Statistics and Scientific Method

Peter J Diggle



Course text

Diggle, P.J. and Chetwynd, A.G. (2011). *Statistics and Scientific Method: an Introduction for Students and Researchers*. Oxford: Oxford University Press.

AIMS Rwanda, March 2022

イロト 不得 トイヨト イヨト 二日

Additional references

Box, G.E.P, Hunter, W.G. and Hunter, J. S. (1978). *Statistics for Experimenters: An Introduction to Design, Data Analysis and Model Building* New York: Wiley

Cox, D.R. (1972). Regression models and life tables (with Discussion). *Journal* of the Royal Statistical Society B, **34**, 187–220.

Diggle, P.J. (1990). *Time Series: a biostatistical introduction*. Oxford: Oxford University Press.

Diggle, P.J., Heagerty, P., Liang, K.Y. and Zeger, S.L. (2002). *Analysis of Longitudinal Data (second edition)*. Oxford: Oxford University Press.

Diggle, P.J. and Giorgi, E. (2019). *Model-based Geostatistics: Methods and Applications in Global Public Health*. Boca Raton: CRC Press

Diggle, P.J., Menezes, R. and Su, T.-L. (2010). Geostatistical analysis under preferential sampling (with Discussion). *Applied Statistics*, **59**, 191–232.

Eglen, S.J. and Wong, J.C.T. (2008). Spatial constraints underlying the retinal mosaics of two types of horizontal cells in cat and macaque. *Visual Neuroscience*, **25**, 209–213.

Graff-Lonnevig, V. and Browaldh, L. (1990). Twelve hours bronchodilating effect of inhaled Formoterol in children with asthma: a double-blind cross-over study versus Salbutamol. *Clinical and Experimental Allergy*, **20**, 429–432.

Grimmett, G. and Stirzaker, D. (2020). *Probability and Random Processes: Fourth Edition*. Oxford: Oxford University Press

Lee, P.M. (2012). Bayesian Inference: an Introduction (4th edition). Chichester: Wiley

McCullagh, P. and Nelder, J.A. (1989). *Generalized Linear Models* (second edition). London : Chapman and Hall.

Mercer, W.B. and Hall, A.D. (1911). The experimental error of field trials. *Journal of Agricultural Science*, **4**, 107–132.

Pawitan, Y. (2001). In All Likelihood: Statistical Modelling and Inference Using Likelihood. Oxford: Oxford University Press

Priestley, M.B. (1981). *Spectral Analysis and Time Series*. London: Academic Press.

Wakefield, J. (2007). Disease mapping and spatial regression with count data. *Biostatistics*, **8**, 158–183.

Overview

◆□▶ ◆□▶ ◆臣▶ ◆臣▶ = 臣 = のへで



・ロト・西ト・モン・ビー もくの



Isaac Newton (1643-1727)

An experiment to illustrate Newton's Law



▲□▶ ▲□▶ ▲ □▶ ▲ □▶ ▲ □ ● のへで

Experimental results

The data reproduced below are the results obtained by one student from 22 experimental runs.

t(sec)	d(cm)	t(sec)	d(cm)	t(sec)	d(cm)
0.241	10	0.358	40	0.460	70
0.249	10	0.395	45	0.485	75
0.285	15	0.435	50	0.508	80
0.291	20	0.412	50	0.516	85
0.327	25	0.451	55	0.524	90
0.329	30	0.444	60	0.545	90
0.334	30	0.461	65		
0.365	35	0.481	70		

Points for discussion

- choice of design points, d?
- inputs and outputs?
- sources of systematic variation in results?
- sources of random variation in results?

Experimental results in graphical form



- How would you describe the relationship between distance and time?
- Why did I put time on the *y*-axis of the graph, rather than on the *x*-axis?

Newton's law states that the vertical distance d travelled in time t by a body initially at rest and falling under the influence of gravity is given by the formula

$$d=\frac{1}{2}gt^2$$

◆□ > ◆□ > ◆三 > ◆三 > ・三 ・ のへで

where g is a constant.

Newtonian mechanics and experimental data



• The experimentally observed relationship is non-linear, as predicted by Newton's law.



• Transformation of the data, from distance *d* to the new variable $x = \sqrt{d}$ makes the relationship linear, also as predicted by Newton's law.



• But Newton's law does not fit the data!



• We need to add an intercept to the straight line



- What does the intercept represent?
- What does the variation of the data around the fitted line represent?

A statistical model for the experiment

Start with Newtonian mechanics,

$$d=\frac{1}{2}g\times t^2$$

Now transform to $x = \sqrt{d}$ and Y = t, and put Y on the left-hand side of the equation,

$$x = \sqrt{g/2} \times Y$$
 $Y = \beta \times x$ $(\beta = \sqrt{2/g})$

Now incorporate the effects of the experimenter's reaction-time,

 $\mathbf{Y} = \alpha + \beta \mathbf{x} + \mathbf{Z}$

◆□ > ◆□ > ◆三 > ◆三 > ・三 ・ のへで

 $\mathbf{Y} = \alpha + \beta \mathbf{x} + \mathbf{Z}$

- α represents mean reaction time
- $\beta = \sqrt{2/g}$ is the quantity of scientific interest
- Z is a random error, which varies independently between different runs of the experiment

◆□ > ◆□ > ◆三 > ◆三 > ・三 ・ のへで

• Design

What data should be collected in order to answer a scientific question as precisely as possible?

Modelling

How can the variation in the data be described mathematically, so that the description:

- is not demonstrably inconsistent with the data
- incorporates the underlying science, to the extent that this is well understood
- is a simple as possible subject to the above two constraints

Inference

Given the data and the model, what scientific conclusions can be drawn?

- **9** Graphical presentation of data is almost always useful.
- **2** Statistical models should:
 - respect the data;
 - respect the underlying science.
- **③** Transforming the data can help to achieve both goals.
- The results of a statistical analysis should always be interpreted in relation to the original scientific question.
- The fundamental role of statistical method is to enable scientists to answer their questions as precisely as possible.
- The fundamental ingredients of statistical method are: design; modelling; inference.

Uncertainty

イロン イロン イヨン イヨン

Tossing a coin



Is the coin fair?

Tossing a coin



Is the coin fair?



Is the coin fair?

Variation in experimental results can arise in two qualitatively different ways:

- systematic: a change in the experimental conditions produces a different result
- **stochastic**: replication under identical conditions produces a different result

Even a well-designed experiment can therefore lead to uncertainty in how to interpret the results, but:

- stochastic does not necessarily mean completely random
- statisticians use probability to measure uncertainty

Example: weather forecasting

How uncertain are we about the weather later today? tomorrow? next week? next year?

A meteorological time series

- maximum daily temperatures (degrees C) at Bailrigg (Lancaster) field-station, September 1995 to August 1996
- note that an unusually cold Christmas 1995 was followed by a mild period in January-February



Mathematical

A branch of pure mathematics, invented by the Russian mathematician Andrey Kolmogorov (1903 – 1987) to formalise the everyday notion of uncertainty.



Frequentist

Probability as the limit of a proportion: eg genetic variation,...

Subjectivist

Probability as personal belief: eg personal decision-making

Kolmogorov's axioms

$$W$$
: all possible outcomes of a study
 $A, B, ...$: any particular outcomes
 $P(\cdot) =$ "the probability of"

Image: P(A) ≥ 0

- $\bigcirc P(W) = 1$
- **(3)** if A, B are mutually exclusive, $P(A \cup B) = P(A) + P(B)$

Addition law

$$P(A \cup B) = P(A) + P(B) - P(A \cap B)$$

Multiplication law

$$\mathbf{P}(\boldsymbol{A} \cap \boldsymbol{B}) = \mathbf{P}(\boldsymbol{A})\mathbf{P}(\boldsymbol{B}|\boldsymbol{A}),$$

where P(B|A) is the probability of B given that A has occurred.

Bayes' Theorem

$$\mathbf{P}(\boldsymbol{A}|\boldsymbol{B}) = \frac{\mathbf{P}(\boldsymbol{B}|\boldsymbol{A})\mathbf{P}(\boldsymbol{A})}{\sum_{j}\mathbf{P}(\boldsymbol{B}|\boldsymbol{A}_{j})\mathbf{P}(\boldsymbol{A}_{j})},$$

where $A_1, A_2, ...$ represent all possible different outcomes

Grimmett and Stirzaker (2020)

◆□▶ ◆□▶ ◆三▶ ◆三▶ ・三 ・ つへぐ

A topical example: Covid testing



How likely is it that a positive test result is correct?

Lateral flow test properties:

- Sensitivity $Se \approx 0.8$
- Specificity $\mathrm{Sp} \approx 0.99$
- p = prevalence of disease

 $Se = 0.8, Sp = 0.99, p = 0.01 \Rightarrow P(D|+) = 0.447$

	Disease status				
Test	D	D			
+	$\boldsymbol{p} imes ext{Se}$	$(1-p) \times (1-\mathrm{Sp})$			
_	p * (1 - Se)	(1 - p) imes Sp			
	р	1 — р			

$$P(D|+) = \frac{P(D \cap +)}{P(+)}$$
$$= \frac{p \times Se}{p \times Se + (1-p) \times (1-Sp)}$$

- **Variation** in the results of an investigation can be of two kinds:
 - systematic variation arises as a direct result of known factors typically a change in the inputs to the investigation
 - stochastic variation arises when all known factors cannot explain the observed variation
- The mathematical theory of probability formalises and quantifies the intuitive idea of uncertainty
- Stochastic is not necessarily completely random: eg coin-tossing vs weather forecasting



◆□▶ ◆□▶ ◆臣▶ ◆臣▶ ─臣 ─ のへで

Asthma

- What is it?
- How do you treat it?
- How do you know if the treatment works?

◆□ > ◆□ > ◆臣 > ◆臣 > ―臣 - のへで

A Simple Comparative Trial

- Formoterol (F) and Salbutamol (S) are two drugs used to treat chronic asthma.
- trial was carried out to compare efficacy of F and S:
 - asthmatic children recruited;
 - Peak Expiratory Flow (PEF) measured after administration of F or S

PEF is the (primary) outcome, F and S are the treatments of interest, time is a potential confounder.

Graff-Lonnevig and Browaldh (1990). Twelve hours bronchodilating effect of inhaled Formoterol in children with asthma: a double-blind cross-over study versus Salbutamol. *Clinical and Experimental Allergy*, 20, 429–432.

Data from the Asthma Trial

```
Formoterol (sorted)
```

220 250 310 310 320 330 340 370 380 385 400 410 410

Salbutamol (sorted)

90 210 260 260 270 290 300 310 350 365 370 380 390



- one very small S result
- considerable overlap between F and S results
- comparison between F and S inconclusive?

• Difference between average PEF under F and under S is 45.4 in favour of F

◆□ > ◆□ > ◆三 > ◆三 > ・三 ・ のへで

- Hypothesise that both treatments are equally effective
 - what would the data have looked like?
 - and what do they actually look like?

Data from the Asthma Trial: continued

Child	Drug F	Drug S	Child	Drug F	Drug S
1	310	270	8	385	370
2	310	260	9	400	310
3	370	300	10	410	380
4	410	390	11	320	290
5	250	210	12	340	260
6	380	350	13	220	90
7	330	365			
A Picture Paints a Thousand Words ...



F and S results are (positively) correlated... why?

- Formoterol (F) and Salbutamol (S) are two drugs used to treat chronic asthma.
- trial was carried out to compare efficacy of F and S:
 - 13 asthmatic children recruited;
 - Peak Expiratory Flow (PEF) measured after administration of F;
 - PEF measured after administration of S;
 - one week between two measurements (ordering chosen at random)

PEF is the (primary) outcome, F and S are the treatments of interest, time is a (potential) confounder.

40 50 70 20 40 30 -35 15 90 30 30 80 130

• Difference between average PEF under F and under S is 45.4 in favour of F

◆□ > ◆□ > ◆三 > ◆三 > ・三 ・ のへで

- Hypothesise that both treatments are equally effective
 - what would the data have looked like?
 - and what do they actually look like?

Study design is important.

Decide what comparisons are of interest

• F vs S

Design the experiment to eliminate extraneous sources of variation

▲□▶ ▲□▶ ▲□▶ ▲□▶ □ ● ● ●

- pairing eliminates variation between children
- randomization eliminates (on average) the time effect

- Iooks like F beats S
- but by enough to be useful? 45.4litres/min? or $45.4\pm$?

n = 13	how much data do we have?
$\bar{d} = 45.4$	how big is the average experimental effect?
SD = 40.6	how variable is it?
$SE = SD/\sqrt{n} = 11.3$	how precisely have we estimated it?

A convenient rule of thumb:

$$ar{d} \pm 2 imes SE$$

▲□▶ ▲□▶ ▲□▶ ▲□▶ □ ● ● ●

Two Statistical Pioneers



Sir Ronald Aylmer Fisher (1890–1962): statistician and geneticist



Sir Austin Bradford Hill (1897–1991): medical statistician

A block is a group of experimental units which are thought to be relatively homogeneous, i.e. which will give relatively similar results.

Examples:

- paired PEF values
- siblings
- communities

Blocking

Whenever possible, design your experiment so that comparisons of interest can be made within the same block.

How should you choose which units receive which experimental treatment?

Randomisation:

Within any block, allocate treatments to experimental units at random;

- to avoid conscious or unconscious bias in the allocation of experimental treatments to experimental units;
- and (sometimes but more often than you might think) to ensure validity of statistical inferences

▲□▶ ▲□▶ ▲□▶ ▲□▶ □ ● ● ●

Box, Hunter and Hunter (1978)

Don't do t-tests on non-Normal data?

 $t = \bar{d}/SE = 45.4/11.2 = 4.05$ P($|t_{12}| > 4.05 = 0.0016$

Design-based sampling distribution of the *t*-statistic



Before you run an experiment:

- list the possible sources of variation in the results;
- design the experiment to eliminate extraneous sources of variation from comparisons of interest;
- if it is not possible to eliminate an extraneous source of variation, use random allocation to avoid bias

Example: in the paired experiment

- pairing eliminated extraneous variation between children;
- crossover eliminated extraneous variation between time-periods
- randomisation eliminated any possible residual bias

Designing an agricultural field trial



The research farm at Rothamsted, Hertfordshire, UK.

Agricultural field trials are conducted to compare the yields of different varieties of crop-plants under realistic conditions

Experimental units are contiguous plots of land, typically long, narrow strips within a square or rectangular field.

Design questions include:

• how to orient the strips (eg North-South or East-West)?

▲□▶ ▲□▶ ▲□▶ ▲□▶ □ ● ● ●

which treatments to apply to which strips?

Aim: compare yields of four varieties of wheat

Experimental material:

- experiment to be run on a 100 metre by 100 metre square field.
- experimental units to be 100 metre by 5 metre strips, hence 20 units in all.

Context: there is a suspected North-South fertility gradient over the field.

▲□▶ ▲□▶ ▲□▶ ▲□▶ □ ● ● ●

Two possible layouts of the experiment



イロン イロン イヨン イヨン

э.

South

1. What statistic(s) might we use to compare the yields of the four varieties of wheat?

2. Would you prefer to have the strips running North-South or East-West?

3. How would you allocate varieties, A, B, C, D say, amongst the 20 strips?

4. How might you alter your design if the precise objective was to find which varieties give the largest yields:

- on fertile ground?
- and on infertile ground?

- Careful design of an investigation can:
 - **1** simplify the analysis of the resulting data
 - **2** enable more precise conclusions to be drawn
- **②** Two fundamental principles of good statistical design are:
 - **blocking:** minimise the adverse effects of known sources of extraneous variation by comparing like-with-like
 - andomisation: minimise the adverse effects of unknown sources of extraneous variation and/or subconscious bias on the investigator's part by random allocation of experimental units to treatments
- Randomisation also allows simple questions to be answered without having to assume that the data follow a particular probability distribution.

Inference

◆□ > ◆□ > ◆臣 > ◆臣 > ―臣 - のへで

Key ideas in statistical inference

Populations and samples

- population: the set of (biological, environmental,...) material to which you hope your work is relevant
- sample: the (much smaller) set of material on which you actually base your work

Design and inference

- design: how you choose your sample
- inference: how you use results from your sample to reach conclusions about the population



random sampling from a finite population

Population size N, sample size n, each member of the population has a known probability, p_i say, of being included in the sample

- **simple** random sampling: all *p_i* equal
- stratified random sampling: all p_i equal within pre-defined sub-populations
- Quasi-random sampling from an infinite population This is the implicit assumption in most lab-based studies
- convenience sampling

Use whatever is closest to hand...not recommended

Parameter estimation

- parameter: an unknown constant whose value is scientifically interesting
- point estimate: a best guess at the true value
- interval estimate: a range of values which is, in some sense, likely to include the true value of a parameter (again, informed by your data)

Hypothesis testing

- hypothesis: a statement about a parameter
- statistical test: a way of assessing whether your data are reasonably consistent with a pre-specified hypothesis

Prediction

• a probability statement about an unobserved outcome

A point estimate is of little value without some indication of its precision

Interval estimation is a compromise between:

- the width of your interval;
- the likelihood that your interval will include the right answer.
- A confidence interval is defined as follows:
 - choose an acceptable level of confidence, say 100p% (conventionally p = 0.95, but it's your choice)
 - construct an interval so that 100*p*% of the time, it will include the true value of the parameter

Exercise: all other things being equal, if you increase *p*, will the resulting confidence interval become wider or narrower?

Statistical significance is not the same as clinical/practical significance: comparing five anti-hypertensive drugs



- which drugs give a significant reduction in blood pressure?
- which drugs give a useful reduction in blood pressure?
- which drugs need further investigation?

Probably the most widely used formula for a confidence interval is

$$ar{x} \pm 2\sqrt{s^2/n}$$

Data independently replicated under identical conditions

$$x_1, x_2, ..., x_n$$

- \bar{x} is the sample mean
- *s*² is the sample variance
- *n* is the sample size
- or replace ± 2 by a number depending on n ... see next slide

Estimating a mean (2)

Strictly, you should use $\pm c_n$, where c_n depends on n, but is approximately 2 for reasonably large n



Where does the formula for the confidence interval come from?

・ロト ・ 同ト ・ ヨト ・ ヨト

э

Estimating a mean (3)

The distribution of the sample mean changes with sample size, *n*:



(ロ) (部) (目) (日) (の)

Estimating a mean (4)

The variance, V, of the sample mean changes with sample size, n:



$$\log(V) = a - \log(n) \implies V = A \times n^{-1}$$
$$\implies SE \propto \times n^{-0.5}$$

The Central Limit Theorem

Suppose outcomes $Y_1, ..., Y_n$ are an independent random sample from any distribution with mean μ and variance σ^2 , and write \bar{Y} for the sample mean, $\bar{Y} = \sum_{i=1}^{n} Y_i$.

Theorem In the limit, as $n \to \infty$,

$$rac{ar{\mathbf{Y}}-\mu}{\sqrt{\sigma^2/n}}\sim \mathrm{N}(\mathbf{0},\mathbf{1}).$$

Or in words...sample means are approximately Normally distributed

Also called the Gaussian distribution, after the German mathematician Carl Friedrich Gauss (1777 – 1855).



A statistical model is a specification of the joint distribution of a set of random variables, $Y = (Y_1, ..., Y_n)$, indexed by a set of parameters, $\theta = (\theta_1, ..., \theta_p)$; a shorthand notation for this is $[Y; \theta]$.

The likelihood function is $L(\theta) = [y; \theta]$, where y is the observed value of Y

The maximum likelihood estimator, $\hat{\theta}$, maximises $L(\theta)$

For *n* large,
$$[\hat{\theta}] = MVN(\theta, V(\theta))$$
, where $V(\theta) = \left[-\frac{\partial^2 L(\theta)}{\partial \theta^2}\right]^{-1}$

Bayesian estimation requires you to specify a prior distribution, $[\theta]$, from which you can deduce a posterior distribution,

$$[\theta|y] = \frac{L(\theta)[\theta]}{\int L(\theta)[\theta] d\theta}$$

Pawitan (2001), Lee (2012)

Conclusions

- sample means are approximately Normally distributed (symmetric, bell-shaped histogram)
- larger samples lead to more precise estimates
- the variance of an estimate is inversely proportional to the sample size, *n*
- the standard error of an estimate is therefore inversely proportional to \sqrt{n}
- law of diminishing returns doubling the sample size does not double the precision of your estimate
- likelihood-based methods provide a generally applicable, efficient and principled approach to model-based inference

- How you choose your sample is important
- Why statistical significance is not the same thing as clinical/practical significance
- How to estimate a mean from an independently replicated sample
- The \sqrt{n} law of statistical precision
- The Central Limit Theorem...sample means are approximately Normally distributed

▲□▶ ▲□▶ ▲□▶ ▲□▶ □ ● ● ●

Modelling

イロン イロン イヨン イヨン

• models are devices to answer questions, and should:

- be not demonstrably inconsistent with the data;
- incorporate the underlying science, where this is well understood
- be as simple as possible, within the above constraints

"Too many notes, Mozart"

Emperor Joseph II

"Only as many as there needed to be"

Mozart (apochryphal?)

Empirical or mechanistic approach to model-building?



• data are pairs of values x and y (input and output)

0



data are pairs of values x and y (input and output)the best-fitting line?

∃ ► < ∃ ►</p>



data are pairs of values x and y (input and output)

 some lines are "obviously" better fits to the data than others, but which is the "'best-fitting" line?



- data are pairs of values x and y (input and output)
- dashed vertical lines are residuals
- the best-fitting line makes the sum of squared residuals as small as possible

くぼう くほう くほう

Exercise: why are residuals measured vertically?
Glyphosate data

- Glyphosate is a powerful weed-killer, its presence in the water-supply is potentially harmful to irrigated crops.
- Experiment conducted to investigate how average root-length of batches of 15 safflower plants is affected by glyphosate added to distilled or tap water

<i>x</i> (ppm)	0.000	0.000	0.053
y (distilled)	107.0	110.9	106.2
y (tap)	111.0	168.3	105.7
x (ppm)	0.106	0.211	0.423
y (distilled)	97.3	105.9	88.5
y (tap)	116.7	143.7	84.7
v (ppm)	0 845	1 600	3 380
× (ppiii)	0.045	1.005	3.300
y (distilled)	14.4	46.2	30.0
y (tap)	59.3	36.7	38.0

Scientific and statistical objectives

- what can these data tell us about the effect of small concentrations of glyphosate on plant growth?
- how could we build a statistical model to describe the relationship between glyphosate concentration and root-length?

Where do models come from?

- in the gravity experiment, the linear regression model had a mechanistic justification (physics plus physiology)
- In the glyphosate experiment, there is no scientific law to guide us.
- but we may still be able to use the linear regression model to describe the empirical relationship between glyphosate concentration and root-length.

Plotting the glyphosate data



◆ロ ▶ ◆□ ▶ ◆ □ ▶ ◆ □ ▶ ◆ □ ▶ ◆ □ ▶

Notes on data-transformations

- other transformations of the data could have been used
- choice could be determined by empirical and/or scientific considerations

Example.

Suppose x and y follow a power law model,

$$y = ax^b$$

Then, log-log transformation produces a linear model,

$$\mathbf{Y} = \alpha + \beta \mathbf{X}$$

where $Y = \log y$, $X = \log x$, $\alpha = \log a$ and $\beta = b$.

A linear model for the glyphosate data

$$y = \log(\text{root length})$$
 (response)

 $x = \log(1 + \text{glyphosate})$ (explanatory variable)

w = 0/1 = distilled/tap water (factor)

$$\mathbf{y} = \{\alpha_0 + \alpha_1 \mathbf{w}\} + \beta \mathbf{x} + \mathbf{z}$$

- parallel straight-line relationships for distilled and for tap-water
- α₁ = 0 if source of water does not affect average root-length
- β measures effect of glyphosate on plant-growth (on transformed scale)

Residuals vs fitted values plot



fitted values

▲□▶ ▲□▶ ▲ □▶ ▲ □▶ ▲ □ ● のへで

Anscombe's quartet



2

Residuals: Anscombe's quartet



DATA = FITTED VALUE + RESIDUAL

Anscombe's quartet	R code for linear regression			
All four data-sets have:	Call: lm(formula = y	y1 x123)		
 the same best-fitting straight line, 	Coefficients:	Estimate	SE	t value
$\mathbf{y} = \hat{\alpha} + \hat{\beta}\mathbf{x}$	(Intercept)	3.0001	1.1247	2.667
• the same standard	x123	0.5001	0.1179	4.241
errors for \hat{lpha} and \hat{eta}	Residual standard error: 1.237 on 9			
 the same residual sum of squares 	degrees of fre	eedom		

But the residuals themselves tell four different stories.

	res1	res2	res3	res4
1	-0.740	-1.901	0.389	0.000
2	0.179	-0.761	0.229	-0.111
3	1.239	0.129	0.079	-1.751
4	-1.681	0.759	-0.081	0.909
5	-0.051	1.139	-0.230	-1.241
6	1.309	1.269	-0.390	1.839
7	0.039	1.139	-0.540	-0.421
8	-0.171	0.759	-0.689	1.469
9	1.839	0.129	-0.849	-1.441
10	-1.921	-0.761	3.241	0.709
11	-0.041	-1.901	-1.159	0.039

Why are the four sets of residuals different?

- obvious if you have only one explanatory variable
- less obvious when you have many

Analysing residuals

- check that their average value is (close to) zero
- plot them against fitted values
- plot them against explanatory variables in the model
- plot them against explanatory variables **not** in the model (for example, residuals against time-order)

▲□▶ ▲□▶ ▲□▶ ▲□▶ □ ● ● ●

What if our response is a proportion?

- a linear regression may fit over a restricted range of the input variable
- but it usually fails when the observed proportions cover most of the range from zero to one.



input

What if our response is a proportion?

- a linear regression may fit over a restricted range of the input variable
- but it usually fails when the observed proportions cover most of the range from zero to one.



input

What if our response is a proportion?

- a linear regression may fit over a restricted range of the input variable
- but it usually fails when the observed proportions cover most of the range from zero to one.



input

 $\mathbf{Y} = \alpha + \beta \mathbf{x} + \mathbf{Z}$

Or equivalently:

- $\mu = \mathbf{E}[\mathbf{Y}] = \alpha + \beta \mathbf{x}$
- $\mathbf{Y} \sim \mathrm{N}(\mu, \sigma^2)$

The generalisation:

- $h(\mu) = \alpha + \beta x$ (link function)
- $\mathbf{Y} \sim f(\mu, ...)$ (error distribution)

- $h(\mu) = \alpha + \beta x$ (Link function)
- $Y \sim f(\mu, ...)$ (error distribution)
- choice of link function makes linear dependence on explanatory variables less restrictive
- choice of error distribution adds flexibility
- puts a very wide range of statistical methods under a common framework
- encourages open thinking (problem-driven rather than recipe-driven)

McCullagh and Nelder (1989)

Linear model

$$\mu_i = \mathbf{E}[\mathbf{Y}_i] = \alpha + \beta \mathbf{x}_i \qquad \text{Var}(\mathbf{Y}) = \sigma^2$$

Generalised linear model

$$\mu_i = \mathbf{E}[\mathbf{Y}_i] = h^{-1}(\alpha + \beta x_i) \qquad \operatorname{Var}(\mathbf{Y}) = \mathbf{v}(\mu_i)$$

Residuals and standardised residuals

	Linear	Generalised linear
residual	$y_i - \hat{\mu}_i$	$\mathbf{y}_i - \hat{\mu}_i$
standardised residual	$(y_i - \hat{\mu}_i)/\hat{\sigma}$	$(\mathbf{y}_i - \hat{\mu}_i)/\sqrt{\mathbf{v}(\hat{\mu}_i)}$

The statistical modelling cycle

- identify the question
- design the study
- collect the data
- build the model
 - exploration
 - Ø fitting
 - o diagnostic checking
 - repeat (a) to (c) as necessary

◆□▶ ◆□▶ ◆三▶ ◆三▶ 三 シのので

answer the question

Open-ended count response

Default choice is Poisson log-linear model

$$\log(\mu) = \alpha + \beta x$$
 $Y \sim \text{Poiss}(\mu)$

• Data-sets often show evidence of extra-Poisson variation

$$\log(\mu) = \alpha + \beta x$$
 $Var(Y) = \phi \mu : \phi > 1$

◆□ > ◆□ > ◆三 > ◆三 > ・三 ・ のへで

Examples of generalised linear models (2)

Binary or closed-count response (number out of *n*)

• (Binomial) logistic model

$$\log\{\mu/(1-\mu)\} = lpha + eta x \qquad \mu = P(Y=1)$$

• Complementary log-log

$$\log\{-\log(\mu)\} = \alpha + \beta x \qquad \mu = P(Y = 1)$$



• Extra-binomial variation can also arise, and can be handled in the same way as extra-Poisson variation

Examples of generalised linear models (3)

Survival analysis (life-times)

• Exponential log-linear model

$$\log(\mu) = \alpha + \beta x$$
 $\mu = P(Y = 1)$

$$f(\mathbf{y}) = \mu^{-1} \exp(-\mathbf{y}/\mu)$$

Rarely used in practice, because of the problem of censoring

 Standard approach is semi-parametric proportional hazards modelling

Cox (1972)

- statististical models are devices to answer questions
- statististical models can be empirical or mechanistic (or a mix of the two)
- **(3)** the linear regression model is surprisingly flexible
- but for count data, generalised linear models are usually preferable
- residual diagnostic checks on model fit are an important part of the statistical modelling cycle

◆□ > ◆□ > ◆三 > ◆三 > ・三 ・ のへで

Random effects

◆□ > ◆□ > ◆臣 > ◆臣 > ―臣 - のへで

• input variable x

factor, covariate, explanatory variable, ...

• output variable y

response, end-point, primary outcome,...

◆□ > ◆□ > ◆臣 > ◆臣 > ―臣 - のへで

A synthetic example



- relationship between x and y can be captured approximately by a straight line
- scatter about the line is approximately the same at all values of *x*

Interpreting the linear regression model



 $Y_i = \alpha + \beta x_i + z_i$

・ 同 ト ・ ヨ ト ・ ヨ ト

How precisely can we:

- estimate the parameters α and β ?
- estimate the straight-line relationship?
- predict a future value of y?

The data are longitudinal



- simple linear regression software assumes that data are uncorrelated
- in longitudinal studies, with repeated measurements on each subject, this is rarely true
- reported standard errors and p-values are then not correct

Diggle, Heagerty, Liang and Zeger (2002)

Standard fitting method: least squares

- > fit1<-lm(y~1+x)</pre>
- > summary(fit1)

... plus lots of stuff you don't want to know

Coefficients:

Estimate Std. Error t value Pr(>|t|) (Intercept) 3.26789 0.10274 31.81 <2e-16 *** x 0.39286 0.01924 20.41 <2e-16 ***

Residual standard error: 0.7817 on 198 degrees of freedom Multiple R-squared: 0.6779, Adjusted R-squared: 0.6763 F-statistic: 416.7 on 1 and 198 DF, p-value: < 2.2e-16

Correct fitting method: maximum likelihood with random effect

```
> library(nlme)
> fit2<-lme(y~1+x,random=~1|id)</pre>
> summary(fit2)
Linear mixed-effects model fit by REML
. . .
Random effects:
 Formula: ~1 | id
        (Intercept) Residual
StdDev: 0.7477531 0.2730349
Fixed effects: y ~ 1+x
               Value Std.Error DF t-value p-value
(Intercept) 3.267887 0.17100989 179 19.10935
                                                    0
            0.392856 0.00672165 179 58.44637
x
                                                    0
Number of Observations: 200
Number of Groups: 20
                                        ◆□▶ ◆□▶ ◆三▶ ◆三▶ 三 シのので
```

- Random effects can be thought of as missing information on individual subjects that, were it available, would be included in the statistical model
- we model the missing information as a random sample from a distribution (usually, we assume a Normal distribution)
- This induces correlation amongst repeated measurements on the same subjects

▲□▶ ▲□▶ ▲□▶ ▲□▶ □ ● ● ●

Example: some subjects are intrinsically high responders, others intrinsically low responders

 \Rightarrow replace fixed intercept α by a random intercept, $\alpha + U_i$

For our synthetic example, write:

$$Y_{ij} = j^{th}$$
 response from i^{th} subject: $i = 1, ..., n$

 x_{ij} = corresponding value of explanatory variable

A random effects model

•
$$Y_{ij} = \alpha + \beta x_{ij} + U_i + Z_{ij}$$

- U_i = random effect for subject *i*
- Z_{ij} = residual
- all U_i and all Z_{ij} mutually independent

Model implies that different responses on the same subject are positively correlated:

$$\rho = \frac{\operatorname{Var}(U)}{\operatorname{Var}(U) + \operatorname{Var}(Z)}$$

For our synthetic example, write:

$$Y_{ij} = j^{th}$$
 response from i^{th} subject: $i = 1, ..., n$
 $x_{ij} =$ corresponding value of explanatory variable

A fixed effects model

•
$$\mathbf{Y}_{ij} = \alpha_i + \beta \mathbf{x}_{ij} + \mathbf{Z}_{ij}$$

• α_i = intercept for subject *i*

• all Z_{ij} mutually independent

Model implies that different responses on the same subject are uncorrelated

▲□▶ ▲□▶ ▲□▶ ▲□▶ □ ● ● ●

• They both are

- The choice between them depends primarily on why you are analysing the data
 - to adjust for the heterogentiy amonst the particular subjects in your data
 - to undertand the heterogeneity amongst members of the population from which your subjects were drawn
- These two cases call for fixed effects and random effects models, respectively

• A second consideration is statistical efficiency: for large *n*, the fixed effects model has many more parameters than the random effects model, which necessarily implies that the parameter estimates are less precise

 Y_{ij} = mark for student *i* on exam paper j = 1, ..., pQuestion: what overall mark should you give to student *i*?

Fixed effects model: $Y_{ij} = \alpha_i + Z_{ij}$ $Z_{ij} \sim N(0, \tau^2)$, independent Answer: $\hat{\alpha}_i = \bar{Y}_i$, the observed average mark for student *i*

Random effects model: $Y_{ij} = A_i + Z_{ij}$

• $Z_{ij} \sim N(0, \tau^2)$, independent • $A_i \sim N(\alpha, \sigma^2)$, independent

Answer: $\hat{A}_i = c \times \bar{y}_i + (1 - c) \times \bar{y}, \quad c = p/(p + \tau^2/\sigma^2)$ Observed average for student *i* is "shrunk" towards the class average

An RCT of drug therapies for schizophrenia

- Randomised clinical trial of drug therapies
- Three treatments:

haloperidol (standard); placebo; risperidone (novel)

• Dropout due to "inadequate response to treatment"

Treatment	Number of non-dropouts at week					
	0	1	2	4	6	8
haloperidol	85	83	74	64	46	41
placebo	88	86	70	56	40	29
risperidone	345	340	307	276	229	199
total	518	509	451	396	315	269

The schizophrenia trial data



time (weeks since randomisation)

э < 6 b
A summary of the schizophrenia trial data



< □ > < □ > < 三 > < 三 > < 三 > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □

Mean response depends on treatment and time

Two random effects:

- between subjects (high or low responders)
- between times within subjects (good and bad days)

Method of analysis allows for dropouts: maximum likelihood with random effects

▲□▶ ▲□▶ ▲□▶ ▲□▶ □ ● ● ●

PANSS mean response profiles



What's going on?

- dropout is selective (high responders more likely to leave)
- but the data are correlated
- and this allows the model to infer what you would have seen, had there not been any dropouts
- which may or may not be what you want



Diggle, Heagerty, Liang and Zeger (2002)

< ロ > < 同 > < 回 > < 回 >

- Fixed effects describe the variation in average responses of groups of subjects according to their measured characteristics (age, sex, treatment,...)
- Pandom effects describe variation in subject-specific responses according to their unmeasured characteristics
- Both kinds of model can easily be fitted using the open-source software R, or in various proprietary packages
- Oropout in longitudinal studies can have surprising consequences:
- Sandom effects and parameters are different things
 - parameters don't change if you re-run an experiment
 - random effects do

Time

▲□ > ▲圖 > ▲ 臣 > ▲臣 > □ 臣 = の Q @

A meteorological time series

- maximum daily temperatures (degrees C) at Bailrigg (Lancaster) field-station, September 1995 to August 1996
- note that an unusually cold Christmas 1995 was followed by a mild period in January-February



- what are the main features of the data?
- how did I fit the smooth curve to the data?
- what features are and are not explained by the fitted curve?

$$Y(t) = \mu + \alpha \cos(2\pi t/p + \phi) + residual$$

$$= \mu + \beta_1 \cos(2\pi t/p) + \beta_2 \sin(2\pi t/p) + \text{residual}$$

- μ = overall mean value (of time series Y(t))
- p = period
- $\alpha = \text{amplitude}$
- $\phi = \text{phase}$

Usually, the period is known, but the mean, amplitude and phase are not

Use the first form of the model,

$$Y(t) = \mu + lpha \cos(2\pi t/p + \phi) + ext{residual}$$



Use the first form of the model,

$$Y(t) = \mu + lpha \cos(2\pi t/p + \phi) + ext{residual}$$



Use the first form of the model,

$$Y(t) = \mu + lpha \cos(2\pi t/p + \phi) + ext{residual}$$



Use the first form of the model,

$$Y(t) = \mu + lpha \cos(2\pi t/p + \phi) + ext{residual}$$



Use the first form of the model,

$$Y(t) = \mu + lpha \cos(2\pi t/p + \phi) + ext{residual}$$



Use the first form of the model,

$$Y(t) = \mu + lpha \cos(2\pi t/p + \phi) + ext{residual}$$



Use the first form of the model,

$$Y(t) = \mu + lpha \cos(2\pi t/p + \phi) + ext{residual}$$



Use the first form of the model,

$$Y(t) = \mu + lpha \cos(2\pi t/p + \phi) + ext{residual}$$



Use the first form of the model,

$$Y(t) = \mu + lpha \cos(2\pi t/p + \phi) + ext{residual}$$



Use the first form of the model,

$$Y(t) = \mu + lpha \cos(2\pi t/p + \phi) + ext{residual}$$



$$Y(t) = \mu + \alpha \cos(2\pi t/p + \phi) + \text{residual}$$

Lifting

 μ : adjust to match observed and modelled average value

◆□ > ◆□ > ◆臣 > ◆臣 > ─臣 ─のへで

Stretching

 α : adjust to match observed and modelled range

Shifting

 ϕ : adjust to match observed and modelled nadir

Fitting the model

Use the second form of the model,

$$Y(t) = \mu + \beta_1 \cos(2\pi t/p) + \beta_2 \sin(2\pi t/p) + \text{residual}$$

Note that the following quantities are known, i.e. they can be calculated without having to estimate anything

- $x_1(t) = \cos(2\pi t/p)$
- $x_2(t) = \sin(2\pi t/p)$

Re-write the model as

$$\mathbf{Y} = \mu + \beta_1 \mathbf{x}_1 + \beta_2 \mathbf{x}_2$$

After fitting (see next page), amplitude and phase can be recovered using

$$\alpha = \sqrt{\beta_1^2 + \beta_2^2} \qquad \phi = \tan^{-1}(\beta_2/\beta_1)$$

▲□▶ ▲□▶ ▲□▶ ▲□▶ □ ● ● ●

Using the lm() function to fit the model

```
data<-read.table("maxtemp.dat")
y<-data[,4]
day<-1:366
x1<-cos(2*pi*day/366)
x2<-sin(2*pi*day/366)
fit<-lm(y~x1+x2)
summary(fit)</pre>
```

Call: lm(formula = y ~ x1 + x2)Residuals: Min 10 Median 30 Max -7.5921 -1.8240 -0.1475 1.7140 8.5232 Coefficients: Estimate Std. Error t value Pr(>|t|) (Intercept) 11.8467 0.1441 82.22 <2e-16 *** x1 6.2508 0.2038 30.68 <2e-16 *** -3.3177 0.2038 -16.28 <2e-16 *** x2 Signif. codes: 0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1

Residual standard error: 2.756 on 363 degrees of freedom Multiple R-Squared: 0.7687, Adjusted R-squared: 0.7674 F-statistic: 603.1 on 2 and 363 DF, p-value: < 2.2e-16

- relationship between today's and yesterday's temperature?
- relationship between today's and yesterday's residual?



◆□ ▶ ◆□ ▶ ◆ 臣 ▶ ◆ 臣 ▶ ○ 臣 ○ のへで

• how and why are the two relationships different?

Autocorrelation (2)

How does the relationship between residuals today and k days ago change as k increases?







residuals



residuals



Autocorrelation (3)

- lag-k autocorrelation is the correlation between pairs of values from the same time series k time-units apart
- correlogram is a plot of lag-k autocorrelation against k



- dashed lines at $\pm 2/\sqrt{n}$ are pointwise 95% limits for uncorrelated residuals
- but overall pattern is more important than individual numerical values

Imagine that you have data on daily maximum temperatures for several years, up to today.

- 1. How would make a forecast of:
 - tomorrow's temperature?
 - the temperature one month from now?
- 2. In what ways are your two answers different, and why?

 $Y(t) = \alpha + \beta_1 \cos(2\pi t/p) + \beta_2 \sin(2\pi t/p) + \text{residual}(t)$

1. Model consists of a time-varying mean, also called the trend,

$$\mu(t) = \alpha + \beta_1 \cos(2\pi t/p) + \beta_2 \sin(2\pi t/p)$$

and stochastic variation, residual(t), about the trend.

2. Decompose the residual into two terms:

$$residual(t) = S(t) + Z_t$$

• $\operatorname{Cov}{S(t), S(t-u)} = \sigma^2 \exp(-u/\phi)$ (random effect)

• Z_t uncorrelated N(0, τ^2) (noise/measurement error) Diggle (1990), Priestley (1981)





Spatial Data Formats

spatially aggregated; spatially sampled; point process

Example 1a. Wheat uniformity trial



Mercer and Hall (1911)

Example 1b. Cancer atlases

Raw (left panel) and spatially smoothed (right panel) relative risk estimates for

lip cancer in 56 Scottish counties



Wakefield (2007)

Example 2. Galicia biomonitoring study

Lead concentrations measured in samples of moss, map shows locations and log-concentrations



Diggle, Menezes and Su (2010)

Example 3. Retinal mosaics

Locations of two types of light-responsive cells in macaque retina (2 animals)





Eglen and Wong (2008)

◆□ > ◆□ > ◆臣 > ◆臣 > ―臣 - のへで

First law of geography: all things are related but close things are more strongly related than distant things

Spatial dependence may be:

- functional (deterministic)
- statistical (stochastic)

Origin of statistical dependence may be mechanistic (causal) or empirical (descriptive/predictive)

▲□▶ ▲□▶ ▲□▶ ▲□▶ □ ● ● ●

A common scenario in tropical disease epidemiology:

- Objective: understand geographical variation in risk
- Study-design: Identify *n* communities in study-region *A*, by a suitable sampling scheme (not necessarily completely random)
- Data: Measure empirical prevalence and hypothesised risk-factors on each sampled community, also risk-factors on all unsampled communities for which estimates of risk are required
- Model: either a linear or a generalized linear model for prevalence (or transformation thereof)
- Diagnostic checking: all the usual plus check for spatial independence of residuals

Diggle and Giorgi (2019)

A model-based approach: random effects again

$Y_i = \mu(x_i) + S(x_i) + Z_i : i = 1, ..., n$

▲□▶ ▲□▶ ▲□▶ ▲□▶ □ ● ● ●

- $x_i = \text{location}$
- $Y_i = (transformed)$ prevalence
- $\mu(x_i) = u'_i \beta$ (fixed effects)
- S(x) = spatially correlated process, variance σ^2
- Z_i = uncorrelated residuals, variance τ^2
Estimating correlation structure

- Regression residuals, $r_i = Y_i u_i \hat{\beta}$
- Each r_i estimates $S(x_i) + Z_i$
- Do we need the random effect process S(x)?

The variogram

•
$$d_{ij} = ||x_i - x_j||$$
(distance) $v_{ij} = \frac{1}{2}(r_i - r_j)^2$

•
$$V(d_{ij}) = \mathbb{E}[v_{ij}^2]$$

•
$$V(d) = \tau^2 + \sigma^2 \{1 - \rho(d)\}$$

"A small black fly, with thick shoulders and bullet-head, infests the place, and torments the naked arms and legs of the people with its sharp stings to an extent that must make life miserable to them" John Hanning Speke, Uganda, 1864



The first written description of the onchocerciasis vector, *Simulium damnosum* (blackfly)

Onchocerciasis: aka River Blindness



MDA: a tool for control of vector-borne filarial disease





- Ivermectin (Mectizan): annual dose clears microfilarial infections of the blood
- generally considered safe, with no serious side-effects
- mass distribution made possible by donation programme (Merck)
- used in multi-national programmes to combat onchocerciasis and lymphatic flariasis

A prevalence survey data-set: onchocerciasis in Liberia



- prevalence data from 90 villages in Liberia
- sample sizes 40 to 50
- empirical prevalences 0% to 35%

くぼう くほう くほう

Data: $(x_i, d(x_i), n_i, Y_i) : i = 1, ..., n$

Model:

- $\log[p(x_i)/\{1-p(x)_i\}] = d(x)_i'\beta + U_i + S(x_i)$
- $U_i \sim N(0, \nu^2)$, mutually independent
- $S(x) \sim$ Gaussian process, spatially correlated

・ロト ・ 日 ト ・ ヨ ト ・ ヨ ・ つへで

• $Y_i | U_i, S(x_i) \sim \text{Binomial}\{n_i, p(x_i)\}$

Exploratory analysis: empirical logits

- fitting the binomial logistic model is computationally demanding, and requires judgement:
 - convergence of iterative algorithms
 - judicious choice of approximations
- empirical logit transform:

 $Z_i = \log\{(Y_i + 0.5)/(n_i - Y_i + 0.5)\}$

▲□▶ ▲□▶ ▲□▶ ▲□▶ □ ● ● ●

• fit linear model with Z_i as response

Exploratory analysis of onchocerciasis data



• residual variogram after fitting linear trend surface

Exploratory analysis of onchocerciasis data



• fitted Matérn model with $\kappa = 0.5$

- likelihood function involves intractable high-dimensional integral
- need to use Monte Carlo methods
- Monte Carlo maximum likelihood or Bayesian estimation according to choice
- for large data-sets, algorithms need careful tuning to preserve accuracy while remaining computationally feasible

```
library(splancs)
xy<-data.frame(longitude=data$longitude,latitude=data$latitude)
par(pty="s"); pointmap(xy)
poly<-getpoly()</pre>
grid.predict<-as.data.frame(gridpts(poly,xs=0.1,ys=0.1))</pre>
names(grid.predict)<-c("longitude","latitude")</pre>
predict.MCML<spatial.pred.binomial.MCML(fit.bl,</pre>
       grid.pred=grid.predict,predictors=grid.predict,
       control.mcmc=mcmc,scale.predictions="prevalence",
       standard.errors=TRUE,thresholds=0.2,
       scale.thresholds="prevalence")
plot(predict.MCML,type="prevalence")
points(gd,add=TRUE)
polymap(poly,add=TRUE)
```

Liberia: onchocerciasis prevalence map



向下 イヨト イヨト э

Liberia: onchocerciasis exceedance maps

P(prevalence>10%)

P(prevalence>20%)



Summary

- Statistics is fundamentally about understanding variation
- Design: eliminate extraneous sources of variation
- Model: acknowledge remaining sources of variation
- Inference: estimate model parameters efficiently so that interpretations are:
 - valid (honest)
 - efficient (as precise as possible)
- Design and modelling both involve scientific judgement
- Inference should follow automatically (likelihood-based)
- And the inference should answer the original question

Statistics and scientific method



Statistical method :

- a device to answer a question
- a bridge between theoretical and applied science
- a framework to enable principled inference in the presence of uncertainty
- Scientific purpose is more important than data-format
- Analyse problems, not data

Diggle, P.J. and Chetwynd, A.G. (2011). *Statistics and Scientific Method: an Introduction for Students and Researchers.* **Oxford: Oxford University Press.**