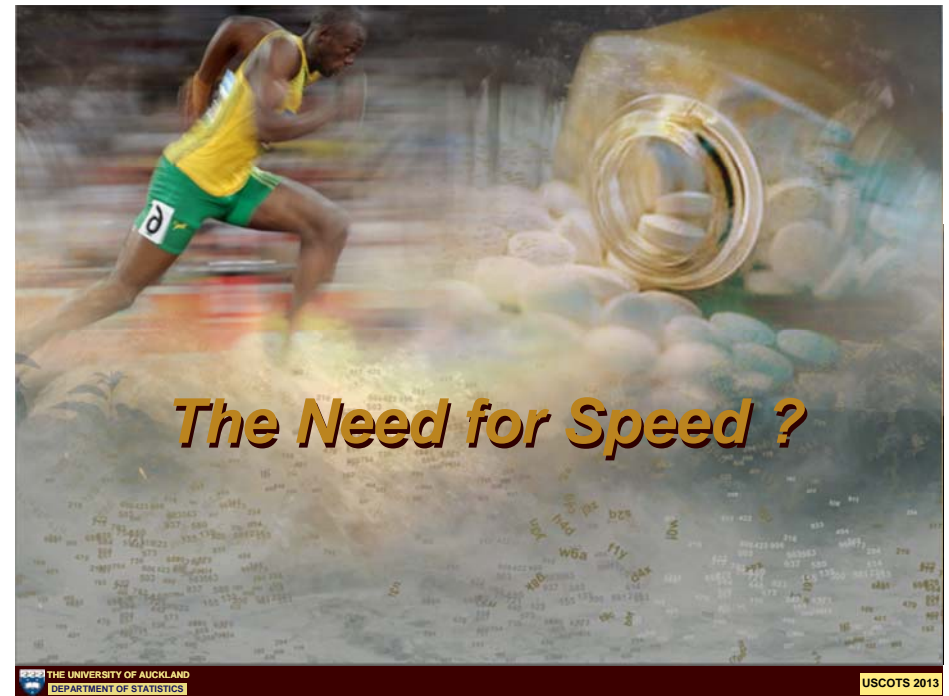# The Need for Speed
## in the Path of the Deluge

### Chris Wild

Department of Statistics
University of Auckland, New Zealand

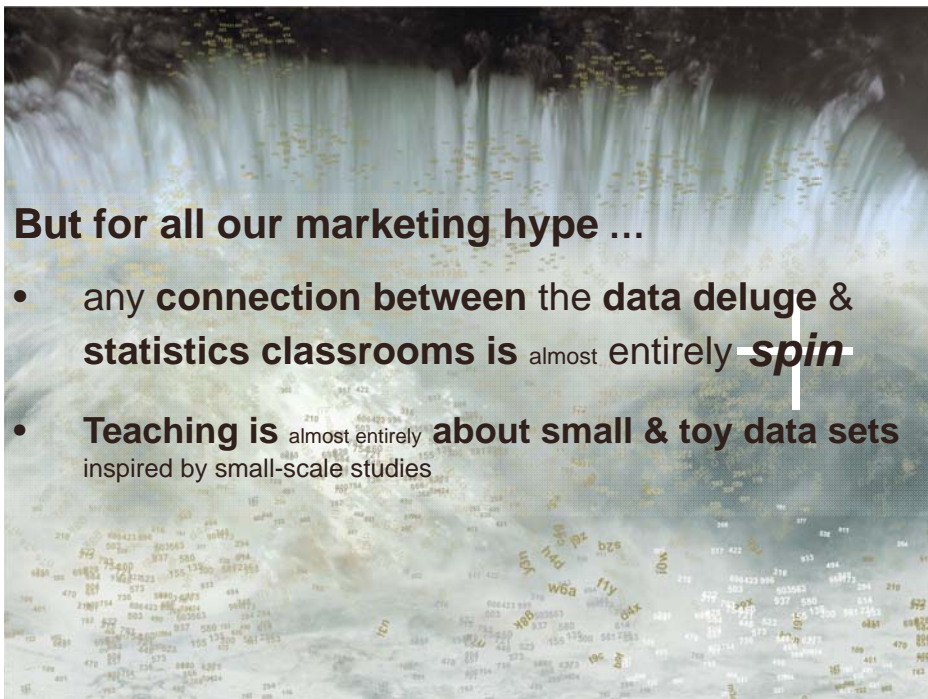*Manuscript:* www.stat.auckland.ac.nz/~wild/TEMP/bootstrap.pdf
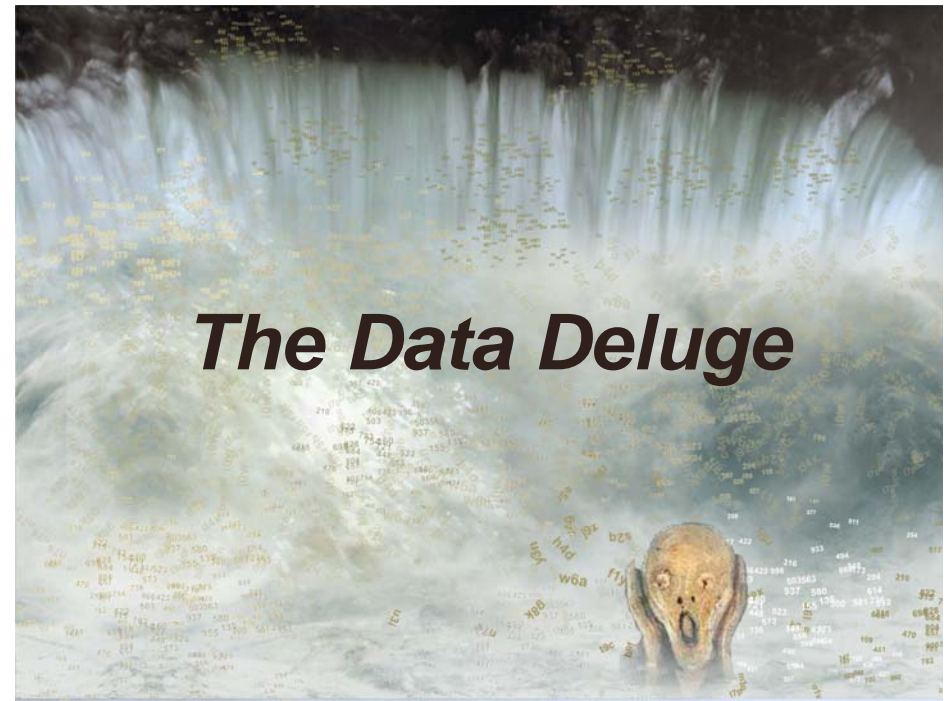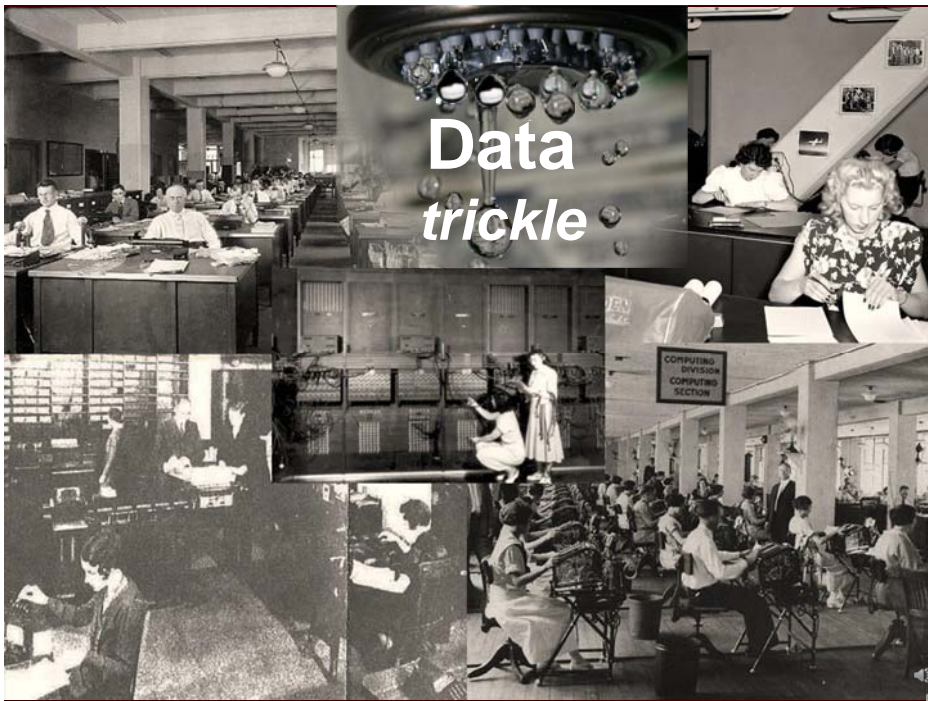
THE UNIVERSITY OF AUCKLAND
DEPARTMENT OF STATISTICS
USCOTS 2013

---

## The Need for Speed ?

THE UNIVERSITY OF AUCKLAND
DEPARTMENT OF STATISTICS
USCOTS 2013

---

First ...

# A Brief History of
# ~~Data~~

---

In the beginning there was ...

## Slide 1

**Data**
*trickle*

## Slide 2

*The Data Deluge*

## Slide 3

**But for all our marketing hype ...**

- any **connection between** the **data deluge** & **statistics classrooms is** almost entirely *spin*

- **Teaching is** almost entirely **about small & toy data sets** inspired by small-scale studies

## Slide 4

**But the data world …**
**is getting a whole lot bigger**

*Further, Faster, Better*
- Data world exploding
- Green shoots in Software
- Vision
- Future is visual
- Accelerators

## But the data world … is getting a whole lot bigger

- **There is an explosion in the …**
  - **quantities of data being collected**
  - **conceptions of what constitutes data**
  - **settings in which it can arise**
  - **ways of looking at it**

---

*Further, Faster, Better*
- Data world exploding
- Green shoots in Software
- Vision
- Future is visual
- Accelerators

## But the data world … is getting a whole lot bigger

**Can't just keep illuminating same small patch**

- **Need to get much …**
  - *further*
  - *faster*
  - & with **better** *comprehension*

---

*Further, Faster, Better*
- Data world exploding
- Green shoots in Software
- Vision
- Future is visual
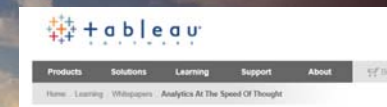- Accelerators

## Green shoots in software …

---

*Further, Faster, Better*
- Data world exploding
- Green shoots in Software
- Vision
- Future is visual
- Accelerators

## Green shoots in software …

- **Hans Rossling**
- **Visualisation generally**
  - **Tableau**

+ableau

Products   Solutions   Learning   Support   About

Home :: Learning   Whitepapers   Analytics At The Speed Of Thought

*"Visualisation is a gateway drug to statistics"*

*"People who look at visualisations will start asking statistically important questions … even without knowing the jargon."*

**-- Martin Wattenberg**
*in interview for New Scientist by Peter Aldhous, 5 February 2011, p44.*

- **AdviseStat**

**AdviseStat, First in the Next Generation of Analytics**

AdviseStat is your statistics expert-in-a-box. It does all of the heavy statistical thinking, and you get to take all of the credit!

**Simple to Work With**

Tell AdviseStat what you want to see, and it will take care of everything else. AdviseStat uses a unique button-tree interface that allows you to command it with a simple sentence:

· "Predict height based on weight."

· "Compare salaries when grouped by gender."

· "Forecast sales."

That's it. AdviseStat will ask a few yes-or-no follow-up questions, and then run the entire analysis itself.

**Does the Thinking For You**

AdviseStat is the first artificially intelligent statistical assistant. Developed by a statistician, it makes its own observations about the data it sees and then chooses the most appropriate methods. All the transformations are done for you in the background. It makes corrections, like imputing missing data, all on its own. It is the closest thing to having a statistician there with you.

**Guides You, and Explains Your Results**

When the analysis is complete, most statistical software will spit out a couple coefficients and graphs, without any explanation or context. AdviseStat gives you a full whitepaper with text that's been customized to your specific findings. It starts by explaining about the type of analysis it chose, then talks about the corrections it made and the significant results it found.

**Raises Red Flags**

The great thing about an expert statistics program is that it can keep you out of pitfalls you might not even have known were there.

AdviseStat catches subtle statistical mistakes and other lurking factors that even expert statisticians would have a hard time finding, such as:

*Further, Faster, Better*
- Data world exploding
- **Green shoots in Software**
- Vision
- Future is visual
- Accelerators

**Leland Wilkinson**



# A Growing Gap

Practice & potential

GAP

GAP

Education

Time

THE UNIVERSITY OF AUCKLAND
DEPARTMENT OF STATISTICS

USCOTS 2013



# Opening Up the World of data

## Need to open it up *wider, faster*



*Further, Faster, Better*
- Data world exploding
- Green shoots in Software
- **Vision**
- Future is visual
- Accelerators

*Uncovering the stories* in a sea of data

*Ideal*

# Data analysis

Slide 1:

"the future of improved statistical understanding is *visual*"

"*The Eyes have it*"

Slide 2:

"the future of improved statistical understanding is *visual*"

improved **understanding of ...**  →  data

→  inferential concepts

Slide 3:

**VISION Statement for Early Statistics**

- To *create excitement* about
  - "What I can do with data &
  - What data can do for me"

Slide 4:

**Some things**

- **Some things change**
  - Exploding world of data
  - Need to convey more of this more quickly

The Data Deluge

- **Some things stay the same**
  - Available time
  - Inability to hold more than 4-7 ideas in working memory

- So **something's got to give** !!
  - Details of how we construct things
    - Segmentation th... behind them are

And software (& simulation & visualisation) **provides the key**

THE UNIVERSITY OF AUCKLAND
DEPARTMENT OF STATISTICS

USCOTS 2013

## Slide 1

Further, Faster, Better
- Data world exploding
- Green shoots in Software
- Vision
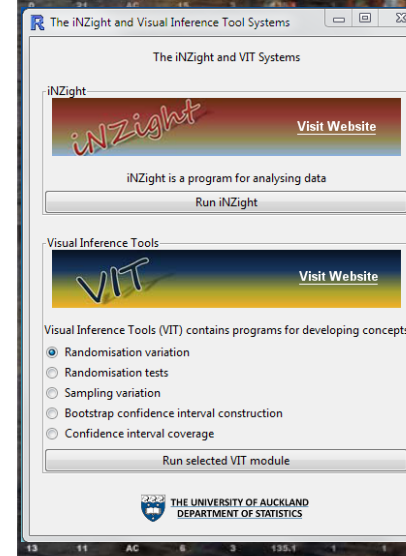- **Future is visual**
- **Accelerators**

***Uncovering the stories***
in a sea of data

**A belief that this is possible …**
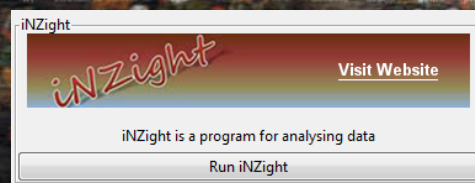
Data analysis

## Slide 2

Further, Faster, Better
- Data world exploding
- Green shoots in Software
- Vision
- **Future is visual**
- **Accelerators**

led to work on **tools to facilitate** this future

The iNZight and Visual Inference Tool Systems

The iNZight and VIT Systems

iNZight

iNZight

**Visit Website**

iNZight is a program for analysing data

Run iNZight

Visual Inference Tools

VIT

**Visit Website**

Visual Inference Tools (VIT) contains programs for developing concepts
- Randomisation variation
- Randomisation tests
- Sampling variation
- Bootstrap confidence interval construction
- Confidence interval coverage

Run selected VIT module

THE UNIVERSITY OF AUCKLAND
DEPARTMENT OF STATISTICS

## Slide 3

Further, Faster, Better
- Data world exploding
- Green shoots in Software
- Vision
- **Future is visual**
- **Accelerators**

**A belief that this is possible**
led to work on **tools to facilitate** this future

iNZight

iNZight

**Visit Website**

iNZight is a program for analysing data
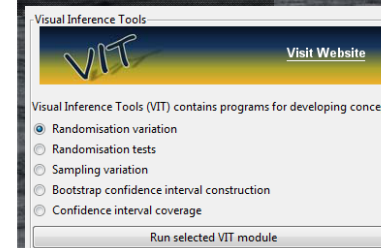
Run iNZight

- a simple **data exploration system**
  - to allow **uncovering** of **stories** in **several dimensions**
  - using very simple graphical forms
    - Stealing good ideas and metaphors from …
      e.g. Hans Rossling; Jim Ridgway, Smart Centre, Durham U.
  - Facilitating …
  ***"exploring data at the speed of your thoughts"***

(Similar to Tableau Software slogan)

## Slide 4

Further, Faster, Better
- Data world exploding
- Green shoots in Software
- Vision
- **Future is visual**
- **Accelerators**

**A belief that this is possible**
led to work on **tools to facilitate** this future

Visual Inference Tools

VIT

**Visit Website**

Visual Inference Tools (VIT) contains programs for developing concepts
- Randomisation variation
- Randomisation tests
- Sampling variation
- Bootstrap confidence interval construction
- Confidence interval coverage

Run selected VIT module

**Visualisation modules**

- for building **conceptual understanding**
  - To help **dispel** *the Dark Magic*

***"I have a feel for what's going on here"***

## Slide 1: Time Series

*Further, Faster, Better*
- Data world exploding
- Green shoots in Software
- Vision
- Future is visual
- **Accelerators**

iNZight

# Time Series

## Slide 2: Shift of focus

*Further, Faster, Better*
- Data world exploding
- Green shoots in Software
- Vision
- Future is visual
- **Accelerators**

iNZight

# Shift of focus

**From:**

- "*What are all the things I have to do to get the output I'm meant to produce?*"

**To:**

- "*What are the questions?*"
- "*What can I **see**?*"
- "*What does that tell me?*"

Using *"enabling software"* to "**cut out the middle man**"

## Slide 3: But …

*Further, Faster, Better*
- Data world exploding
- Green shoots in Software
- Vision
- Future is visual
- Accelerators
- **Excitement & "safer sex"**

# But …

- **with arousal to the pleasures of discovery** …



CAUTION SAFE SEX

- **need fundamental statistical "safe sex" messages**

## Slide 4

*Further, Faster, Better*
- Data world exploding
- Green shoots in Software
- Vision
- Future is visual
- Accelerators
- **Excitement & "safer sex"**

**Looking at the world using data is like looking through a window with ripples in the glass**

"*What I see … is not quite the way it really is*"

We see **a blend of fact & artefact**

**There are stories … & there are "stories"**

Data analysis



**Stories vs "stories"**

Fact

*versus*

Artefact



**Stories vs "stories"**

Fact

*versus*

Artefact

*not always an improvement!*



**Main causes of artefact**

*Further, Faster, Better*
- Data world exploding
- Green shoots in Software
- Vision
- Future is visual
- Accelerators
- Excitement & "safer sex"

Bias

Random Error

Confounding

## Slide 1

*Further, Faster, Better*
- **Data world exploding**
- **Green shoots in Software**
- **Vision**
- **Future is visual**
- **Accelerators**
- **Excitement & "safer sex"**

# We see a blend of fact & artefact

- ## Statistical study design
  - employing random sampling or randomised experimentation (in conjunction with other tricks of the trade)

  used to minimize artifact

- ## Inference based on randomness theory
  - most ***obviously relevant*** and ***valid*** where (all of) the randomness is introduced by the study design
  - other uses are ***based on models***

    that ***assume random mechanisms*** are at work somewhere in the process that generated the data

    *Hugely more complex & difficult thinking*

THE UNIVERSITY OF AUCKLAND
DEPARTMENT OF STATISTICS                                    USCOTS 2013

## Slide 2

# Start with randomness by design
(The basic concepts can transfer to modelling contexts later)

- Confidence interval ideas
  - arise most directly and simply in sampling contexts

- Significan...

  > ***"The traditional road to statistical knowledge is blocked, for most, by a formidable wall of mathematics."***
  >
  > – Efron & Tibshirani (1993)

  - arise mo...
    settings

- In both cases
  - can motivate and convey all of the essential ingredients ***entirely visually*** (without a formula in sight)

THE UNIVERSITY OF AUCKLAND
DEPARTMENT OF STATISTICS                                    USCOTS 2013

## Slide 3

# Start with randomness by design
(The basic concepts can transfer to modelling contexts later)

- C...
  > ***Approach***
  > - ***Get basic ideas in place first***
  >   *- intuitively, visually*
  > - ***Can mathematize later***

  ...xts

- Si...

  ...nment settings (experiments)

- In both cases
  - can motivate and convey all of the essential ingredients ***entirely visually*** (without a formula in sight)
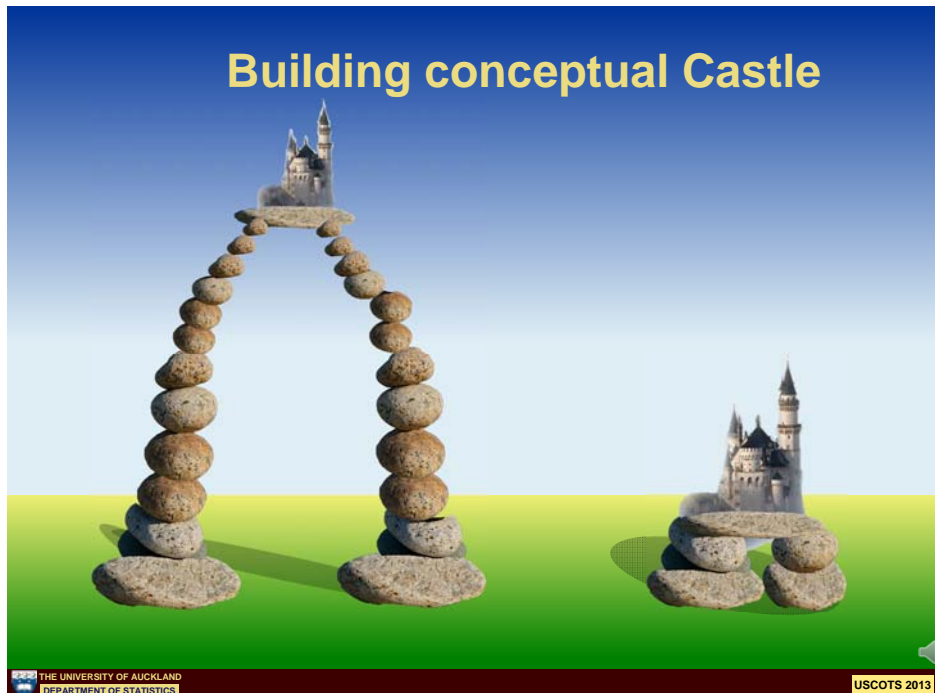
THE UNIVERSITY OF AUCKLAND
DEPARTMENT OF STATISTICS                                    USCOTS 2013

## Slide 4



**Building conceptual Castle**

THE UNIVERSITY OF AUCKLAND
DEPARTMENT OF STATISTICS                                    USCOTS 2013

## Slide 1

**Building conceptual Castle**

Inferences from Normal Theory

Inferences from Bootstrap & randomisation

## Slide 2

**Principle**

*The inferential method should mirror the*

- **I'm going to lead down some sets of conceptual sequences**
- **It's for you to judge how well they hang together**

1. ...his principle
2. ...themselves to visual treatments
3. Connect better to intuition

Will spend rest of the talk exploring and demonstrating this assertion

## Slide 3

**Statistical inference
& "sampling variation"**

## Slide 4

**Original version**

Use *boxplots*
*Track*
- medians
- quartiles

Makes …

# Intuition

*"Where is truth likely to lie?"*

**I** got this

*Truth* is seldom further from **my data median** than this

**Problem: *I don't actually see width of this "uncertainty" band***

**Why?:** I only see **one** frame of sampling variation movie

**So:** We need some sort of estimate of the width of uncertainty band *from the single sample itself*

*How ????*

*Seems impossible!!*

---

# Enter Brad Efron & the ^simple Bootstrap

Efron (1979)

*(1996 photo)*

**I wonder** if "*sampling with replacement from the sample*" will **mimic** the process of "*sampling from the population*"

---

# Bootstrap

- **Bootstrap Sampling with replacement**
  - **What is it?**
  - **What does it do?**
  - **How could we use it?**
  - **Why might it work?**
- **Constructing Bootstrap intervals**
- **Does it work?**

---

*Bootstrap*
- Sampling with replacement
  - What is it?
  - What does it do?
  - How could we use it?
  - Why might it work?
- Constructing Bootstrap intervals
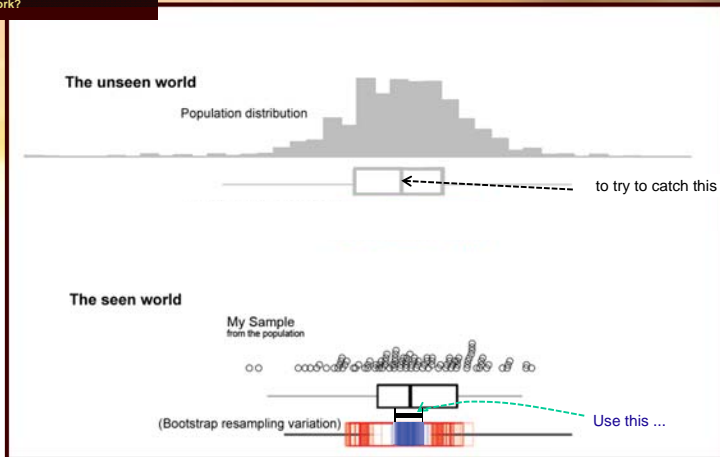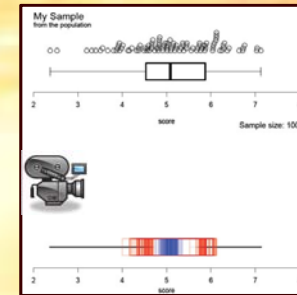- Does it work?

# Bootstrap

*"Re-sampling with replacement"*

- *What is it?*
- *What does it do?*

Play

# How could we use it?

The unseen world

Population distribution

to try to catch this

The seen world

My Sample
from the population

(Bootstrap resampling variation)

Use this ...

---

# Bootstrap

**How could we use it?**

PHILIP MARLOWE
CRIMINAL & CIVIL
INVESTIGATIONS

My Sample
from the population

score

Sample size: 100

score

---

# Bootstrap

**"Re-sampling with replacement"**

- *Why might it work?*

**Plausibility**

---

**Sampling**
**from population**

Looks very similar

**Bootstrap**
**Re-Sampling**
**from sample**

*Single Means*

Slide 1 (top-left):
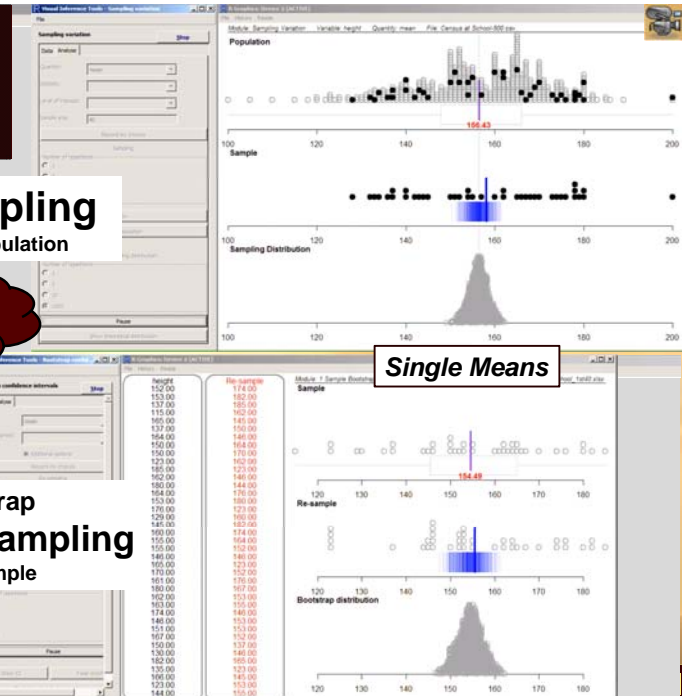Bootstrap
- Sampling with replacement
  - What is it?
  - What does it do?
  - How could we use it?
  - Why might it work?
- Constructing Bootstrap intervals
- Does it work?

**Sampling** from population

Looks very similar

**Bootstrap Re-Sampling** from sample

*Regression slopes*

THE UNIVERSITY OF AUCKLAND
DEPARTMENT OF STATISTICS

Slide 2 (top-right):
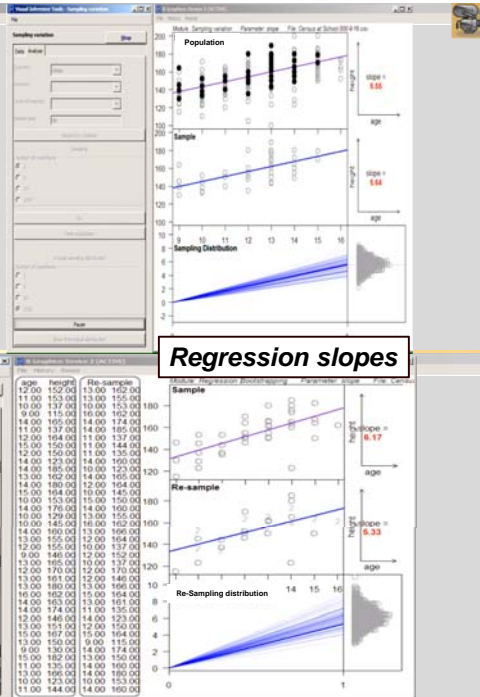Bootstrap
- Sampling with replacement
  - What is it?
  - What does it do?
  - How could we use it?
  - Why might it work?
- Constructing Bootstrap intervals
- Does it work?

**Sampling** from population

Looks very similar

**Bootstrap Re-Sampling** from sample

*Diffs in Propⁿs*

THE UNIVERSITY OF AUCKLAND
DEPARTMENT OF STATISTICS

Slide 3 (bottom-left):
Bootstrap
- Sampling with replacement
  - What is it?
  - What does it do?
  - How could we use it?
  - Why might it work?
- Constructing Bootstrap intervals
- Does it work?

**Sampling** from population

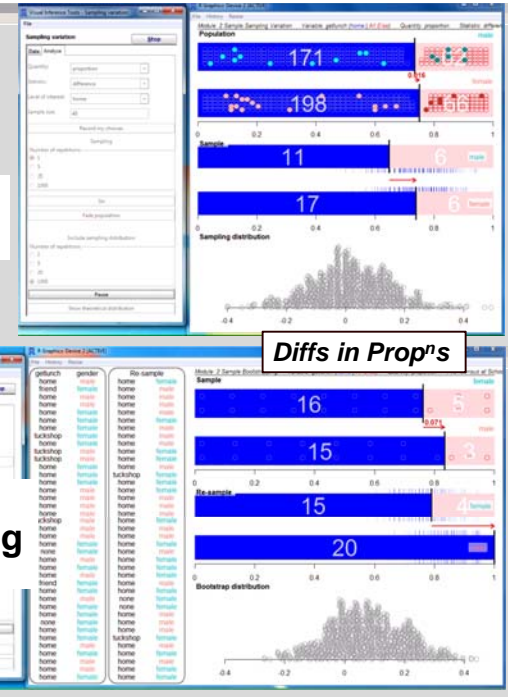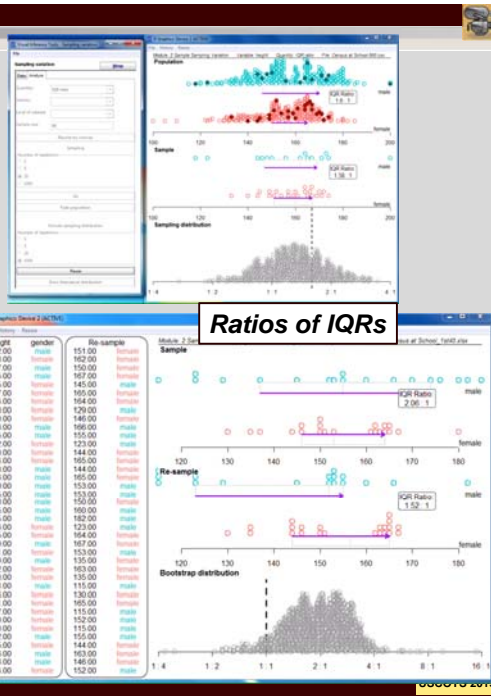Looks very similar

**Bootstrap Re-Sampling** from sample

*Ratios of IQRs*

THE UNIVERSITY OF AUCKLAND
DEPARTMENT OF STATISTICS

Slide 4 (bottom-right):
Bootstrap
- Sampling with replacement
  - What is it?
  - What does it do?
  - How could we use it?
  - Why might it work?
- **Constructing Bootstrap intervals**
- Does it work?

**Bootstrap intervals**

**Using re-sampling to construct an interval**

Construction

- *How is it done?*

Centre    Diffs    Reg Coef

THE UNIVERSITY OF AUCKLAND
DEPARTMENT OF STATISTICS

USCOTS 2013

## Slide 1: Bootstrap intervals

Looks plausible …

Does it work ??

"Simulate & see"

to try to catch this
Use this …

The seen world

My Sample
from my population

(Bootstrap resampling variation)

Play

## Slide 2: Embed in Discovery learning for Methodology

Need

Idea

Does it work?

Simulate and see

*Doesn't*    *Works*

Use it

encounter situations where doesn't work

How does this differ from grown-up statistics?

Only really omits *"Does it work in Asymptopia?"*

## Slide 3: Bootstrap – one basic, widely applicable idea

**Not a huge step now …**

- to accepting that
  - we can & should

  put uncertainty intervals around most everything

- And further …

## Slide 4: More complex: Scatterplot with smoother



FEV by age

## Slide 1

FEV by age

## Slide 2

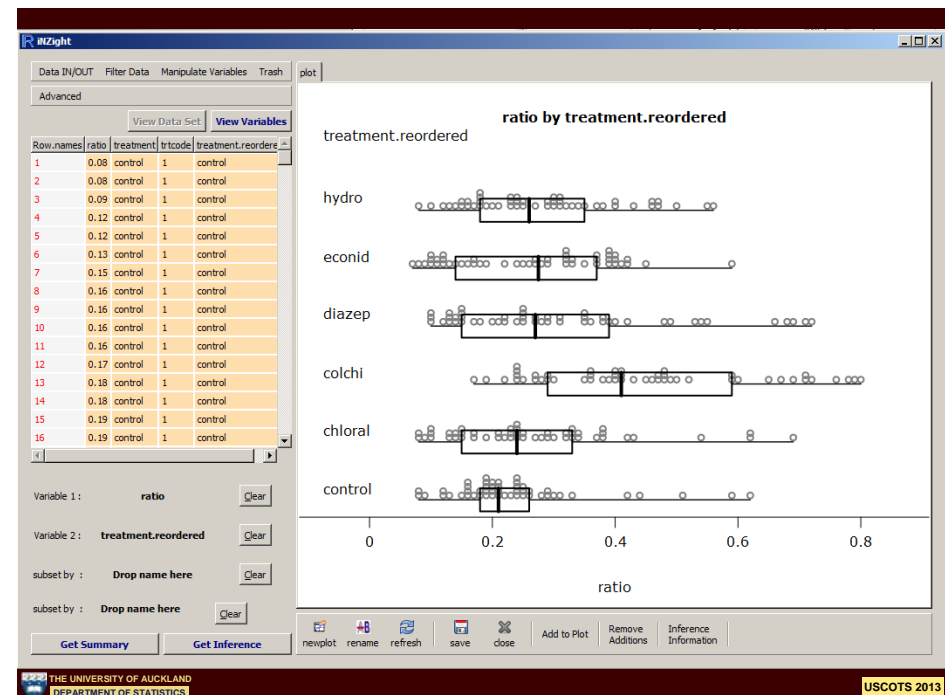# Randomised experiments and Randomisation tests

## Slide 3

# Cell-ratio data

[Source: Auckland's Cancer Research Unit]

- Experimental units are samples of human blood cells
  - Grown in cell culture and then treated with one of
    - *chloral hydrate, hydroquinone, diazepam, econidazole, and colchicine*
      - some of which are known to be potent carcinogens.
      - The carcinogens act by breaking chromosomes, and thus disrupts cell division
      - Broken fragments of chromosome are left as micronuclei, and the average ratio of the size of a micronucleus to its parent cell nucleus is measured.
        - » The more carcinogenic the chemical, the higher the ratio tends to be
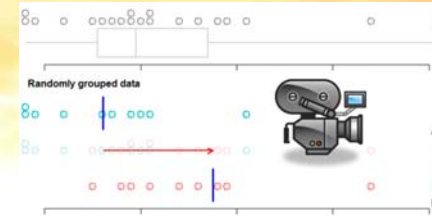  - A random mechanism determined which samples of cells received which chemical treatment

## Slide 4

## Experiential Context

- Simple randomised experiment
  - e.g. drug vs control
  - Follows experiences on "why do randomised experiments"
    - Following up message that randomisation is best way we know of balancing groups on unforeseen factors so that, apart from treatment, we are
      **comparing like with like** *("fair test")*

---

## Experiential Context

1st new message:

- ***Randomisation is best, but not perfect …***



Randomly grouped data

**"*What I see isn't exactly the way it is*"**
**Why?: *Randomisation variation***
*Randomisation alone can make it look like there is a difference between groups*
*(the apparent differences result simply from who, by chance, ends up in what group)*

We should only be impressed by experimental differences if they are larger than those produced by random labelling alone

---

## Randomisation test visualisations

**Can random labelling alone ("*chance alone*") produce differences as large those as I'm seeing?**

**2 groups**

**3 groups**

Play

**Orig data**

**Proportions**

---

## Epilogue

- **Got a quite a long way with tiny number of ideas**
- Can broaden applications greatly without difficulty
  - so long as stick to …
    - intervals in contexts that emphasize sampling
      (Expand the complexity of sampling, range of quantities estimated)
    - Significance tests in contexts that emphasize random assignment
      (Expand the complexity of design, range of quantities estimated)
  - But when you want to do …
    - interval estimates of effect sizes in experimental contexts
    - significance testing in sampling contexts
    you suddenly have to wheel in many more concepts

## Current Module Capabilities

| Track/investigate/use behaviour of … | VIT Module Name | | | |
|---|---|---|---|---|
| | *Motivation of need for method*<br>**Randomisation Variation**<br>(*random assignment of group labels*) | *Inferential* **Method**<br>**Randomisation Tests**<br>(*random re-assignment of group labels*) | *Motivation of need for method*<br>**Sampling Variation**<br>(*random sampling from Popn*) | *Inferential* **Method**<br>**Bootstrap Conf. Ints**<br>(*Random re-sampling from sample*) |
| **1-variable** | | | | |
| Numeric | | | | |
| Categorical | | | | |
| **2-variables** | | | | |
| Num \| Cat (2 grps) | Diffs in: Means, Medians;<br>2-sample t-stats; IQR Ratio | Diffs in: Means, Medians;<br>2-sample t-stats; IQR Ratio | Whole boxplots<br>Diffs in: Means, Medians;<br>IQR Ratio | Diffs in: Means, Medians;<br>IQR Ratio |
| Num \| Cat (k grps) | Av. Deviation, F<br>"pseudo-F" | Av. Deviation,<br>"pseudo-F" | | |
| Cat \| Cat (2 grps) | Diffs in: Proportions | Diffs in: Propor | | |
| Cat \| Cat (k grps) | Av. Deviation, Chisq | Av. Deviation, Chisq | | |
| Num \| Num | Reg Slope (default)<br>Paired diffs (option) | | | |

**Columns differ in:**
- **"parent data"** being operated on (popn/sample)
- **operation** performed (sampling/assigning grp labels)
- Motivation vs Method

**Rows differ in:**
- **Quantity/statistic** worked with

**Principle:**
- **maximise "same-ness"** across & down so …
  - "same-nesses" reinforced
  - visual differences relate to essential differences
    - to improve conceptual transfer

---

## Epilogue

- **What we don't have in this development**
  - Beginner-killing abstractions
    - e.g. null hypotheses, parameters vs estimates, test statistics, formal distributions …
  - Dense clouds of details
  - Dependence on (poorly understood) mathematical ideas

- **What we do have**
  - "Concrete" ideas that make sense in the context
    - e.g. "Can random labelling do this?"
  - Fast access to a wide range of important applications
  - Substantial body of intuition and experience as a foundation to then build abstractions upon

---

## Addressing the Need for Speed

- **iNZight-type software can facilitate a fast, accessible way in to understanding a much wider spectrum of data types**

- **VIT-type software can facilitate a fast , accessible way in to understanding basic inferential conceptions**

**My "vision" is**

- **initially create an appreciation of a very wide array of data types and what they can tell you**

  - **and only then back fill the details** (for those who need it)

---

## The fundamental things apply …

Bias

Random Error

Confounding

Thank you