**Case study**

Official statistics in UK Government

# Official statistics in UK Government



In-house software

https://ukgovdatascience.github.io/rap_companion/why.html#the-current-statistics-production-process
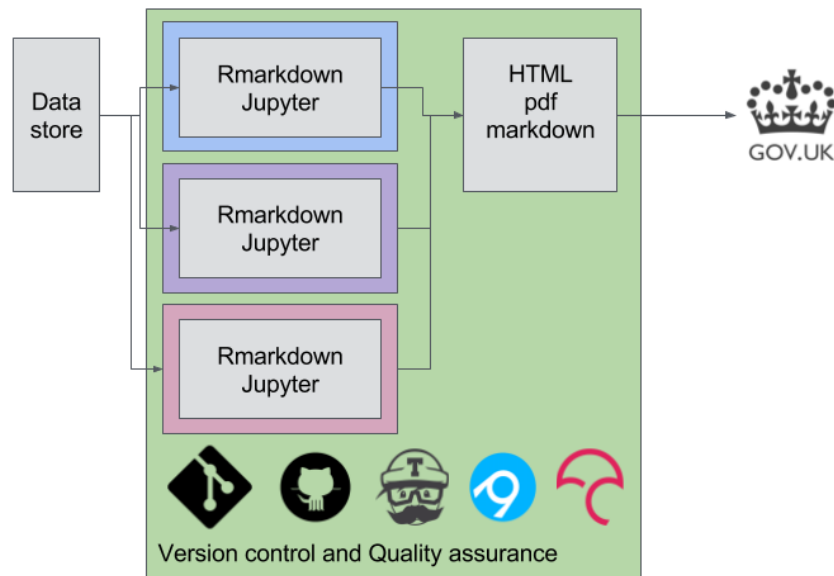
# Manual processes add risk

- Errors in spreadsheets are common

- Manual processes risk introducing human error

- Checking and peer review are not embedded in the process

- Challenging to reproduce previous work

Institute
and Faculty
of Actuaries

# Reproducible analytical pipelines



https://ukgovdatascience.github.io/rap_companion/why.html#desired-reproducible-analytical-pipeline

Institute
and Faculty
of Actuaries

**"** The potential **time savings** for analysts are enormous, freeing them up to focus on the interpretation of the results. The other huge benefit comes from building a process that is fully **transparent, auditable and verifiable** – reducing risk and improving quality. **"**

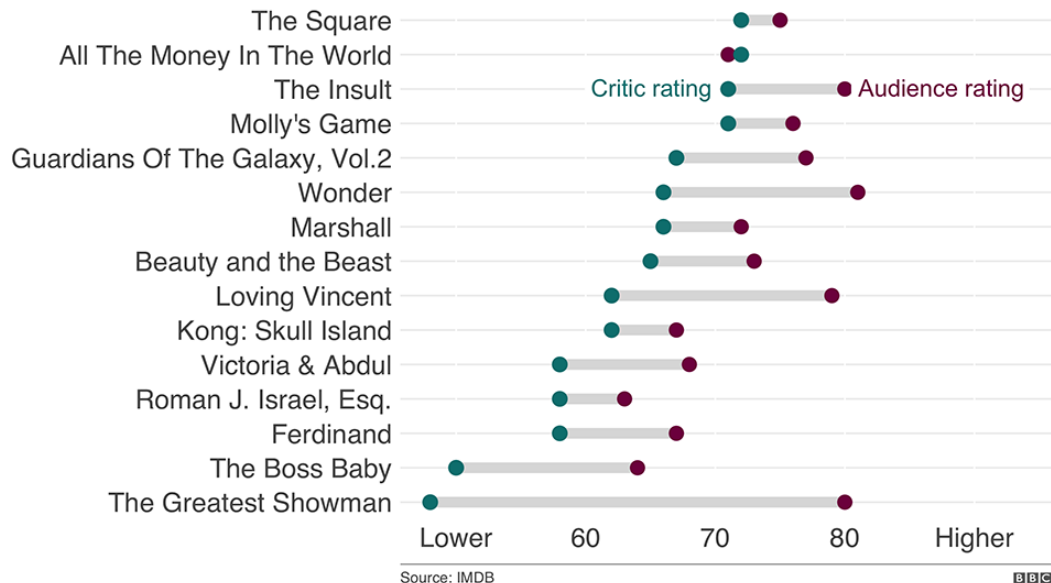Matt Upson and Mat Gregory, Government Digital Service

Institute
and Faculty
of Actuaries

# BBC News website graphics



**How critics and filmgoers disagree**
Difference in average score from critics and audience for 2017's Oscar-nominated films

The Square
All The Money In The World
The Insult — Critic rating / Audience rating
Molly's Game
Guardians Of The Galaxy, Vol.2
Wonder
Marshall
Beauty and the Beast
Loving Vincent
Kong: Skull Island
Victoria & Abdul
Roman J. Israel, Esq.
Ferdinand
The Boss Baby
The Greatest Showman

Lower    60    70    80    Higher

Source: IMDB

https://www.bbc.co.uk/news/entertainment-arts-43146027 (edited to fit)

Institute and Faculty of Actuaries

**"** [This approach] **saves a huge amount of time and effort**, in particular when working with data that needs updating regularly, with **reproducibility** a key requirement of our workflow.  In short, it was a game changer… **"**

BBC Visual and Data Journalism team

Institute and Faculty of Actuaries

# Reproducible work

**Reproducibility** is the process of making code and data available so that others can easily replicate, verify and build on your analysis

Institute and Faculty of Actuaries

# Building blocks of a reproducible workflow

Data

Analytic code and automated checks

Documentation

Computational environment

Packaged in a standard way

Institute
and Faculty
of Actuaries

# Why is this important for actuaries?

- Enables more efficient working

- Allows analysts to focus on the bigger picture

- Easier collaboration

- Helps meet compliance requirements – internal and TAS

- A step towards automation

# Reproducible actuarial pipeline
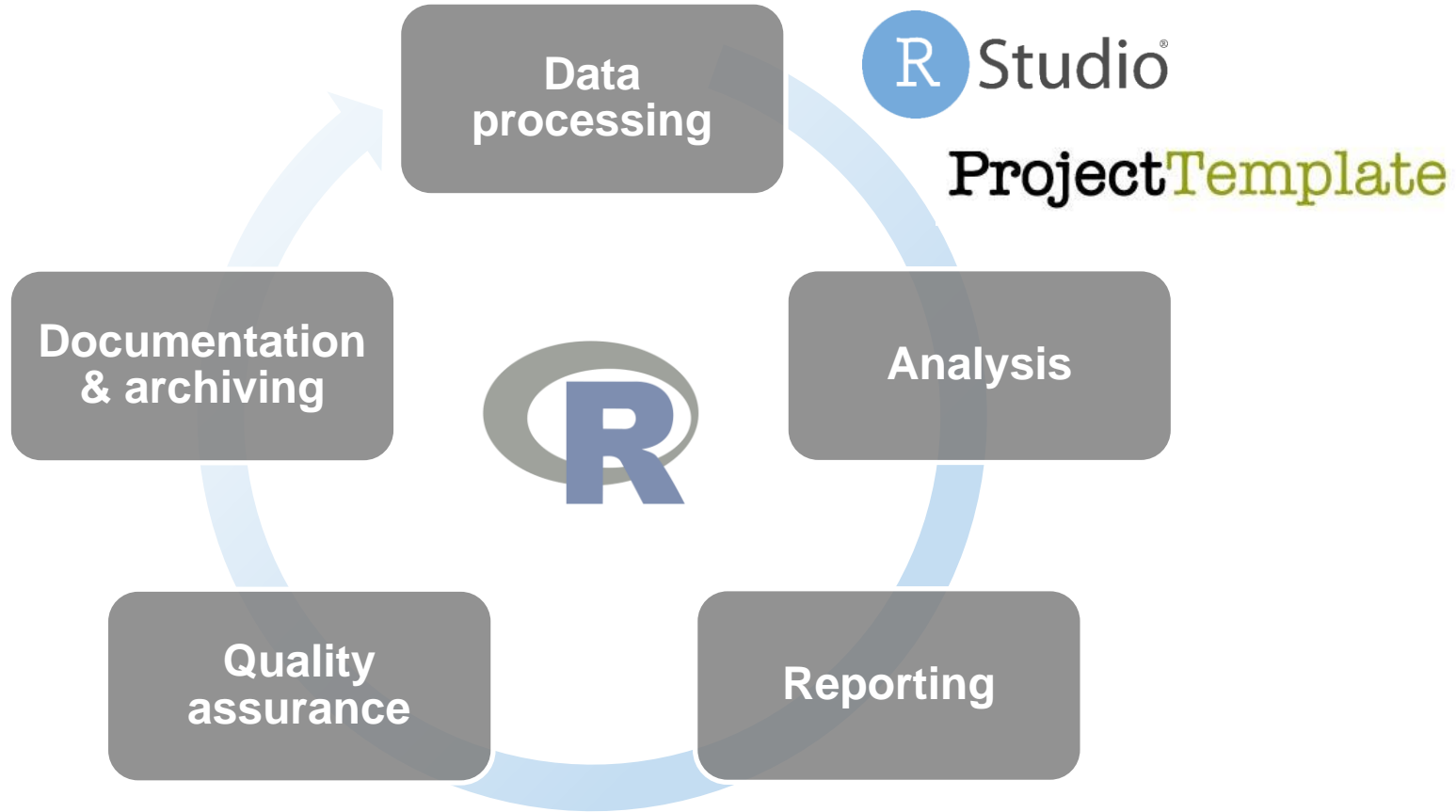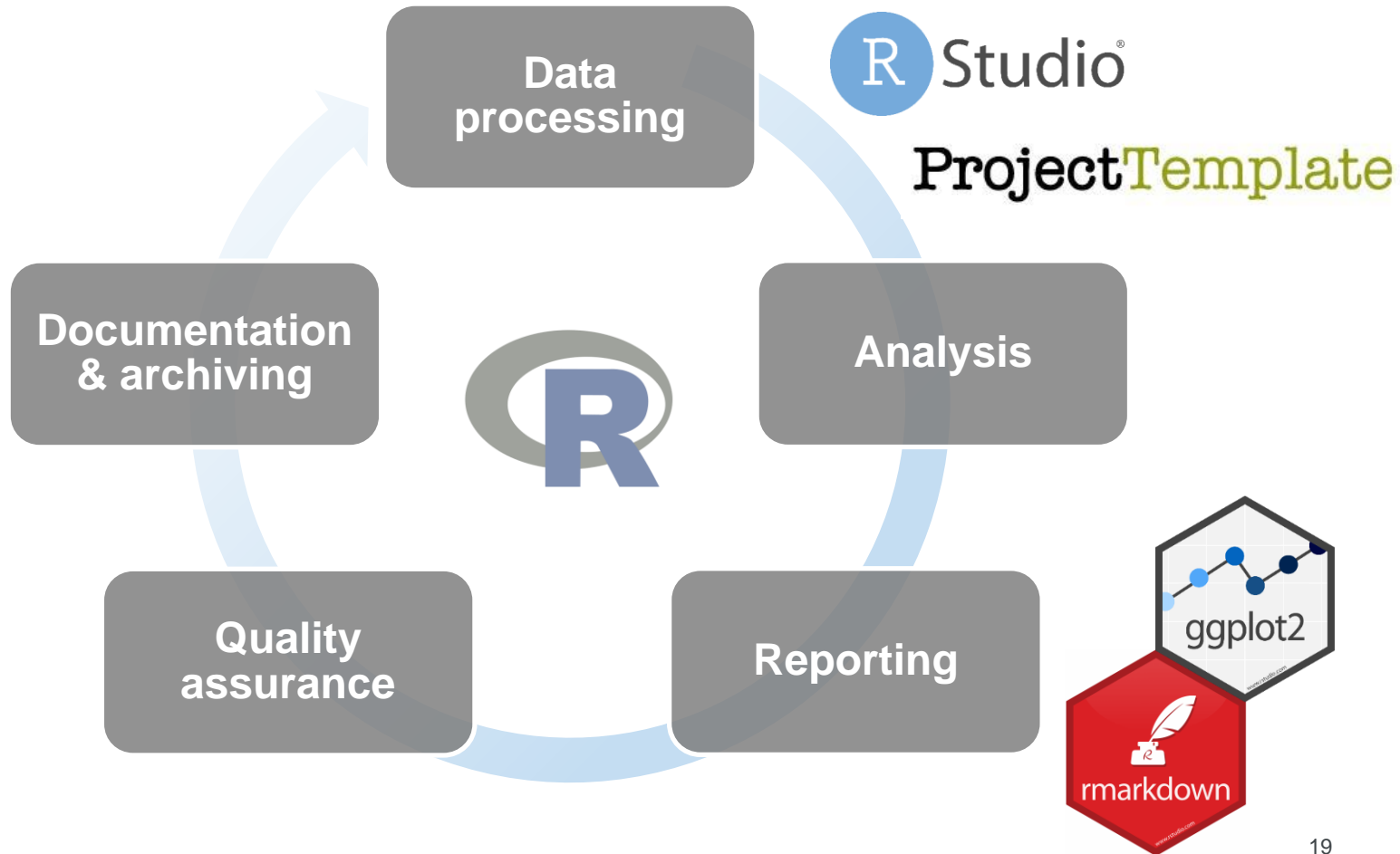
An example using R

# Why R?

- Stable, up-to-date and free

- Open source with an active support community

- Well suited to building reproducible pipelines and reporting

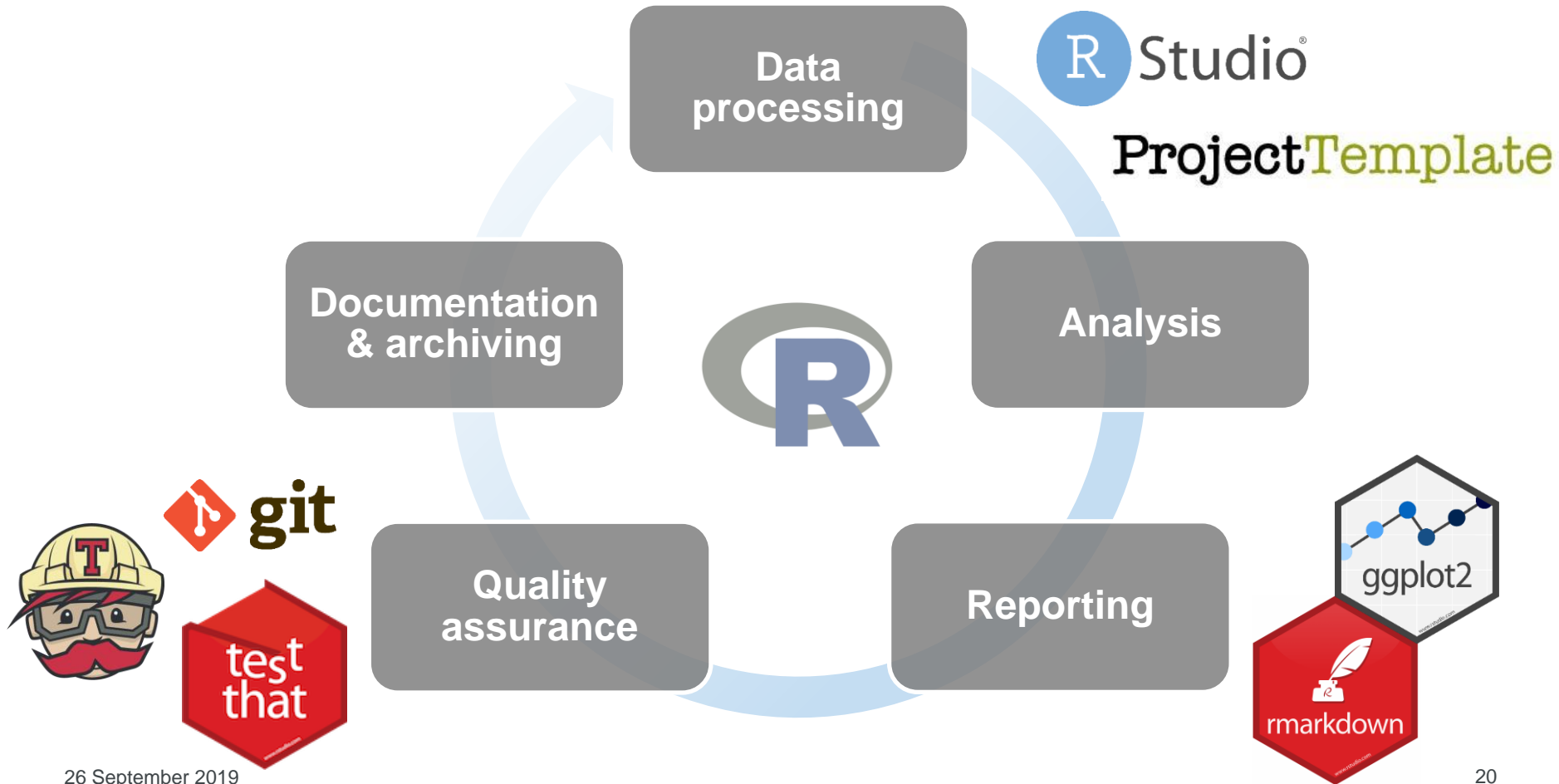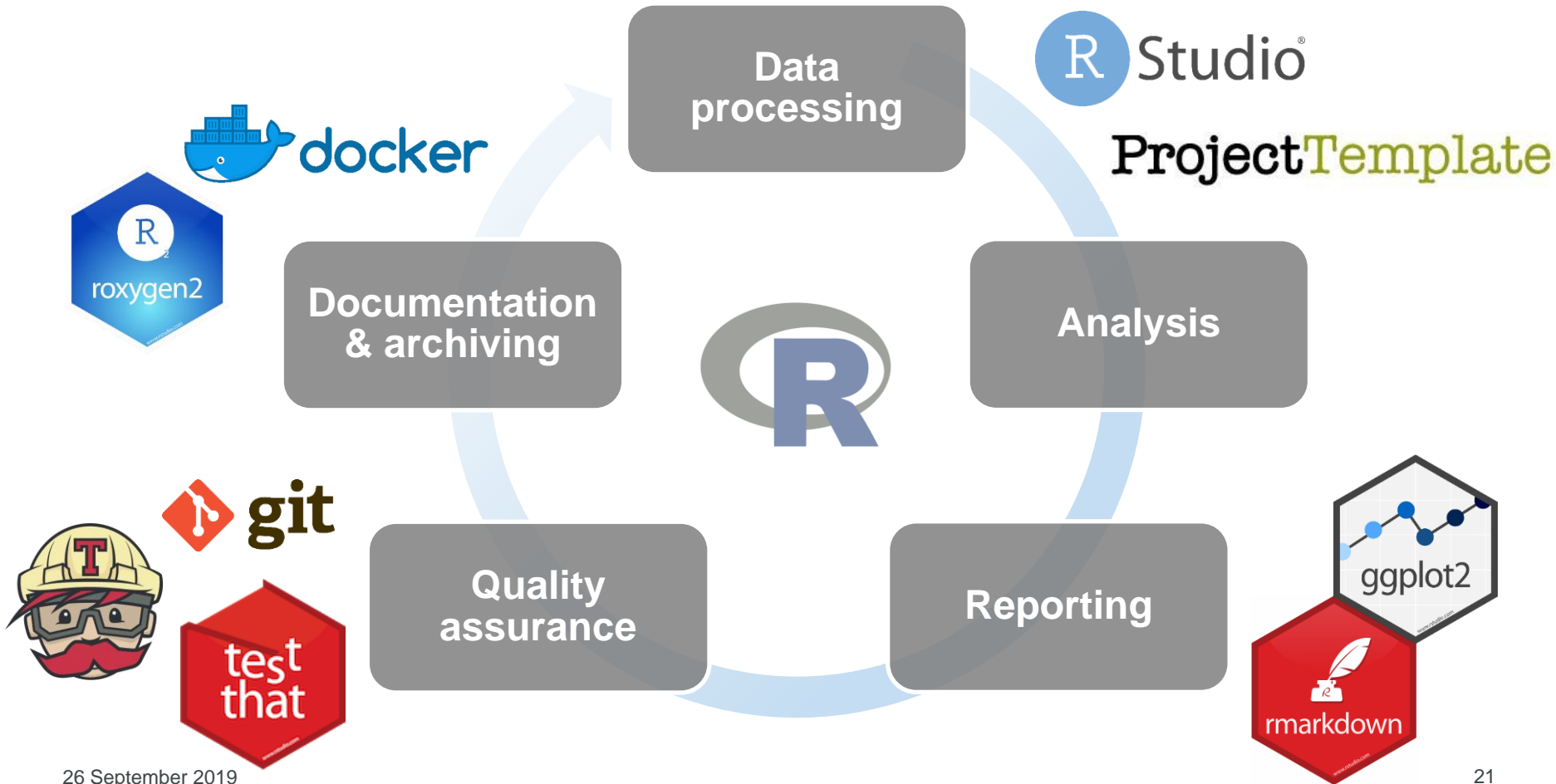- Next generation of actuaries will learn R under the 2019 curriculum

Institute
and Faculty
of Actuaries

> " Where does that figure come from? "

RStudio Source Editor

● ● ●  report.Rmd* ✕

Knit ▾    Insert ▾    Run ▾

```
1  ---
2  title: "`r params$title`"
3  author: "`r params$org`"
4  date: "`r params$date`"
5  output:
6    pdf_document:
7      latex_engine: xelatex
8  geometry: "a4paper, margin=2.5cm"
9  fontsize: 12pt
10 params:
11   title: "Cashflow analysis"
12   org: "Organisation name"
13   author: "Your name FIA"
14   date: "September 2019"
15 ---
16
17 ```{r setup, include=FALSE}
18 # Knitr setup
19 knitr::opts_chunk$set(echo = TRUE)
20 knitr::opts_knit$set(root.dir= normalizePath('..'))
21 ```
22
23 ```{r include=FALSE}
24 # Load project
25 library("ProjectTemplate")
26 load.project()
27
28 # Run code in analysis.R script
29 source("src/analysis.R")
30 ```
31
32 ## Summary
33
34 Using a discount rate of `r format(100*disc)`% p.a. the present value of the projected cashflows is &pound;`r format(round(pv_central,-3), big.mark=",")`.
35
36 95% of model outcomes have a present value in the range &pound;`r format(round(pv_lower,-3), big.mark=",")` to &pound;`r format(round(pv_upper,-3), big.mark=",")`.
37
38 ## Cashflow analysis
39
40 Lorem ipsum dolor sit amet, ad mea sumo vocibus graecis, at mea soleat doctus, usu elit dicta ne. Aliquid
```

Cashflow analysis

*Organisation name*

*September 2019*

a. the present value of the projected cashflows is £8,851,000.

a present value in the range £8,696,000 to £9,007,000.

ad mea sumo vocibus graecis, at mea soleat doctus, usu elit dicta ne, aperiri deniebas quo no. Vel ea assueverit disputando, qui.

usto ridens oporteat ea vim, noster minimum reformidans. sed dicant audiam explicari, vim at viris libris mnesarchum. comsan officiis, habemus accusata periculis an eam. Sapite dissentias temporibus.

v projection

> " Update the report using 2.75% "
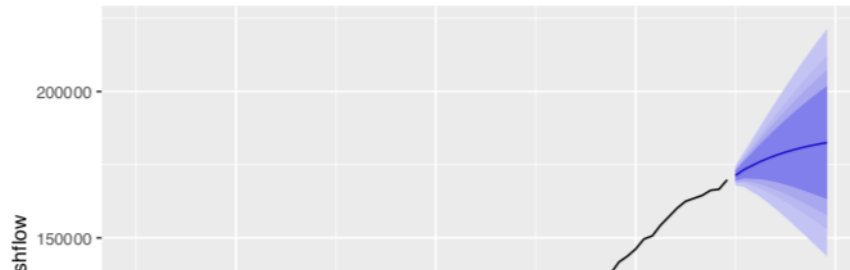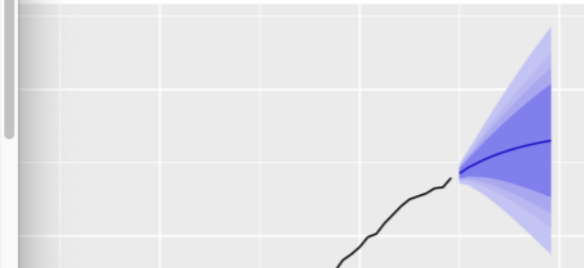
# Cashflow analysis

*Organisation name*

*September 2019*

## Summary

Using a discount rate of 3% p.a. the present value of the p...
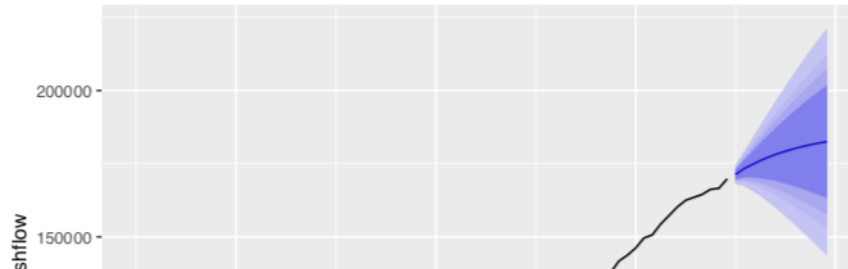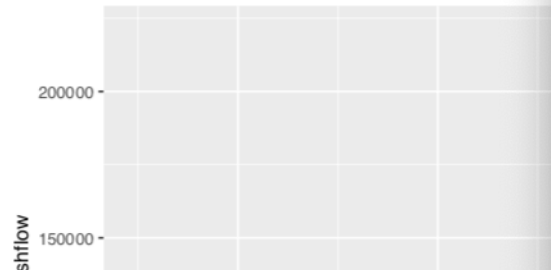
95% of model outcomes have a present value in the range...

## Cashflow analysis

Lorem ipsum dolor sit amet, ad mea sumo vocibus graeci...
dicta ne. Aliquid salutatus vix et, aperiri definiebas quo r...
Dicat alterum posidonium te qui.

Vim vocibus assueverit in, iusto ridens oporteat ea vi...
Scribentur mediocritatem, cu sed dicant audiam explicari...
Id quod consul est. Eu mei accumsan officiis, habemus a...
entem definitiones ut mel, pri te dissentias temporibus.

12 month cashflow projection

200000

shflow
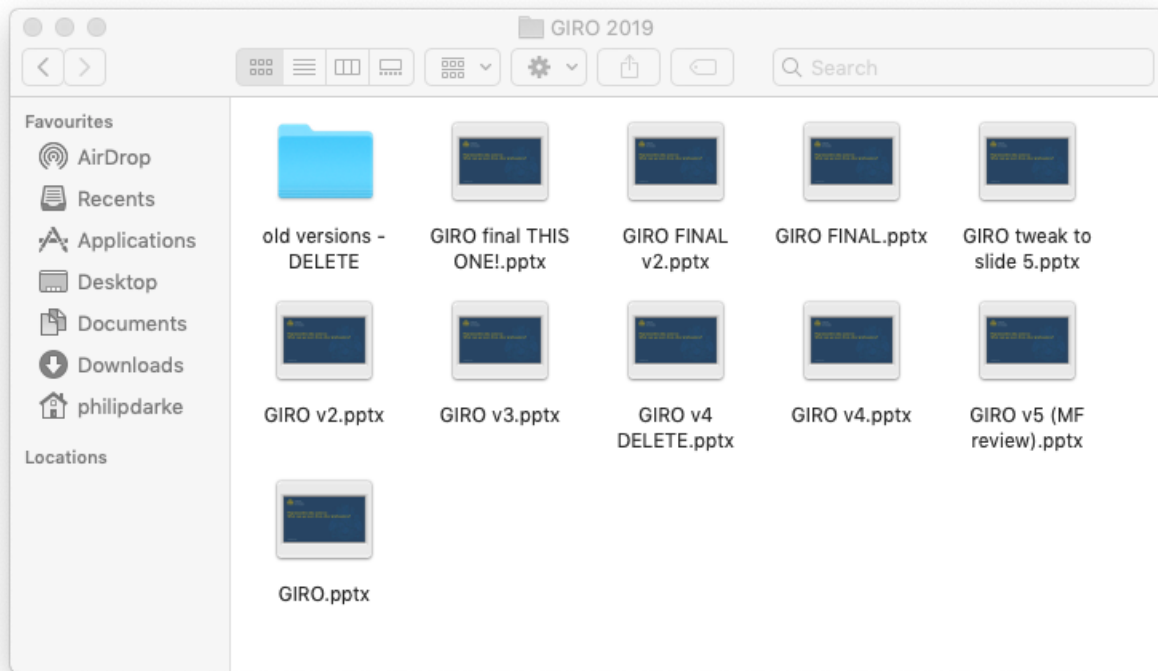
150000

---

RStudio Source Editor

analysis.R

Source on Save    Run    Source

```r
1   # Load project
2   library("ProjectTemplate")
3   load.project()
4
5   # Model cashflows as an ARIMA(2,1,0) time series
6   cashflow_model <- arima(cashflows, order = c(2,1,0))
7
8   # Create a 12 month forecast
9   forecast <- forecast(cashflow_model, 12, level = c(80, 90, 95, 99))
10
11  # Plot the forecast
12  cf_plot <- autoplot(forecast) +
13    xlab("Year") +
14    ylab("Cashflow") +
15    ggtitle("12 month cashflow projection")
16
17  # Hold cashflow forecasts in a data frame
18  forecasts <- data.frame(lower = c(cashflows, forecast$lower[,3]),
19                          central = c(cashflows, forecast$mean),
20                          upper = c(cashflows, forecast$upper[,3]))
21
22  # Set discount rate
23  disc <- 0.03
24
25  # Discount cashflows
26  pv_lower <- discount(forecasts[["lower"]], disc, 12)
27  pv_central <- discount(forecasts[["central"]], disc, 12)
28  pv_upper <- discount(forecasts[["upper"]], disc, 12)
29
```

# Collaboration and keeping an audit trail

# Collaboration and keeping an audit trail

# Challenges

- Relies on open source software

- Timing consuming to set up

- Training requirements

Institute
and Faculty
of Actuaries

Build a simple reproducible pipeline at

[philipdarke.com/reproducible-actuarial-work](philipdarke.com/reproducible-actuarial-work)

Institute
and Faculty
of Actuaries

# Applying these techniques in your work

- Take an existing process

- Develop a minimal viable solution (see the exercises)

- Pilot it and let others contribute

- Share what you learn

Institute
and Faculty
of Actuaries

# Questions

# Comments

The views expressed in this presentation are those of invited contributors and not necessarily those of the IFoA. The IFoA do not endorse any of the views stated, nor any claims or representations made in this presentation and accept no responsibility or liability to any person for loss or damage suffered as a consequence of their placing reliance upon any view, claim or representation made in this presentation.

The information and expressions of opinion contained in this publication are not intended to be a comprehensive study, nor to provide actuarial advice or advice of any nature and should not be treated as a substitute for specific advice concerning individual situations. On no account may any part of this presentation be reproduced without the written permission of the authors.

Institute
and Faculty
of Actuaries

# Useful tools for building a reproducible workflow

**RStudio** is a free and widely used development environment for R that integrates with the tools below.

**ProjectTemplate** automates the menial parts of statistical analysis and provides a standard way of working in R.

**R Markdown** is a notebook interface that allows code to sit alongside narrative text and can be used for reporting as part of a reproducible framework with **ggplot2** for creating charts and visualisations.

**Git** is a version control system for managing code and audit trails – it can be used privately in an organisation or with a web-based service such as **GitHub**.

**testthat** is a formal automated testing ("unit testing") package for R.

**TravisCI** integrates with GitHub to automatically run your tests when code is updated.

**roxygen2** automates the production of documentation for your code in R.

**Docker** packages dependencies inside a container which can run consistently on any infrastructure (also see **checkpoint**/**packrat** or consider creating a R package).

# References and resources

- RAP companion https://ukgovdatascience.github.io/rap-website/

- RAP Udemy video course https://www.udemy.com/course/reproducible-analytical-pipelines/

- Blog post on the use of R at the BBC https://medium.com/bbc-visual-and-data-journalism/how-the-bbc-visual-and-data-journalism-team-works-with-graphics-in-r-ed0b35693535

- Accompanying exercises https://philipdarke.com/reproducible-actuarial-work/

- Icons made by Smashicons and Dimitry Miroliubov from www.flaticon.com

Institute
and Faculty
of Actuaries

# Get in touch

**Dr Matthew Forshaw** is a Lecturer in Data Science at Newcastle University, and Data Skills Policy Leader at The Alan Turing Institute working on the Data Skills Taskforce. He is the Programme Director of Newcastle's Industrial MSc in Data Science.
mattforshaw.com

**Philip Darke** is an actuary with over 10 years' consulting experience at Mercer and a PhD researcher in data science at the EPSRC Centre for Doctoral Training in Cloud Computing for Big Data at Newcastle University.
philipdarke.com

Institute
and Faculty
of Actuaries