# Reproducible Analytical Pipelines & their value in fundamental Data Science

**Jeroen Minderman**

Senior Data Scientist

**23 January 2024**

# Fundamental needs for Data Science in NSOs

1. **Skills**
   - Data literacy
   - Programming literacy
   - Following/building Good Practice

2. **Buy-in**
3. **Resource** ( ~ buy-in!)

Data Science Campus

# RAP in Data Science: efficiencies

- Reproducible Analytical Pipelines (RAPs)
  - Formalised process for automation of analyses
  - Minimise manual steps, maximise transparency & reproducibility

- Improving *quality*, *trust*, *business continuity*
- Create **efficiency**: saving resource

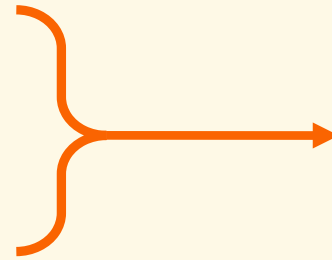https://analysisfunction.civilservice.gov.uk/support/reproducible-analytical-pipelines/

Data Science Campus

# RAP in Data Science: a blueprint for skills

- Data Science ≠ RAP; RAP alone ≠ Data Science
- BUT; for RAP we need, e.g:
  - Understanding process/scope
  - Programming skills; R/Python
  - Focus on application & impact

# RAP in Data Science: a blueprint for skills

- Data Science ≠ RAP; RAP alone ≠ Data Science

- BUT; for RAP we need, e.g:
  - Understanding process/scope
  - Programming skills; R/Python
  - Focus on application & impact

== **Fundamentals of data science**

+

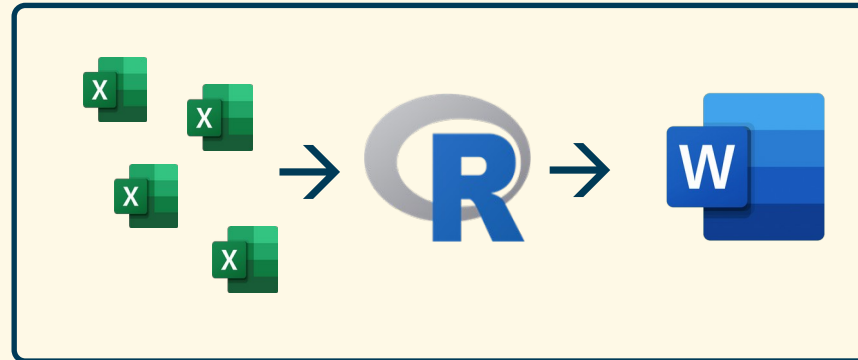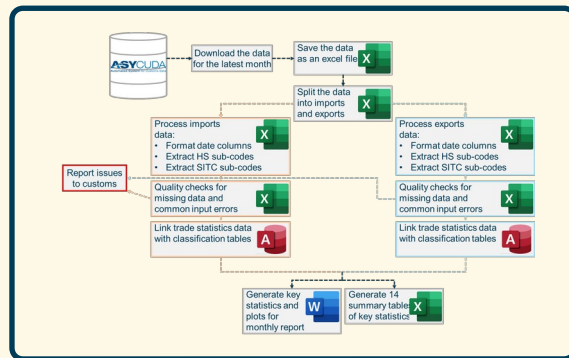efficiency gains

=

**win-win**

# RAP mentoring approach

- Data scientist(s) mentoring small groups in partner NSO's, e.g.
  - Scoping suitable work
  - Flexible & scalable training
  - Support pipeline development

- Longer period with regular check-ins
- ➔ Focus on *application* and *impact*

# Palestine Central Bureau of Statistics (PCBS)

!! Manual & labour intensive process for trade statistics data
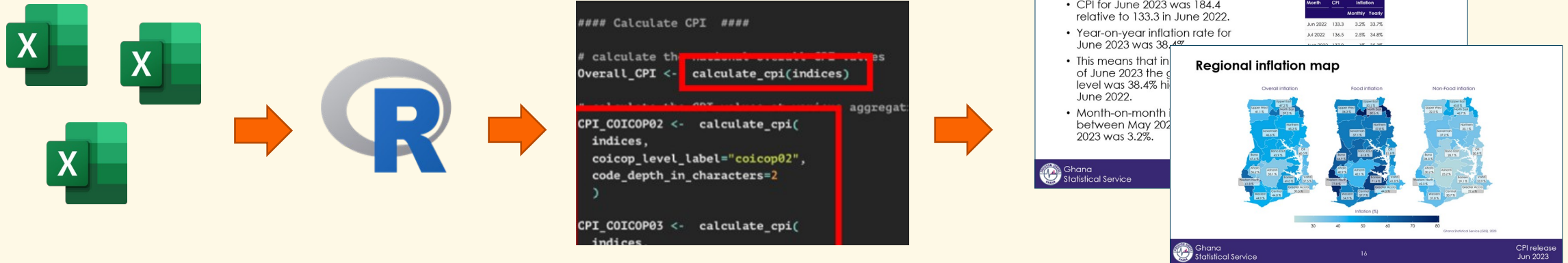
!! Some R training but not applied



✓ Confidence in R, version control → scope for wider DS work

✓ Start on RAP → increased quality & transparency

✓ Wider benefits: support overall process improvements

# Ghana Statistics Service (GSS)

‼ Skilled in R, part automated process for CPI production

‼ Monthly reporting was done manually



✓ Reproducible report generation → transparency and reliability

✓ Built further confidence in R → mentees now supporting others

## Data Science Campus

# RDSA ONS Nature of Crime Automation project

!! Complex and time-consuming/error prone existing process



✓ Massive efficiency gains & increased quality / reliability

✓ Working within organisation → able to advise on standardisation of process & reporting

# RAP strategy

## Standards provide framework for capability & Good Practice

# Summary, suggestions & discussion points

- RAP = efficiencies… but also blueprint for DS skills
  - Mentoring is efficient & scalable means to build both
- "Stepping stone" to Pillars 2 and 3

- Notes –
  - Some *initial* skills beneficial; e.g. precede with training courses?
  - Mentor & mentee ***availability is crucial*** (**buy-in**)
    - E.g. ring-fence part of staff time, but plan continued development
  - Take the long view: initial resource cost →→ efficiencies

Data Science Campus

# **Resources / links**

- ONS [Data Science Campus](#)
- [UK Analysis Function RAP](#)
- [RAP Strategy](#) (UK Analysis Function / ONS)
- RAP [case studies](#)
- [Using RAP to improve statistics](#)
- [Quality Assurance of code for analysis and research](#)

This guidance describes software engineering good practices that are tailored to those working with data using code. It is designed for those who would like to quality assure their code and increase the reproducibility of their analyses. Software that apply these practices are referred to as reproducible analytical pipelines (RAP).

**Data Science Campus**