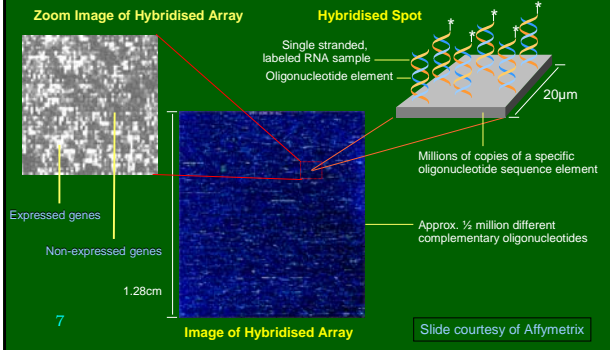
5

## Functional categories of genes in the human genome

Category	Count	Percentage
GO categories	12808	41.7%
molecular function unknown	12808	41.7%
transcription factor	1850	6.0%
kinase	668	2.8%
receptor	1343	5.0%
signaling molecule	376	1.2%
nucleic acid enzyme	2308	7.5%
cytokine	150	0.5%
cytoskeletal structural protein	876	2.8%
extracellular matrix	1437	4.6%
membrane channel	1406	4.5%
ion channel	1406	4.5%
motor	276	0.9%
muscle protein of muscle	296	1.0%
proton pump	902	2.9%
select cell-cell binding protein	24	0.1%
intracellular transporter	350	1.1%
transporter	533	1.7%
select regulatory molecule	98	0.3%
transformase	610	2.0%
glycolytic and cytolitic	313	1.0%
nucleosidase	104	0.3%
lyase	117	0.4%
ligase	356	1.1%
isomerase	1163	3.7%
ly-diolase	1227	4.0%
cell adhesion	571	1.8%
microfilament	338	1.1%
viral protein	100	0.3%
transfer/cancer protein	203	0.7%
transcription factor	1850	6.0%
chaperone	150	0.5%
cytoskeletal structural protein	876	2.8%
extracellular matrix	1437	4.6%
membrane channel	1406	4.5%
ion channel	1406	4.5%
motor	276	0.9%
muscle protein of muscle	296	1.0%
proton pump	902	2.9%
select cell-cell binding protein	24	0.1%
intracellular transporter	350	1.1%
transporter	533	1.7%

Venter et al, Science, 16 February, 2001

## Gene expression using Affymetrix microarrays



7



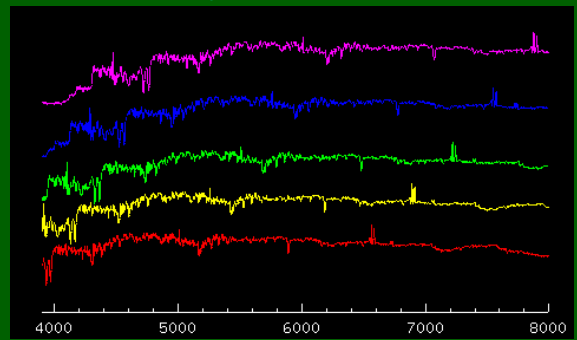
8

## Velocity of recession determines 'colour' through redshift effect

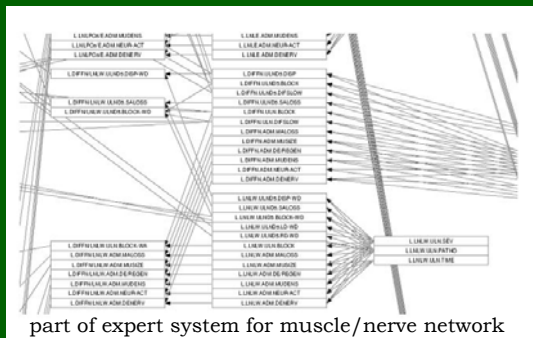


9

## Astronomy: redshifts



## Probabilistic expert systems

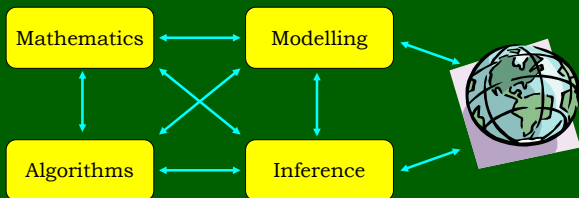


## Complex stochastic systems

Problems in these areas – and many others – have been successfully addressed in a modern statistical framework of structured stochastic modelling

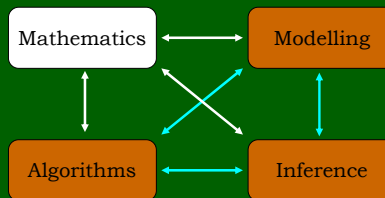
12

## Graphical modelling



13

## 1. Mathematics



14

## Conditional independence

- $X$  and  $Z$  are conditionally independent given  $Y$  if, knowing  $Y$ , discovering  $Z$  tells you nothing more about  $X$ :  $p(X | Y, Z) = p(X | Y)$
- $X \perp Z | Y$



15

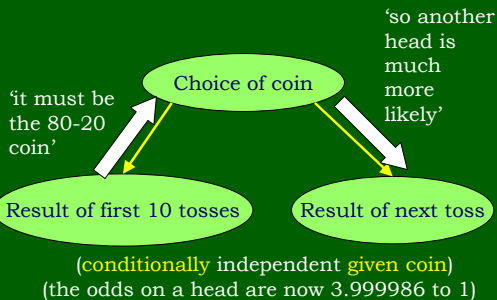
## Coin-tossing

- You take a coin from your pocket, and toss it 10 times and get 10 heads
- What is the chance that the next toss gives head?

Now suppose there are two coins in your pocket – a 80-20 coin and a 20-80 coin – what is the chance now?

16

## Coin-tossing



17

## Conditional independence

as seen in data on perinatal mortality vs. ante-natal care....

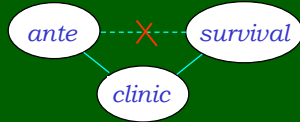
Clinic	Ante	Survived	Died	% died
A	less	176	3	1.7
	more	293	4	1.3
B	less	197	17	7.9
	more	23	2	8.0

Does survival depend on ante-natal care?

18

.... what if you know the clinic?

## Conditional independence

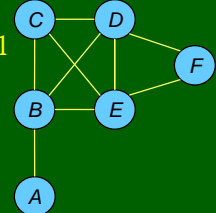


*survival* and *clinic* are dependent  
and *ante* and *clinic* are dependent  
but *survival* and *ante* are CI given *clinic*

19

## Graphical models

- Use ideas from graph theory to
- represent structure of a joint probability distribution
  - by encoding conditional independencies

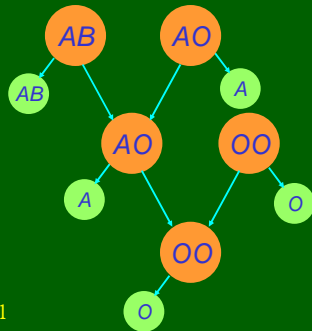


20

## Mendelian inheritance - a natural structured model



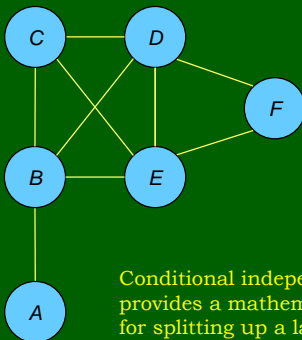
Mendel



## Where does the graph come from?

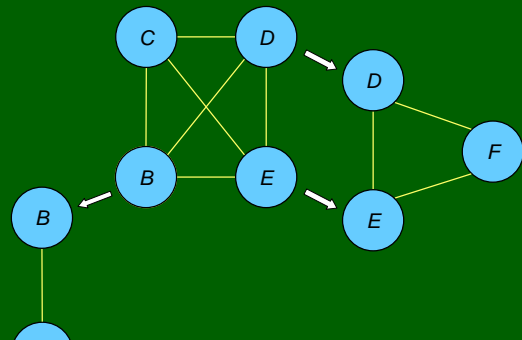
- Genetics
  - pedigree (family connections)
- Physical and biological systems
  - supposed causal effects
- Contingency tables
  - hypothesis tests on data
- Gaussian case
  - graph determined by non-zeroes in inverse variance matrix (i.e. non-zero partial correlations)

22



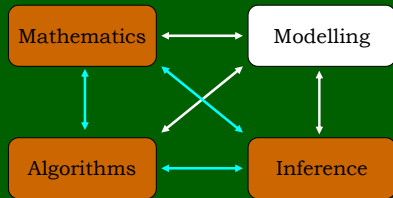
Conditional independence provides a mathematical basis for splitting up a large system into smaller components

23



24

## 2. Modelling



25

## Structured systems

A framework for building models, especially probabilistic models, for empirical data

Key idea -

- understand complex system
- through global model
- built from small pieces
  - comprehensible
  - each with only a few variables
  - modular

26

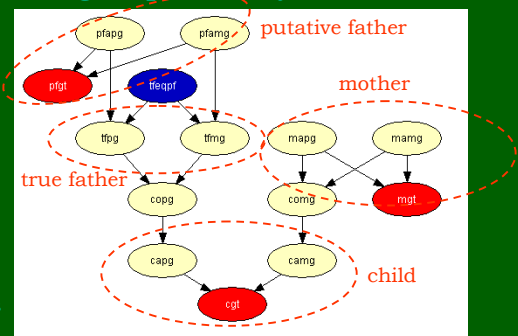
## Modular structure

Basis for

- understanding the real system
- capturing important characteristics statistically
- defining appropriate methods
- computation
- inference and interpretation

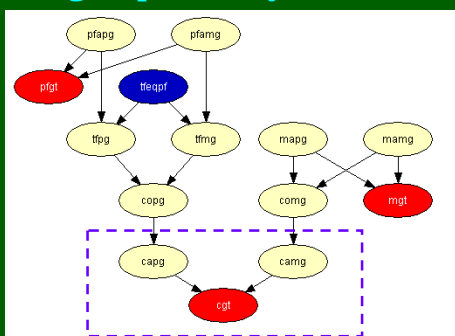
27

## Building a model, for genetic testing of paternity using DNA probes



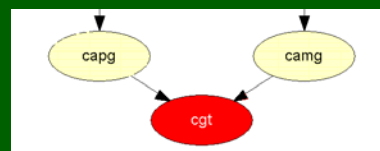
28

## Building a model, for genetic testing of paternity



29

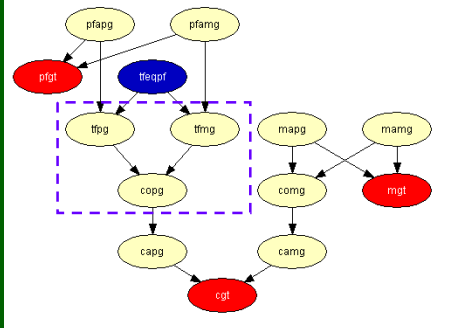
## ... genes determine genotype



e.g. if child's paternal gene is '10' and maternal gene is '12', then its genotype is '10-12'

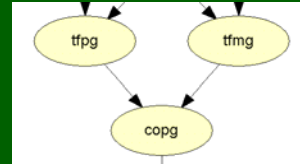
30

## Building a model, for genetic testing of paternity



31

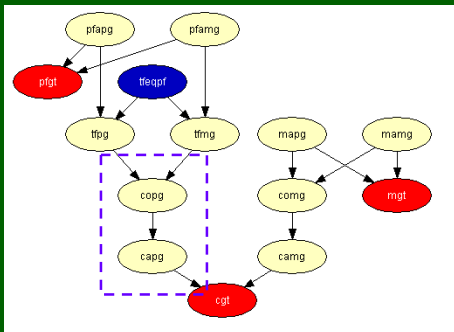
## ... Mendel's law



the gene that the child gets from the father is equally likely to have come from the father's father or mother

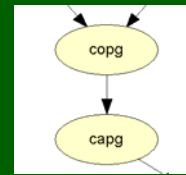
32

## Building a model, for genetic testing of paternity



33

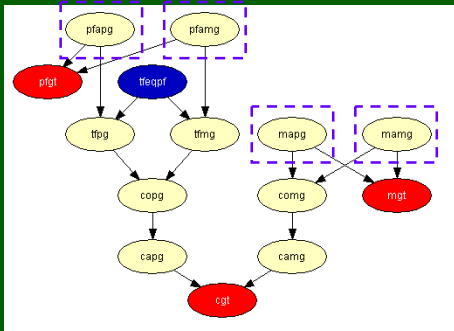
## ... with mutation



there is a small probability of a gene mutating

34

## Building a model, for genetic testing of paternity



35

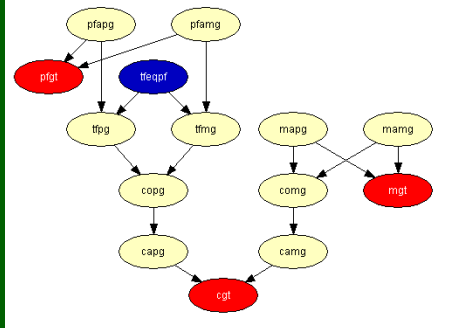
## ... using population data



we need gene frequencies relevant to assumed population for 'founder' nodes

36

## Building a model, for genetic testing of paternity

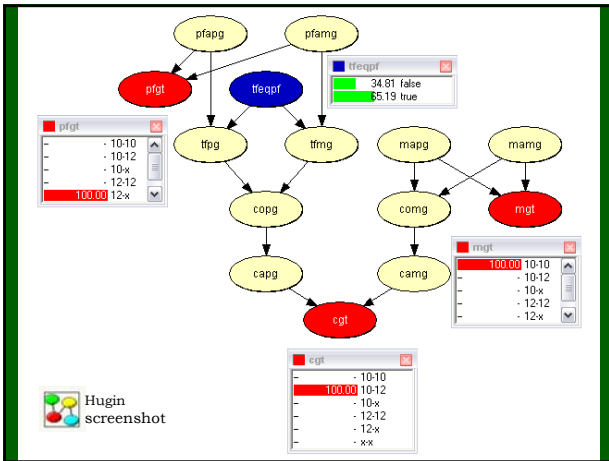


37

## Building a model, for genetic testing of paternity

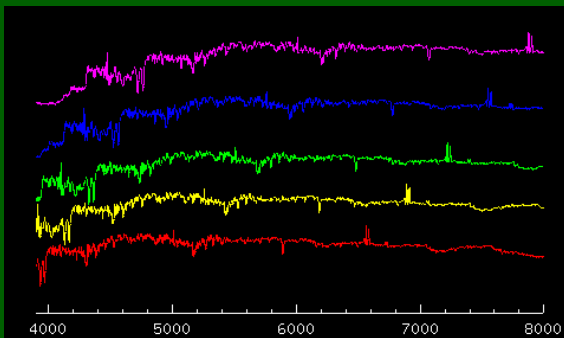
- Having established conditional probabilities within each of these local models....
- We can insert 'evidence' (data) and draw probabilistic inferences...

38



40

## Photometric redshifts



## Photometric redshifts

$$x_{ij} = a_j + b_i + M \left( \sum_{\ell=1}^4 w_{i\ell} M^{-1} \{ \psi_{\ell j}(z_i) \} \right) + \varepsilon_{ij}$$

$$\psi_{\ell j}(z) = \int h_j((1+z)u) \phi_{\ell}(u) du$$

$$M(\cdot) = -2.5 \log_{10}(\cdot)$$

magnitude =  $M(\text{flux})$

42

## Photometric redshifts

$$x_{ij} = a_j + b_i + M \left( \sum_{\ell=1}^4 w_{i\ell} M^{-1} \{ \psi_{\ell j}(z_i) \} \right) + \varepsilon_{ij}$$

Multiplicative model (on flux scale), involving an unknown mixture of templates

$$M(\cdot) = -2.5 \log_{10}(\cdot)$$

magnitude =  $M(\text{flux})$

43

## Photometric redshifts

filter response

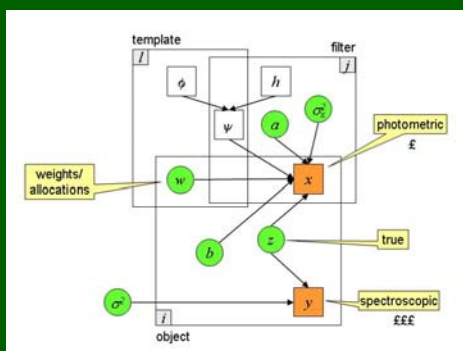
redshift

template

$$\psi_{\ell j}(z) = \int h_j((1+z)u) \phi_{\ell}(u) du$$

44

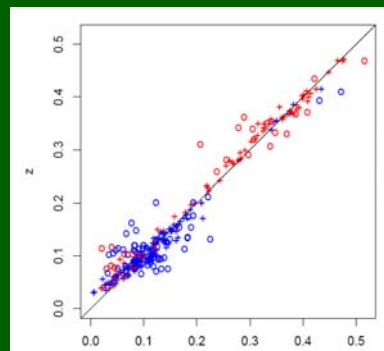
## Photometric redshifts



45

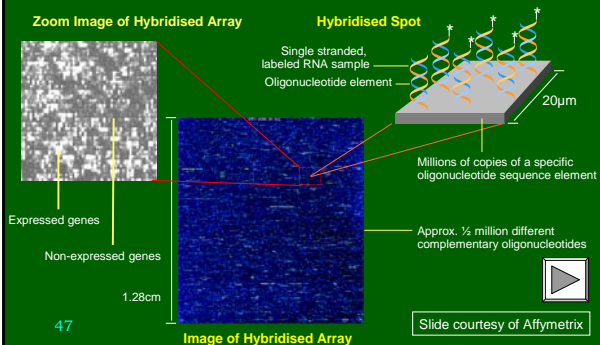
## Photometric redshifts

good agreement with 'gold-standard' spectrographic measurement



46

## Gene expression using Affymetrix microarrays



47

## Variation and uncertainty

Gene expression data (e.g. Affymetrix™) is the result of multiple sources of variability

- condition/treatment
- biological
- array manufacture
- imaging
- technical
- within/between array variation
- gene-specific variability

Structured statistical modelling allows considering all uncertainty at once

48



## Single array model: motivation

(Hein, Richardson, Causton, Ambler & G, 2004)

### Key observations:

PMs and MMs both increase with spike-in concentration (MMs slower than PMs)

Spread of PMs increase with level

Considerable variability in PM (and MM) response within a probe set

Probe effects approximately additive on log-scale

### Conclusions:

MMs bind fraction of signal

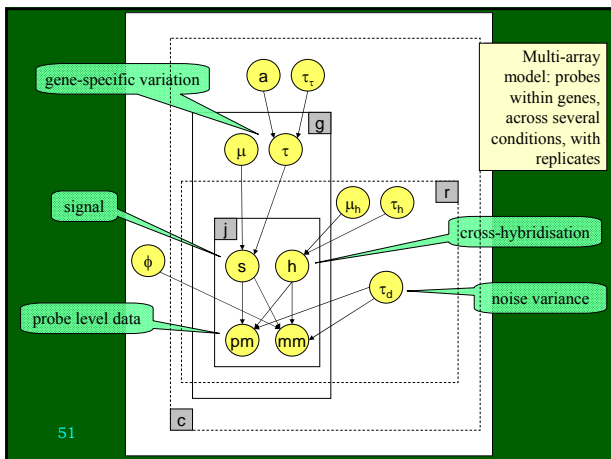
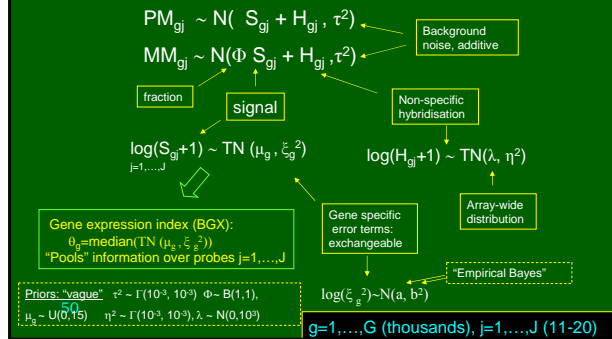
Multiplicative (and additive) error; transformation needed

Varying reliability in gene expression estimation for different genes

Estimate gene expression measure from PMs and MMs on log scale

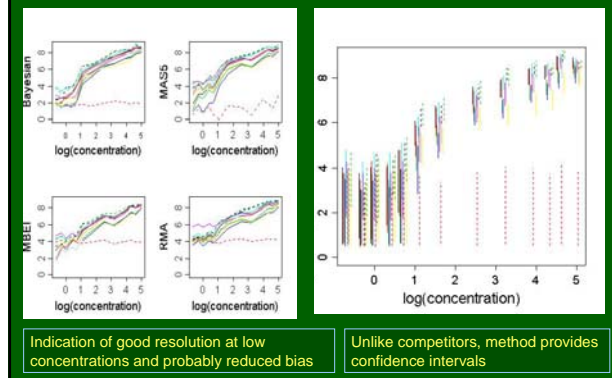
49

## BGX single array model

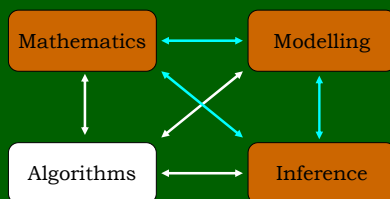


### Single array model performance: 11 genes spiked in at 13 (increasing) concentrations

Hein, Richardson, Causton, Ambler & G, 2004



## 3. Algorithms



53

## Algorithms for probability and likelihood calculations

Exploiting graphical structure:

- Markov chain Monte Carlo
- Probability propagation (Bayes nets)
- Expectation-Maximisation
- Variational (mean-field) methods

Graph representation used in user interface, data structures and in controlling computation

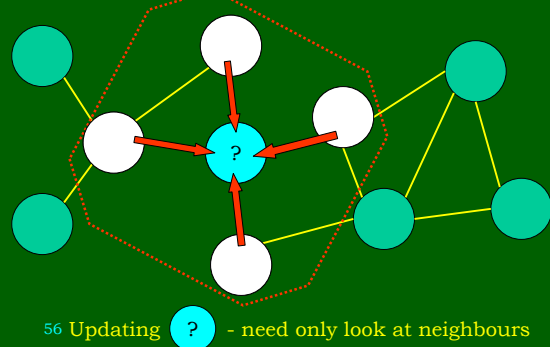
54

## Markov chain Monte Carlo

- Subgroups of one or more variables updated randomly,
  - maintaining detailed balance with respect to target distribution
- Ensemble converges to equilibrium = target distribution (= Bayesian posterior, e.g.)

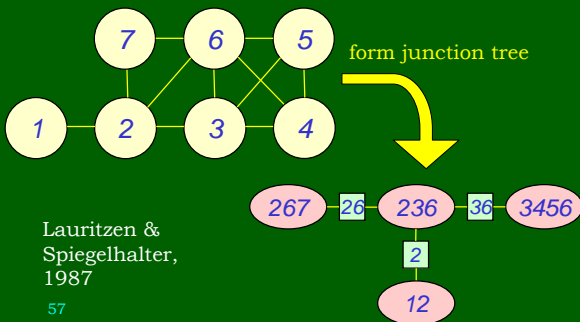
55

## Markov chain Monte Carlo



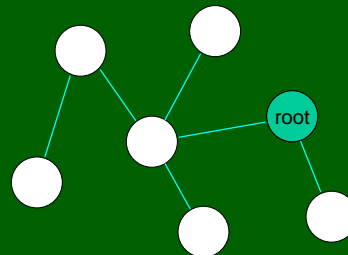
56 Updating ? - need only look at neighbours

## Probability propagation



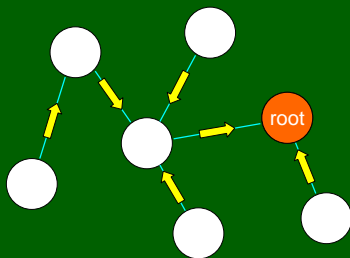
57

## Message passing in junction tree



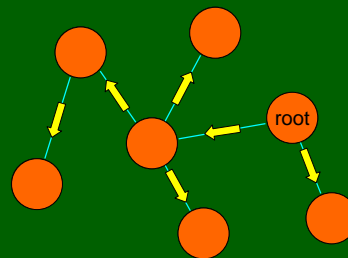
58

## Message passing in junction tree - collect



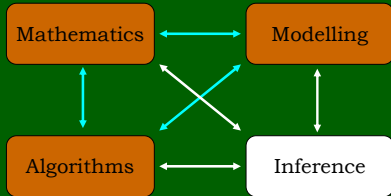
59

## Message passing in junction tree - distribute



60

## 4. Inference



61



Bayesian

62



or non-Bayesian

63

## Bayesian paradigm in structured modelling



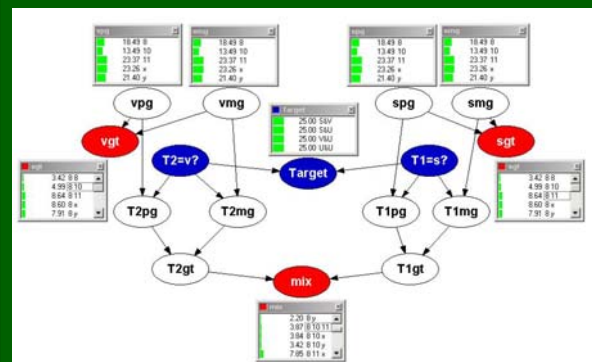
- ‘borrowing strength’
- automatically integrates out all sources of uncertainty
- properly accounting for variability at all levels
- including, in principle, uncertainty in model itself
- avoids over-optimistic claims of certainty

64

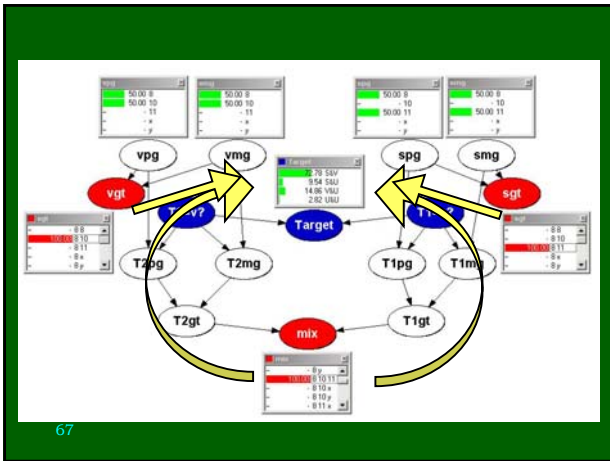
## Bayesian structured modelling

- ‘borrowing strength’
- automatically integrates out all sources of uncertainty
- ... for example in forensic statistics with DNA probe data.....

65



66



67

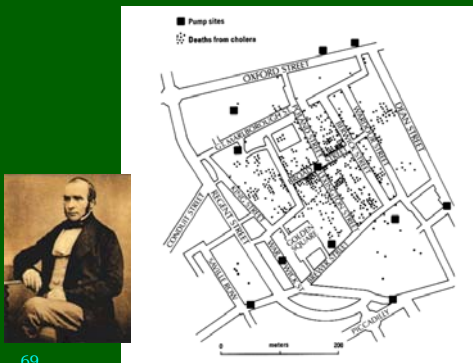
## Bayesian structured modelling

- 'borrowing strength'
- automatically integrates out all sources of uncertainty
- ... for example in hidden Markov models for disease mapping

68

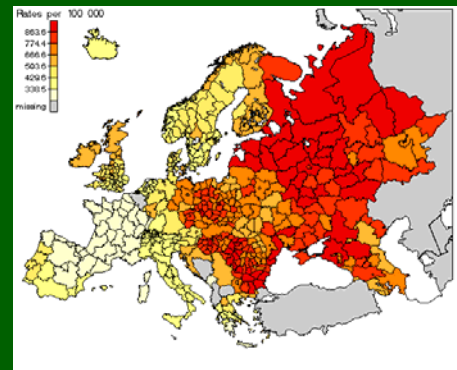


John Snow's 1855 map of cholera cases



69

Mortality for diseases of the circulatory system in males in 1990/1991



70

## Disease mapping

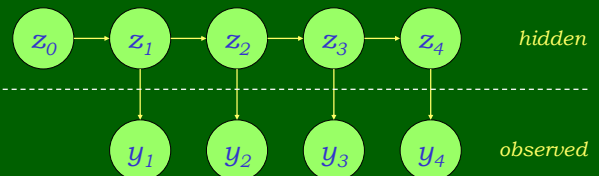
- Observe counts  $y_i$  of cases of rare, non-infectious disease in regions  $i$
- Standard model:  

$$y_i \sim \text{Poisson}(\lambda_i E_i)$$
- where  $E_i$  are adjusted populations at risk
- Relative risks  $\lambda_i$  vary due to unmeasured risk factors, assumed spatially correlated
- Use space as surrogate to separate signal and noise

71

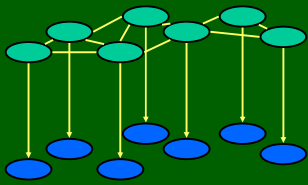
## Hidden Markov models

e.g. Hidden Markov chain



72

## Hidden Markov random fields

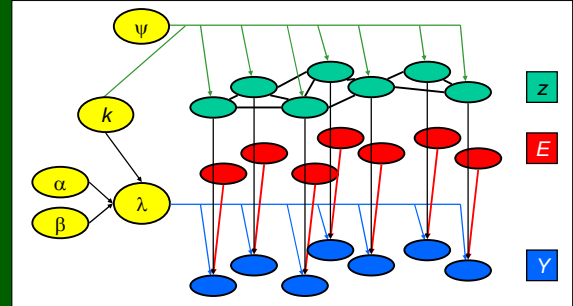


Unobserved dependent field

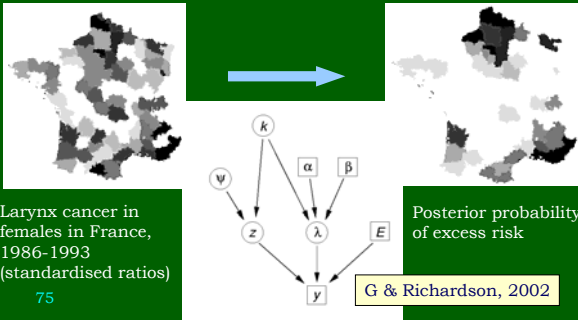
Observed conditionally-independent (e.g. discrete) field

73

## Hierarchical mixture model for disease mapping – a hidden MRF



## Mapping of rare diseases using Hidden Markov model



75

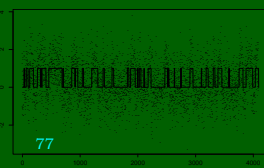
## Bayesian structured modelling

- 'borrowing strength'
- automatically integrates out all sources of uncertainty
- ... for example in modelling complex biomedical systems like ion channels.....

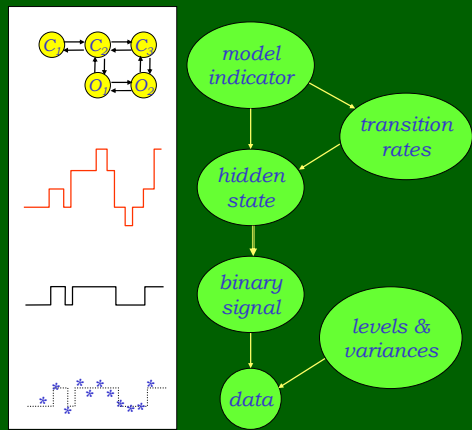
76

## Ion channel model

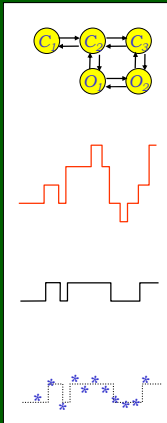
Hodgson and Green, Proc Roy Soc Lond A, 1999



77



78



Unknown physiological states of channel, unknown connections

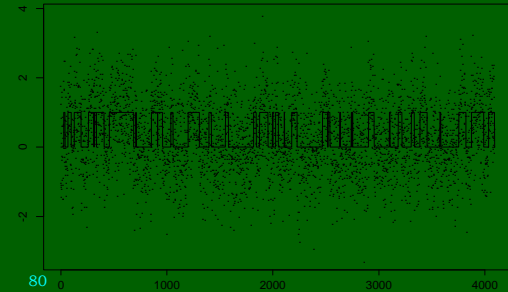
Continuous time Markov chain on this graph, with unknown transition rates

Only open/closed status of states is relevant to observation

We observe only in discrete time, with highly correlated noise

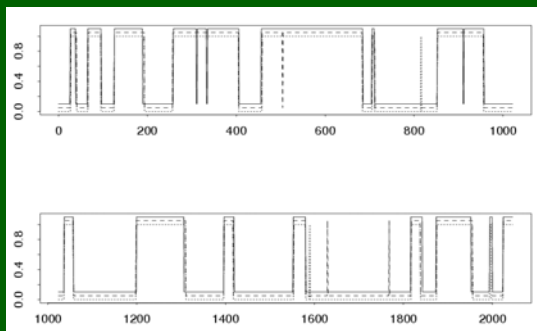
79

## Truth and simulated data



80

## Truth and 2 restorations



## Ion channel model choice

posterior probabilities

Model 11	$C_1 \xrightleftharpoons[\lambda_o]{\lambda_c} O_1$	.405
Model 21	$C_2 \xrightleftharpoons[\mu_c]{\nu_c} C_1 \xrightleftharpoons[\lambda_o]{\lambda_c} O_1$	.119
Model 12	$C_1 \xrightleftharpoons[\lambda_o]{\lambda_c} O_1 \xrightleftharpoons[\nu_o]{\mu_o} O_2$	.369
Model 22	$C_2 \xrightleftharpoons[\mu_c]{\nu_c} C_1 \xrightleftharpoons[\lambda_o]{\lambda_c} O_1 \xrightleftharpoons[\nu_o]{\mu_o} O_2$	.107

82

## Structured systems' success stories include...

- **Genomics & bioinformatics**
  - DNA & protein sequencing, gene mapping, evolutionary genetics
- **Spatial statistics**
  - image analysis, environmetrics, geographical epidemiology, ecology
- **Temporal problems**
  - longitudinal data, financial time series, signal processing

83

## Structured systems' challenges include...

- **Very large/high-dimensional data sets**
  - genomics, telecommunications, commercial data-mining...

84

## Summary

Structured stochastic modelling (the 'HSSS' approach) provides a powerful and flexible approach to the challenges of complex statistical problems

- Applicable in many domains
- Allows exploiting scientific knowledge
- Built on rigorous mathematics
- Principled inferential methods

85



<http://www.stats.bris.ac.uk/~peter>

P.J.Green@bristol.ac.uk

86