



R: A Swiss Army Knife for Market Research

Enhancing work practices with R

Martin Chan – Consultant, Rainmakers CSI
12th September, 2018



This presentation is about using R in **business environments or conditions** where its usage is less intuitively suitable...

... and a story of how R has been **transformative** for our work practices



Who are **we**?
What do we do?

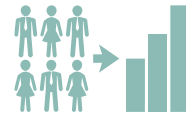


What **How** **Where**



Answer strategic questions to help our clients grow profitably

Inform strategic decisions through creative and analytical thinking



Estimate size of market, and size of opportunity, anticipating trends



Utilise resources available within an organisation (e.g. stakeholder knowledge, existing research)



Develop consumer targeting frameworks – identify core consumers and areas of opportunities



Across multiple industries and markets

- FMCG
- Finance & Insurance
- Travel
- Media



A team with a mix of backgrounds from strategy, brand planning, marketing, research...

(where analytics is a key, but only one of the components of our work)



What's different about the use of **R** in our work?



Challenges of Using R

1

**Nature of data:
disparate, patchy**

2

**Nature of U&A
data**

3

**Client
requirement and
expectations (of
outputs)**

U&A – research with the aim to understand a market and identify growth opportunities by answering questions on whom to target, with what, and how.

(Source: <https://www.ipsos.com/en/ipsos-encyclopedia-usage-attitude-surveys-ua>)



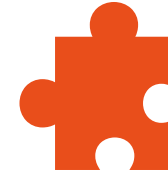
- ① Nature of data:
disparate, patchy



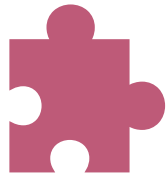
Our data come together in different forms, like pieces of **JIGSAW**



Historical survey data – often designed for different purposes and collected from different samples



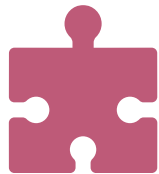
Stakeholder Interviews
(Qualitative)



Population / Demographic data (from census, World Bank research etc.)



Customer Interviews (Qualitative)



Pricing data
(e.g. Euromonitor, Nielsen)



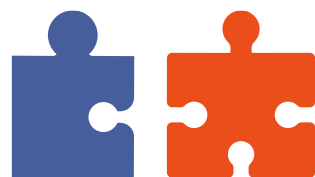
Historical segmentation work
(past work produced by client's suppliers)



Primary research
(e.g. U&A, customer satisfaction surveys)



Disparate and patchy data...



... renders it more difficult to reap the benefits offered by R

Challenges in Using R

1

Ad-hoc data – easier to explore in Excel

Traditionally Excel is an easy tool for these (albeit chaotic) situations when you need to 'play around' to figure out your plan of attack

2

Data is small or 'non-repetitive' – not worth automating

Perception of using R – for automating analysis - is only justifiable if the benefits of reproducibility exceed the costs of setting up code (prominent for small data sets)

3

Poorly structured / stored data – more time-consuming to use R to clean

Condition of the input data (think merged cells in Excel, or data saved in PowerPoint) almost make R even more time-consuming in the short-term



Integrating data types

Integrating analysis from different sources of data under one roof



Package for Qualitative Data Analysis

Qualitative Research

stakeholder interviews, telephone interviews; focus groups, workshops



Output

Understand the “why” and the reasoning behind behaviour, or develop hypotheses / segments – development of a framework of themes and quotes



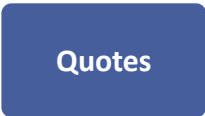
Qualitative Analysis (Without R)



Interview or group takes place, which is recorded and transcribed into a text or Word file



- Transcripts are annotated and interpreted
- Possible themes / hypotheses are documented in a separate document
- Quotes are copied out and pasted into a separate document, with references to the respondent background and interview number



Keep refining **framework** until it is a fair and sensible interpretation of the different themes emerging from the interviews



A framework of themes/hypotheses or a story is developed – and a report is produced using the analysis of the annotations, themes, and quotes

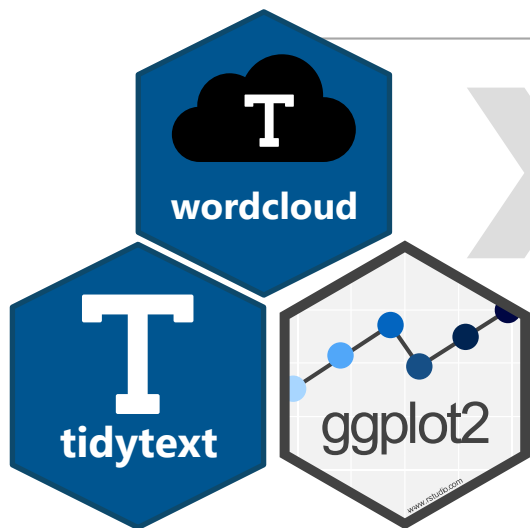


RQDA: What is it used for?

1. A more systematic tool for analysing qualitative data
2. **GUI launched in R for marking up themes**, quotes from qualitative interviews
3. **Creates an analysis output in sqlite database**, allowing further analysis or production of outputs within R
4. **Integrated with R** – use sqlite to extract data, enable exploration of insights through text mining techniques



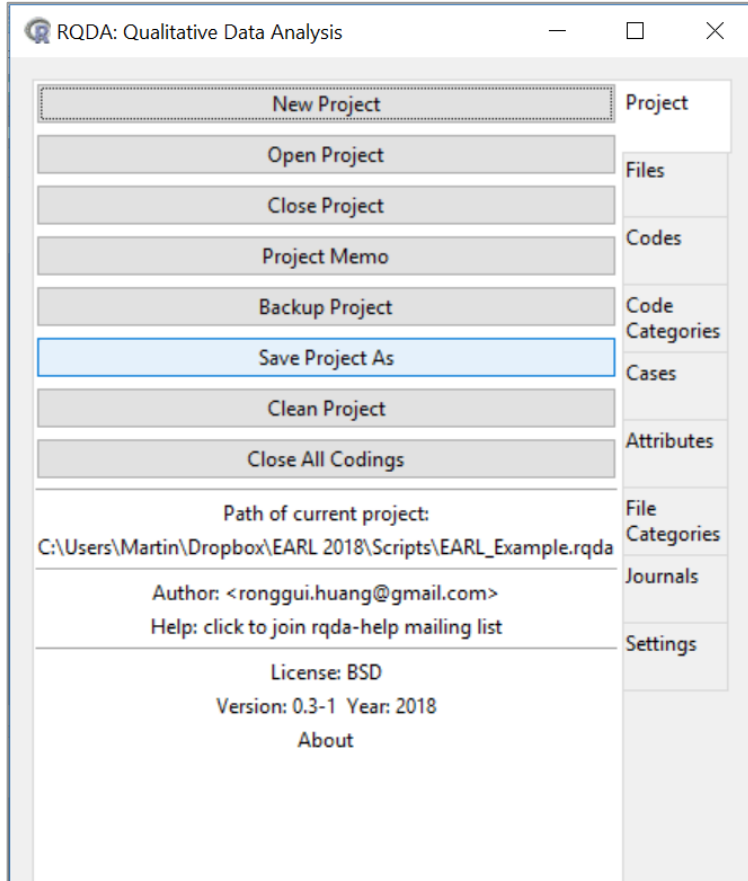
Outputs



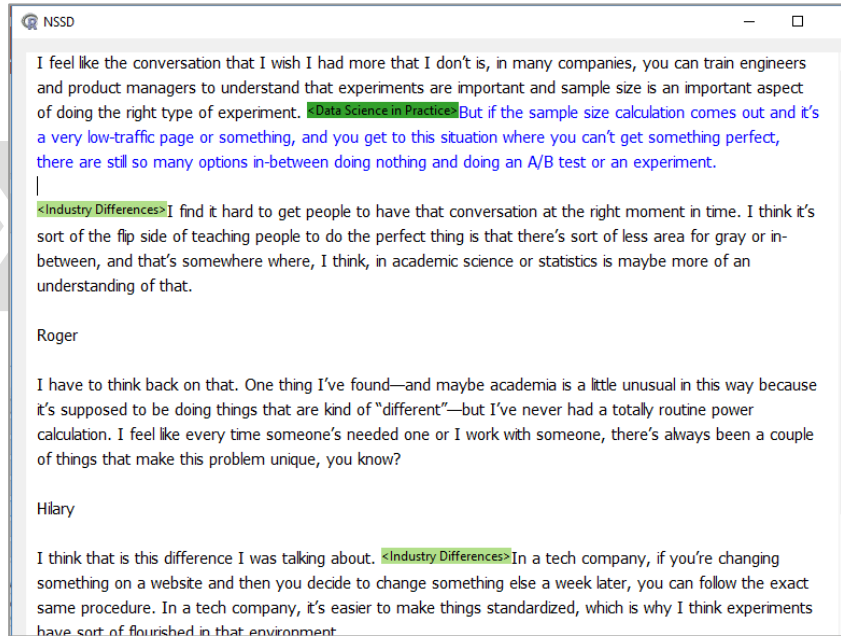
- ① Word / N-gram frequency analysis → **Help generate hypotheses**
- ② Word cloud → **Visual output for communicating strategic output**
- ③ Table output → **Ease of retrieving quotes to be used in a presentation**



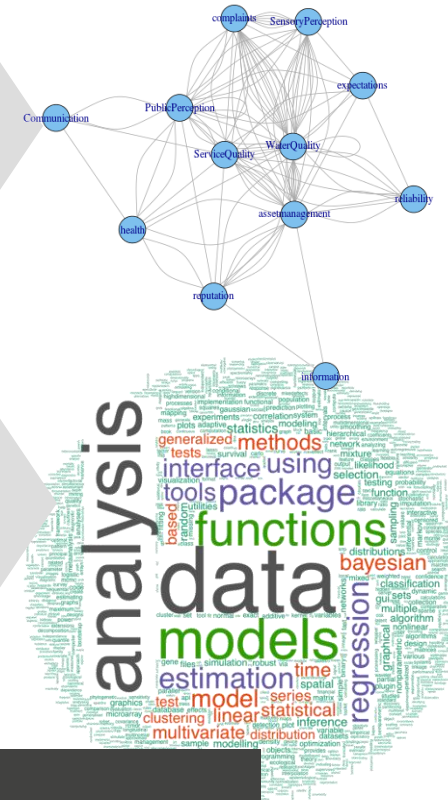
RQDA GUI Interface



Marked up transcript



Outputs



Structure of database

Code ID	Codes / Themes	File	...
1	Data Science in Practice	1, 2, 3	frequency of occurrence, code themes, etc.
2	Industry Differences	1, 2	
3	R versus Python	3,5	
...	



2 Nature of U&A data

U&A – research with the aim to understand a market and identify growth opportunities by answering questions on whom to target, with what, and how.

(Source: <https://www.ipsos.com/en/ipsos-encyclopedia-usage-attitude-surveys-ua>)



Features of a Usage & Attitude Survey (U&A)

Challenges of Using R

1

Wide data

Large number of variables with relatively few cases / observations

2

Design-dependent

Variables must be very closely interpreted relative to the original research design, e.g. questionnaire wording, rating vs ranking – making it difficult to interpret results within R



Example Question 1

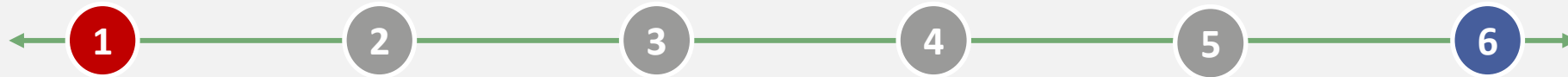
Attitudinal

Q25. For each of the following statements - on a scale of 1 to 6, please select the score which best describes your attitude towards using R packages:

I write my own R packages whenever it is possible to /do so

I look extensively for existing packages as solutions before considering to produce my own

Q25A



⋮

Q25P

I believe R packages should contain the minimum number of functions needed to solve a problem

I believe R packages should be comprehensive enough to cover the most likely tasks that I will do when solving a problem



Example Question 2

Usage

Q27. For each of the following tasks, please select the packages that you are likely to use in carrying out the task... **SELECT A MAXIMUM OF THE TEN MOST LIKELY PACKAGES**

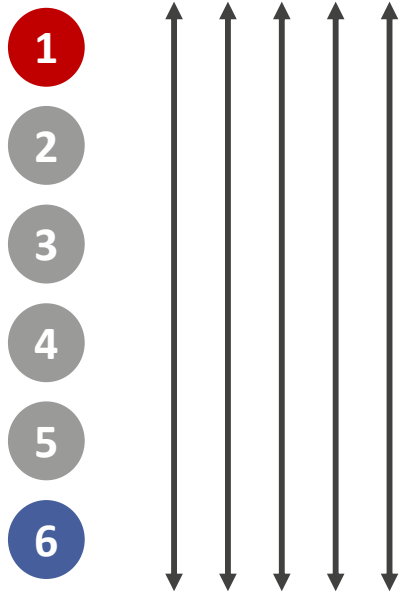
Tasks

1. Building a linear model with an interactive dashboard output
2. Writing a package to store customised analysis functions for a project
3. Producing publishing quality visualisations
- ...
12. Analyse survey data

Packages

1. dplyr
2. tidyr
3. data.table
4. tibble
5. shiny
6. plotly
7. ggplot2
- ...
110. Other (Please enter)

Entered through search box



Our challenge

1 The meaning represented by each variable is highly dependent on the **exact wording on both ends of the scale**

2 **Sparse data** (particularly in example 2)

3 Handling both the **variable and value labels** effectively through R when there are often a large number of these variables, which are not supported by data frames / tibbles

This explains the prevalence of **GUI analysis software** such as SPSS and Q

```


241 q27 1 21 q27_1_21
242 q27 2 21 q27_2_21
243 q27 3 21 q27_3_21
244 q27 4 21 q27_4_21
245 q27 5 21 q27_5_21
246 q27 6 21 q27_6_21
247 q27 7 21 q27_7_21
248 q27 8 21 q27_8_21
249 q27 9 21 q27_9_21
250 q27 10 21 q27_10_21
[ reached getoption("max.print") -- omitted 1070 rows ]

```

Number of variables from one question



The solution to overcoming this is to use a mixture of the following that helps us automate certain workflows:

-  dplyr
- for loops / apply functions
- merge functions
- User-defined functions

Generalised Approach



Create empty list



Loop* through scenarios / categories

Operations to wrangle / analyse data, e.g. `group_by()`, `summarise()`, `mutate()`, `spread()`, `gather()`

Output: either as a tibble / data.frame or **mschart** object

Assign output to member of list



→ Create a combined output using `merge_recurse` or `merge_all`

→ Export Excel output using `writexl::write_excel()`

→ Export PowerPoint output using `mschart` package

+ timestamp to outputs to document analysis

*Or in the form of combining `sapply` and a UDF, if possible – for computation efficiency; but for loops benefit from readability

Survey Analysis (With R)

Import



Import files as SAV (SPSS) or CSV files; sometimes serialised as RDS to save loading time and memory

Prepare



Extract, cleaning, and structuring data to prepare for analysis

Analyse



apply functions, loops, and custom functions for generating the outputs – usually in the form of summary tables

Lepton providing an alternative to R packages for managing ad-hoc functions

Run different iterations

Output



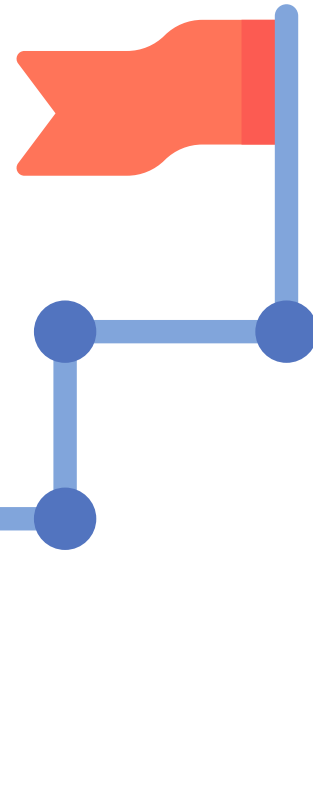
Summary tables, e.g. contingency tables



VBA for automating formatting



PowerPoint outputs – occasionally using SVG (vector) for visual quality



Outcome

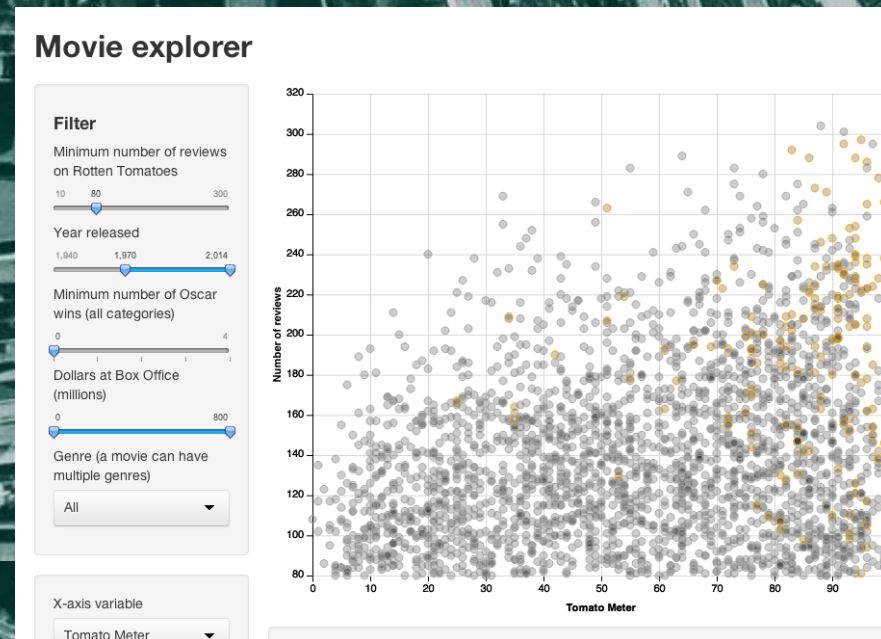
- **Reproducibility** – can trace back to exactly what you have done
- **A quick and organised way of exploring data through iterations** (using custom functions)



-
-
- 3** **Client requirement and expectations (of outputs)**



In many business environments there is still a preference of (perceived) **simple and reliable outputs** over interactive dashboard outputs like Shiny or Tableau





There are many reasons for this – sometimes it may be the **nature of the data** (small, ad-hoc datasets that require a lot of cleaning), but also driven by **certain specific needs**:

- 1** Clients can easily adapt content for their own use when sharing with internal stakeholders
- 2** No requirement to know any code (access to data not dependable on internal skill or supplier)
- 3** Perceived to have fewer 'moving parts', e.g. what do I when there is an error
- 4** Output is not merely data visualisation – often a requirement to heavily annotate or to construct a conceptual framework
- 5** Familiarity...



Correlations: Age (Indices) and intuitive - to investigate

Reach (as % of all 15+ adults) Results not entirely

Sample Weight	Reach (%)
London_15_24	0.1575727
London_25_34	0.3078648
London_35_44	-0.0950570
London_45_54	-0.2475772
London_55_64	-0.3399903
London_65_74	-0.1429540
London_75_84	0.3111817

Correlations: Age (Indices) and intuitive - to investigate

Results not entirely

Reach (as % of all 15+ adults)



Sample Weight	Reach (%)
London_65	-0.2855705
London_All_men 25_34 yrs	0.6457612

Familiarity (Style)

“Content looks great, could you put this in a format that looks less academic?”



We work with partners and clients from very varied backgrounds of technical knowledge and use different kinds of software.



For better or worse, Excel and PowerPoint are still the dominant vessels for communicating analysis and findings.

Thankfully, R has the versatility to accommodate these needs!

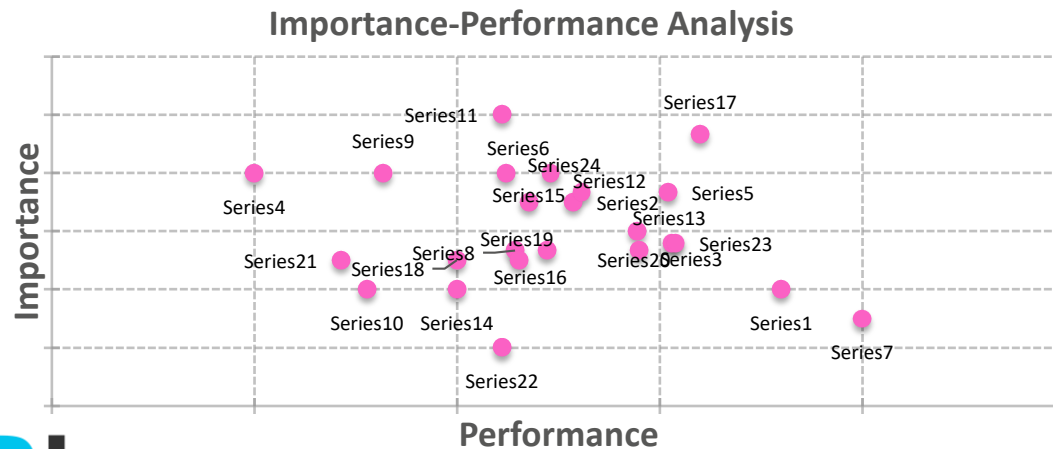


Automating PowerPoint charts with **officeR** and **mschart** packages (David Gohel)



Applications in automating Importance Performance Analysis in PowerPoint

- Scatterplots with variable labels



Generalised Approach



Read in a pre-set PPTX template and assign to an object



Create function to manipulate data, and return a **mschart object** using functions from the **mschart** package



Create function that adds the chart object to the PPTX object

Write PPTX object



1 Create a list of data frames for 'converting' into chart objects

2 Read in Template

3 Function for appending slides to PPT object (within R)

4 Loop sequence that creates formatted chart objects from 'output list'

5 Assigning them to 'output_combined_doc'

6

```
1 ##### Commence PowerPoint production #####
2 output_list <- list(NULL) # Initialise a list
3
4 segment_char <- c("Segment A","Segment B","Segment C")
5
6 for(i in 1:length(segment_char)){
7   str <- segment_char[[i]]
8
9   produce_tables(segment_char[[i]]) %>% # custom function for analysis / dplyr pipeline
10    # takes segment name (character) as argument
11    append_to_list(output_list,str) # assign table to list
12 }
13
14 source_doc <- read_pptx("Source/Source_PowerPoint.pptx") # Read in PowerPoint Template File
15 output_combined_doc <- source_doc # initialise an "output" document
16
17 # Function for adding newly created slides to PowerPoint file
18 gen_chart <- function(file,chart){
19   add_slide(file, layout="Content_1",master = "Empty slide with chart") %>%
20   ph_with_chart(chart=chart)
21 }
22
23 for(i in 1:length(output_list)){
24   segname <- names(output_list)[[i]]
25   title <- paste("Satisfaction Importance Matrix -",segname)
26
27   output_list[[i]] %>%
28   separate('.',sep="/",into="label",extra="drop") %>% # Text to columns abbreviating effect
29   drop_na() %>%
30   ms_scatterchart(x="satisfaction",y="Importance",group="label") %>%
31   chart_labels(title=title) %>%
32   chart_data_labels(position="b", show_legend_key = FALSE,show_serie_name = TRUE) %>%
33   chart_labels_text(values=fp_text(color="black",font.size=10,font.family="Calibri")) %>%
34   chart_data_fill(values="#4F81BD")%>%
35   chart_data_stroke(values="#4F81BD")%>%
36   chart_data_size(values=6) -> p_chart
37
38   output_combined_doc <- output_combined_doc %>%
39   gen_chart(p_chart)
40 }
41
42 print(output_combined_doc,
43       target=timed_fn("Output/Importance Satisfaction Slides",".pptx"))
```



```
1 Sub ScatterLabelsTweak()  
2  
3 Dim sld As Slide  
4 Dim shp As Shape  
5 Dim sr As Series  
6 Dim chrt As Chart  
7 Dim i, j, k, m As Long  
8  
9 r = 28.3464567 'r converts cm to points  
10  
11 For Each sld In ActivePresentation.Slides  
12     For Each shp In sld.Shapes  
13  
14         If shp.HasChart Then  
15             shp.Height = 12.83 * r  
16             shp.Width = 22.73 * r  
17             shp.Left = 1.34 * r  
18             shp.Top = 3.88 * r  
19  
20             shp.Chart.ChartTitle.Font.Size = 14  
21             shp.Chart.ChartTitle.Font.Name = "Calibri"  
22  
23             With shp.Chart.Axes(xlValue) 'Y-axis  
24                 .TickLabelPosition = xlNone  
25                 .AxisTitle.Font.Size = 14  
26                 .AxisTitle.Font.Name = "Calibri"  
27             End With  
28             With shp.Chart.Axes(xlCategory) 'X-axis  
29                 .TickLabelPosition = xlNone  
30                 .AxisTitle.Font.Size = 14  
31                 .AxisTitle.Font.Name = "Calibri"  
32             End With  
33
```

PowerPoint VBA

- Scatter plot style formatting
- Axes formatting
- Repositions and resizes charts
- Loops through entire PowerPoint document

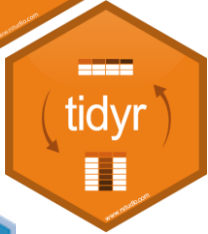
```
34  
35 shp.Chart.HasLegend = False  
36  
37 j = shp.Chart.FullSeriesCollection.Count  
38 Debug.Print j  
39 For i = 1 To j  
40     shp.Chart.FullSeriesCollection(i).Format.Fill.Visible = msoTrue  
41     shp.Chart.FullSeriesCollection(i).Format.Line.Visible = msoTrue  
42  
43     shp.Chart.FullSeriesCollection(i).Format.Fill.ForeColor.RGB = RGB(94, 98, 68)  
44     shp.Chart.FullSeriesCollection(i).Format.Line.ForeColor.RGB = RGB(94, 98, 68)  
45     shp.Chart.FullSeriesCollection(i).HasLeaderLines = True  
46     k = shp.Chart.SeriesCollection(i).Points.Count  
47     Debug.Print k  
48     For m = 1 To k  
49         shp.Chart.SeriesCollection(i).Points(m).DataLabel.Font.Size = 8  
50         shp.Chart.SeriesCollection(i).Points(m).DataLabel.Font.Name = "Calibri"  
51     Next  
52 End If  
53  
54 Next shp  
55 Next sld  
56  
57 End Sub
```



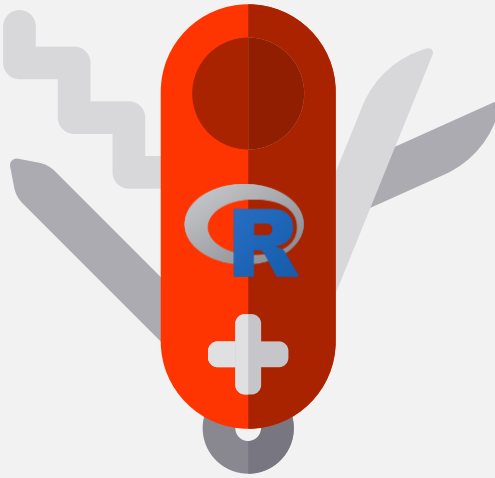
Summary



Data cleaning and manipulation that is readable and reproducible
Also easy to cache in functions – significant improvement to Excel



Extremely versatile tools to work with the kind of data we receive
Multi-level, loop-level, multiple choice type data



Sheer versatility and efficiency offered by user-defined functions, apply functions and loops

PowerPoint chart automation with **OfficerR** and **mschart**



(if necessary – publishing quality vector image outputs with SVG using ggplot2)

Making Qualitative Data Analysis systematic with **RQDA** – significant improvement on traditional analysis



Ability to automate or conduct reproducible survey analysis
Survey analysis with questionr or survey



Thank you!

martin.chan@rainmakerscsi.com

Rainmakers CSI Ltd
71 Gloucester Place
London W1U 8JW
T +44 (0) 20 3691 8100
E info@rainmakerscsi.com
rainmakerscsi.com