## Slide 1

# *Interactive applications written in R to accelerate statistical learning*

**Chris Wild**

Department of Statistics
University of Auckland, New Zealand

https://www.stat.auckland.ac.nz/~wild/

## Slide 2

# What are we doing?

## Slide 3

# Applications

- iNZight
- iNZight Lite
- VIT: Visual Inference Tools
- Mortality Calculator
- Table Maker
- Others

## Slide 4

# inzight

- An interactive data analysis system that has R "unseen under the hood"

https://www.stat.auckland.ac.nz/~wild/iNZight/

**https://www.stat.auckland.ac.nz/~wild/iNZight/user_guides/interface**

## Slide 1 (top-left)

**1956**

(playing across time)

## Slide 2 (top-right)

- An interactive data analysis system that has R unseen "under the hood"
- Interactivity from John Verzani's gwidgets2
  - High-level uses RGtk2 which uses Gtk+
- It has R inside it but users do not see R
  - Windows version also packages Gtk+ inside it
- Caters for beginners through to quite advanced modelling
  - Youngest users in early levels of high school

THE UNIVERSITY OF AUCKLAND
DEPARTMENT OF STATISTICS

useR! 2015

## Slide 3 (bottom-left)

o⁰⁸⁸⁸ inzight   Get iNZight   User Guides   Support   About   Related

Easily **explore data** and **discover trends**
*without* learning complex software

Life Expecancy versus Income from 1952 till 2012

[2008]

Region
America
East Asia & Pacific
Europe & Central Asia
Middle East & North Africa
South Asia
Sub-Saharan Africa

**Download Now**
for Windows
(Mac or Linux downloads)
Latest Version: **2.2** (what's new?)
Release Date: **16 June 2015**
Price: **100% FREE!**

Or try our online application:
**inzight Lite**

**Why??**
***BYOD***

THE UNIVERSITY OF AUCKLAND
DEPARTMENT OF STATISTICS

useR! 2015

("Bring your own device"-driven need for tools to work on tablets and even phones)

## Slide 4 (bottom-right)

o⁰⁸⁸⁸ inzight lite   About   File ▾   Visualize   Row operations ▾   Manipulate variables ▾   Advanced ▾

**inzight lite**

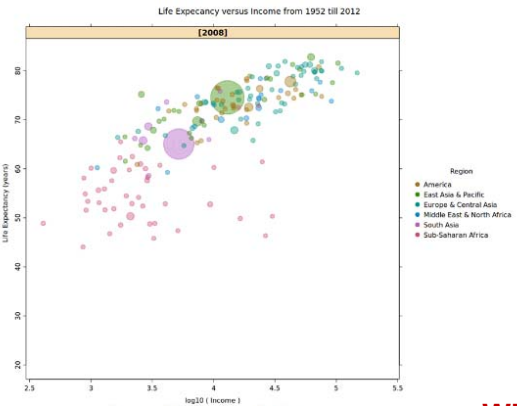**inzight** is a simple data analysis system which was initially designed for high school students to help explore data fast and easy without having to learn complex statistical software. By popular demand, it has been extended to support 3D graphics, multivariate analysis, and time series analysis. **inzight** lite is an online version of the full software, which goes a long way to make it more accessible to a wide range of users.

**inzight** lite lets you import your own data set or explore one of the many example data sets; Even if you don't have a formal background in computer programming or statistics, you can conduct statistical analysis on the data, and modify it to explore hidden secrets behind the data. If you are an expert programmer or statistician, you can contribute to the project by sending us feedback about our source code on github (click "R Source Code" at the bottom of the screen).

This project is led by Professor Chris Wild and has been primarily supported by the Department of Statistics at the University of Auckland, with additional support from Statistics New Zealand and the NZ Ministry of Education

iZight Project | R Source Code | Contact Us
opyright 2015 iNZight | All Rights Reserved

Statistics New Zealand   MINISTRY OF EDUCATION   THE UNIVERSITY OF AUCKLAND   NEW ZEALAND

THE UNIVERSITY OF AUCKLAND
DEPARTMENT OF STATISTICS

useR! 2015

http://new.censusatschool.org.nz/
(demo in which CaS data tools call up iNZight Lite for analysis)



**Visual Inference Tools**

- VIT **(Visual Inference Tools)** is about conceptual development

  https://www.stat.auckland.ac.nz/~wild/VIT/



**Visual Inference Tools**



**Visual Inference Tools**

**Slide 1:**

Module: 2 Sample Sampling Variation    Variable: getlunch (*home* | *All Else*)    Quantity: proportion    Statistic: differen

**Population**    male

171

0.016

198    female

Sample

12    4 male

16    8 female

Sampling distribution

**Slide 2:**

# Why are we doing it?

**Slide 3:**

# The data world …
## has gotten a whole lot bigger

**Slide 4:**

# The data world …
## has gotten a whole lot bigger

## Can't just keep illuminating same small patch

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

## Slide 1

### "Middleware"

(Not in the technical sense)

- software aimed at …
  - allowing student to experience
    - as much as possible of "discovery in the data world"
    - in the least possible time
  - Minimal learning curves, everything happens instantly & you don't have to remember anything

THE UNIVERSITY OF AUCKLAND
DEPARTMENT OF STATISTICS

useR! 2015

## Slide 2

### *Initial data analysis experiences should feel like this!*



## Slide 3

### Developing in concert …

**Software** ⟷ Educational **Experiences**

Desired **Capabilities**

*(Trade-offs everywhere)*

THE UNIVERSITY OF AUCKLAND
DEPARTMENT OF STATISTICS

useR! 2015

## Slide 4

### How are we doing it?

THE UNIVERSITY OF AUCKLAND
DEPARTMENT OF STATISTICS

useR! 2015

## How are we doing it?

- Sitting in behind iNZight and VIT are sets of R packages
  - (Can be run directly from R)
  - User interfaces use John Verzani's **gwidgets**
    - which uses Gtk+ via RGtk2
  - Each module has a "receiver function" which manages the mapping of user choices to R calls
- iNZight Lite is same with R Shiny user interface
  - Connecting to exactly the same receiver-functions

## Mortality "Calculator"

## Census At School Table Maker

**(1 & 2-way tables)**

## Shiny applications

- iNZight Lite (versions), mortality calculator, probability visualisations, …

- Each application
  - Is in its own docker container with R shiny server (the free one) with R and relevant libraries.
  - Docker container inside a virtual machine
  - Replicate virtual machines if needed to meet demand
- We want others to be able …
  - just pick these things up and put down on own servers with a minimum of effort and knowledge

- **Volunteers anyone???**

## My projects …

*(Visual Inference Tools)*

**iNZight**

**VIT**

**Data Analysis**  ⟷  **Conceptual development**

*Research*

**MOOC**

*"Data to Insight"*
*(prototype for "getting further faster")*

*(Trade-offs everywhere)*

---

## MOOC — DATA to insight — D2i DATA→insight

**Week 1: Introduction**; gee whiz; software; data
**Week 2: Boot Camp** (Basic Training) – you'll see bits of this
**Week 3 & 4: Relationships between variables**
- Relationships between **categorical variables** – you'll see bits of this
- Relationships between **numeric variables**
  - Trend, scatter & outliers; Clusters
  - Prediction with uncertainty
  - Association & Correlation
  - Trends: Lines, curves & smoothers
  - Large data-set problems and solutions
  - Overprinting, jitter & transparency, granularity & point size, running quantiles
  - **More variables** with size, colour and subsetting

**Week 5: Why "what I see is never quite the way it really is"**
- Measurement and "selection" bias
- Sampling error and sampling variation
- Causation and confounding

**Week 6: Estimation with confidence via bootstrap**
**Week 7: Designed Experiments and randomisation tests**
**Week 8: Time series** stressing seasonal series with forecasting & comparing related series

https://www.stat.auckland.ac.nz/~wild/d2i/4StatEducators/

---

https://www.stat.auckland.ac.nz/~wild/wildaboutstatistics/ (index to youTube channel)

**Wild About Statistics**

**Week 3: RELATIONSHIPS**

- **Introduction to Relationships** (*Why we care; Outcome & predictor variables*) [2:52]
- **Relationships between Categorical Variables** (*Exploration using separate bar charts and side-by-side bar charts*) [6:22]
- **Changes across subgroups** (*Exploring effects of a 3rd and 4th variable on a relationship via subsetting, tiling & movement*) [4:52]
- **Relationships between numeric variables** (*Scatter plots; Trend, scatter & outliers; Clustering*) [5:17]
- **Trend, Scatter & Outliers** (*Examples; Prediction & prediction intervals; Training the eye*) [6:42]

**Week 4: MORE RELATIONSHIPS**

- **More Relationships** (*Introduction to week's coverage*) [1:32]
- **Lines, curves and smoothers** (*Lines, curves & smoothers; Least squares*) [4:01]
- **Overcoming Perceptual Problems** (*Problems with large datasets; Overprinting; Jitter; Varying transparency & point size; Running quantiles; Tile-density plots*) [7:05]
- **Diving deeper with more variables** (*Additional variables using colour, subsetting & movement; different trends per group*) [5:20]
- **Our Changing Health and Wealth** (*Case study using up to 6 variables at once by employing colour, size, subsetting, matrices of tiled plots and movement*) [5:41]

**Week 5: WHY WHAT WE SEE IS NEVER QUITE THE WAY IT REALLY IS**

- **Why what I see is never quite the way it really is** (*Intro to week; Facts & artefacts*) [3:04]
- **Bad Data** (*"Measurement" issues/ biases; "Selection" biases; missingness*) [7:02]
- **Causes and Confounding Variables, Part I** (*Confounding & adjustment*) [6:38]
- **Causes and Confounding Variables, Part II** (*Confounding & adjustment*) [3:57]
- **Random Error, Part I** (*Random variation/error; effect of sample size*) [7:06]
- **Random Error, Part II** (*Random variation/error; effect of sample size; biases*) [6:05]

---

## What do I want from this session??

https://www.stat.auckland.ac.nz/~wild/

# *Potential collaborators!!*

- Any aspects of any of these projects
- New R packages to link to
- Skills we don't have
- … ????

Initial data analysis experiences should feel like this!

Thank you


"Don't make students crawl over broken glass ... before a desire has been aroused for what's on the other side"