

## **Training research output checkers**

Felix Ritchie (University of the West of England, United Kingdom)

*Felix.Ritchie@uwe.ac.uk*

### ***Abstract and Paper***

As National Statistics Institutes (NSIs) support research use of their microdata, checking those research outputs for disclosure risk becomes an important part of the NSI's risk management. Where the data are distributed, NSIs rely upon giving researchers guidance on how to reduce risk by not producing tables with small numbers, for example (see Eurostat, 2016 for an example of researcher guidance).

For more sensitive microdata, NSIs increasingly provide access through controlled facilities, on-site or remote. In such facilities data is usually very detailed, with minimal de-identification carried out. Hence, most NSIs operate a system of manually checking research outputs for disclosure risk.

Checking regular NSI outputs is a skilled task, but there is much training available and a large literature. However, checking research outputs is more complicated. Statistically, • researchers produce a much wider range of statistics • researchers are less concerned with consistency of outputs • the manipulation of data by researchers makes outputs inherently less sensitive

and hence research outputs are lower risk. On the other hand, • researchers are less likely to have had extensive training in disclosure control • researchers are less likely to share the NSIs perspective on risk avoidance • researchers may view NSI attempts to check output as blocking research rather than supporting it

The people element is less important when operating a strictly rules-based output checking model. In practice, very few NSIs are strictly rules-based. Most are notionally rules-based but with some form of 'flexibility' allowed (the 'ad hoc' model); a smaller proportion follow the principles-based approach to output disclosure control. For all flexible systems, getting buy-in from researchers makes a substantial difference to the operation of the facility. As a result, checking the outputs of researchers requires both statistical skills, personal management abilities, and a good awareness of institutional context. This is a very different skillset from the checkers of NSI outputs.

However, there is little training for or guidance to support the research output checker. Some statistical elements are covered in Brandt et al (2010, and its various updates and incorporations into other texts) and Greci et al (2019); Ritchie and Welpton (2015) focus on procedures and on people management. Eurostat offers a course in output checking for research [check content]. But generally research output checkers are trained on a 'grandfather' basis: learning on the job via more experienced colleagues.

The grandfathering approach has many advantages, not the least of which area a focus on practical skills and the strong institutional awareness. The disadvantages are that the more experienced colleagues may not be available, that there are fewer checks against poor institutional practices developing, and that there are fewer opportunities to learn from others. In addition, only a small



number of individuals in a very small number of organisations are likely to be output checkers, and the isolation can lead to this being seen as a backwater in one's career.

These have all been issues in the UK in recent years, and the potential risks are likely to expand as the number of research centres checking their own outputs is expected to increase rapidly. Hence, in 2019 ONS commissioned the University of the West of England (UWE) to review research output checker training and propose ways to develop both training and the output checking community.

The UWE team proposed a solution based upon three elements • A new training program focused equally on statistical skills, decision-making under uncertainty, and people management • Using the training course to build up an output-checking community • Integrating and developing new literature specifically for this new course

The new courses are due to start in May 2019. By the end of Summer, we expect the pilots have been completed and the user community and websites to be up and running. In this presentation we will report on the experience, feedback, and lessons learned, and look forward to comments from colleagues in other NSIs.

# Training research output checkers

Felix Ritchie\*

\* University of the West of England, Bristol. [Felix.ritchie@uwe.ac.uk](mailto:Felix.ritchie@uwe.ac.uk)

**Abstract:** National statistics institutes (NSIs) and other organisations increasingly allow users access to confidential microdata within strictly controlled environments. These are then checked for disclosure risk on leaving the environment by the researchers and/or by NSI staff.

Statistical disclosure control (SDC) for regular NSI outputs is well understood; focused on tabulations of high-dimensional tables, there are training courses and much supporting material available to staff. Checking research outputs is a rather different task. These outputs are likely to be unstructured, complex and using transformed data; subjective judgement is crucial, and hence assessment requires a different set of skills. As research output checking has only developed over the last ten years or so, there is relatively little support for staff or formal training.

In 2019 the UK Office for National Statistics (ONS) commissioned its first course in research output checking. The course follows the ‘user centred’ approach of other UK confidentiality training, and places equal emphasis on statistical ability, subjective decision-making and people management. This paper describes the objectives, structure and ethos of the course, and initial feedback from the pilots.

## 1 Introduction

There is an insatiable demand for access to microdata for research and analysis. For government organisations, particularly national statistical institutes (NSIs), this demand is increasingly met by providing ‘secure use files’ (SecUFs): that is, the data holder allows researchers to use very sensitive data but only in an environment under the data holder’s control. For most organisations these days, this means on-site or remote access to a ‘research data centre’ (RDC), where the researchers have almost complete freedom to carry out analysis as detailed in the approved project scope, but with entry and exit from the facility very strictly controlled.

In such facilities data is usually very detailed, often with only direct identifiers removed. Hence, checking outputs for residual risk is an important part of the data holder’s risk management portfolio. Because of the nature of research outputs, this ‘statistical disclosure control’ (SDC) needs to be done manually in most cases.

SDC is also applied to the official statistics produced by NSIs and other government organisations. Official statistics are typically high-dimensional linear combinations;

that is, tabulations, indexes, means etc which are broken into many different categories. SDC checking of official statistics is a skilled task, but there is much training available. There is also a large literature stretching back forty years, and relatively little controversy about the rules to be applied.

Research output checking requires a different skill set to SDC for official statistics. Research outputs have a wider scope, are more likely to transform data as part of the analysis, and select subsets of data based on subjective choices. Judgement, experience and other personal qualities are therefore likely to be more important than any specific rule for clearing outputs.

However, there is little training for, or guidance to, support the research output checker. Some statistical elements are covered in guidelines for research centre managers (Brandt et al 2010; Bond et al, 2015; Greci et al, 2019). Ritchie and Welpton (2015) focus on procedures and on people management. Eurostat offers a course in output checking for research<sup>1</sup>. But generally research output checkers are trained on a ‘grandfather’ basis: learning on the job via more experienced colleagues.

The grandfathering approach has many advantages, not the least of which offer? a focus on practical skills and the strong institutional awareness. The disadvantages are that the more experienced colleagues may not be available, that there are fewer checks against poor institutional practices developing, and that there are fewer opportunities to learn from others. In addition, only a small number of individuals in a very small number of organisations are likely to be output checkers, and remaining in such a niche role can be seen as limiting in terms of career development.

These have all been challenges in the UK in recent years when making microdata available for research, and are likely to expand as the number of research centres checking their own outputs is expected to increase rapidly. Hence, in 2019 ONS commissioned the University of the West of England (UWE) to review research output checker training and propose ways to develop both training and the output checking community.

The UWE team proposed a solution based upon three elements

- A new training program focused equally on statistical skills, decision-making under uncertainty, and people management
- Using the training course to build up an output-checking community
- Integrating and developing new literature specifically for this new course

The new course has been piloted over Summer 2019 and is now being offered to all UK data centres. This paper reports on the aims and objectives of the course, its ethos and structure, and lessons learned from the pilots.

---

<sup>1</sup> Catalogue: <https://ec.europa.eu/eurostat/documents/747709/6103606/2019-ESTP-catalogue-final.pdf>

The next section describes the problem of research output checking. Section 3 covers the course aims and pedagogical approach. Section 4 reviews the initial pilots. Section 5 concludes. For simplicity we shall assume throughout that the data holder is an NSI.

## **2 Research output checking**

SDC for research outputs is complicated. Statistically,

- researchers produce a much wider range of statistics, many of which have no inherent disclosure risk
- researchers are less concerned with consistency of outputs and so disclosure by differencing is less likely
- the transformation and manipulation of data by researchers makes outputs inherently less sensitive

and hence research outputs are lower risk. This also makes the assessment of risk in research outputs much more subjective. Checking tabulations for an official statistic is largely a matter of identifying the target population, the sample, unique values, and the potential for differencing. In contrast, checking a tabulation produced by a researcher may involve substantial effort to understand how the researcher arrive at the table – and this is assuming that the researcher has labelled or described the output correctly and usefully.

Moreover,

- researchers are less likely to have had extensive training in disclosure control
- researchers are less likely to share the NSIs perspective on risk avoidance
- researchers may view NSI attempts to check output as blocking research rather than supporting it

Output checking processes are typically ‘rules-based’ or ‘principles-based’ (Lowthian and Ritchie, 2017). In the former, clearance is a yes/no process, and the skill of the checker comes from being able to understand quickly and clearly which rules to apply. In the principles based approach, there are no absolute rules: checkers apply a strict set of rules in most cases, but allow researchers to ask for exceptions (ONS, 2019). This requires more judgment from the checker.

In practice, very few NSIs are strictly rules-based (Alves and Ritchie, 2019). Most are notionally rules-based but with some form of ‘flexibility’ allowed (the ‘ad hoc’ model); a smaller proportion follow the principles-based approach to output disclosure control. For all flexible systems, getting buy-in from researchers makes a substantial difference to the operation of the facility.

However, people management skills are also important in rules-based systems. Even if the rules are crystal clear and researchers have been fully trained, errors are likely to occur, particularly when researchers are working to deadlines.

As a result, checking researcher outputs requires statistical skills, personal management abilities, and a good awareness of institutional context. This is a different skillset from the checkers of official statistics, which focuses on consistent presentation of high-dimensional tables.

### **3 Course aims, structure and ethos**

#### **3.1 Aims**

The aim of the course is to ensure that attendees have the statistical skills to make sensible and efficient clearance decisions, but that they also have the personal and professional skills to ensure that they carry out their work effectively in collaboration with colleagues and researchers.

Specific learning objectives were therefore identified as, for each attendee,

1. Ensure confidence in dealing with known and unknown types of output
2. Understand the subjectivity of the process:
  - a. how to use judgement
  - b. respecting the views of others and coming to agreement
  - c. learn how to develop internal consistency
3. Develop a strategy for dealing with difficult outputs or unknown types
4. Develop skills in dealing with researchers
5. Integrate into and help develop output checking community

Course attendees are all required to have already attended and passed the Safe Researcher Training (SRT) accredited under the Digital Economy Act 2017, which includes an introduction to SDC. As course attendees are expected to be more motivated than SRT attendees, they can be reasonably expected to carry out pre-reading.

#### **3.2 Ethos**

The underlying ethos of the course is that of the ‘EDRU’ (evidence-based, default open, risk managed, user-centred) model (Ritchie and Green, 2016; Hafner et al, 2017). This approach emphasises the importance of psychology in process design; for a training course, the relevant elements are:

- the motivation of attendees
- how information is absorbed and maintained
- using group interaction to build self-learning

This leads to active learning through small-group exercises being the primary pedagogic method. Knowledge acquisition is strictly limited to information which can be reliably retained. This approach has already been successfully trialled and adopted in the SRT (Green et al, 2017).

### **3.3 Structure**

The course time was split 50-50 between statistical skills and people management. Attendees were split into teams. As the first exercise, attendees were asked to demonstrate that they had done the pre-reading via preparing summaries of the concepts in the pre-reading.

The statistical element involved giving attendees paper handouts representing dummy 'output requests' from a researcher. The exercises were of increasing complexity, starting with simple tabulations and ending with outputs which the checkers were unlikely to know in advance. Teams were given variable amounts of time to assess the output, as well as draft bullet-pointed feedback for the researcher. The teams, as expected, tended to identify many potential problems with the outputs. With assessments completed, a group discussion was used to focus minds on what could be reasonably expected in a practical environment. This is the evidence-based, default-open, risk-managed part of the EDRU approach. The user-centred element was addressed by considering how the researcher feedback could be used to encourage positive behaviours.

The people management exercise involved facilitated discussions on the problems caused by researchers. In various exercises, attendees were asked to identify major problems: 'major' could be interpreted as 'most frequent', 'most important', or 'most time-consuming'. Attendees were then asked to identify the source of problems and possible solutions. Again, the emphasis was on directing researchers towards positive behaviours.

### **3.4 Testing**

The course was followed up by a short test consisting of a number of outputs to be assessed, plus an essay-style question.

## **4 Lessons learned from the pilots**

Two pilot sessions were held (five and fifteen attendees, respectively). One consisted of a group of output checkers from ONS. The larger group included output checkers from other RDCs and SDC experts.

Most attendees had done the pre-reading; this simplified presentation. Those who hadn't were suitably embarrassed.

The statistical section worked relatively well from the start and stayed largely unchanged. The technique that proved most useful for assessing ‘likely’ disclosure in an output was the chain-of-events argument: “What set of circumstances would have to occur for something to be discovered from this output? And how likely is that chain of events?”

The people management section was unsatisfactory in the first pilot. The tutor reflected that the exercises were repetitive and didn’t lead to a clear outcome; feedback from attendees was also mixed and suggested that the tutor was pushing a ‘right’ answer. The exercises were revised and re-ordered for the second pilot and worked better. After the second pilot, reflecting on the ‘summaries’ presented after each people-management exercise suggested they were unnecessary and perhaps patronising, and still repetitive. These were replaced by a single summary.

Anonymous feedback was collected by the course organiser and via an online feedback form. These highlighted that

- Most users enjoyed the course and found it helpful (although one attendee felt that the material was too simplistic to be useful)
- Respondents generally reported feeling challenged and made to reflect on their implicit assumptions
- Some felt uncomfortable about being challenged on their organisation’s rules (note: the tutor focused on why the organisation had such a rule, not the rule itself)
- The group-discussion format helped understanding and forming opinions

On the key learning objectives (1) and (2):

- most attendees felt that they now had a better understating of the subjectivity of decision-making
- attendees generally were more confident in dealing with researchers, but some felt that the questions raised had made them less confident

One limitation on the course is the quality of the statistical examples. Creating meaningful examples that can still be analysed in a short space of time takes a considerable amount of time. While most attendees thought the exercises worked, a few felt that they were too simple and there was not enough progression.

## **5 Conclusion**

Overall, the course has met its broad aim: to give output-checkers (a) a more nuanced view of the process, and (b) the knowledge and confidence to operate in a more evidence-based, user-centred risk-managed way.

There remain areas for improvement. The people management will be tightened further; it was described by one attendee as a 'sharing session'. Future courses will have slightly more time devoted to the statistical section.

An Expert Group was identified to provide comments and suggestions for the course, including new examples. The Group did not meet in advance of the pilots, so that they could attend the sessions with no preconceptions. The Group will hold its first meeting in Autumn 2019, with most of the members having attended one of the pilot sessions. A particular role envisaged for the Group is to develop a wider range of statistical examples for both the course and the test.

NSIs interested in joining the Expert Group or finding more information about the course are welcome to contact the author.

## Acknowledgements

I am grateful to Adil Deedat, Lizzie Green and Bill South for reviewing this draft. All errors and omissions are those of the author.

## References

- Alves K. and Ritchie F. (2019) *Runners, repeaters, strangers and aliens: operationalising efficient output disclosure control*. Working papers in Economics no 1904, University of the West of England.  
<https://www2.uwe.ac.uk/faculties/BBS/Documents/UWE-working-paper-Runners-repeaters-strangers-and-aliens-Operationalising-efficient-output-disclosure-control.pdf>
- Brandt M., Franconi L., Guerke C., Hundepool A., Lucarelli M., Mol J., Ritchie F., Seri G. and Welpton R. (2010), *Guidelines for the checking of output based on microdata research*, Final report of ESSnet sub-group on output SDC.  
<http://eprints.uwe.ac.uk/22487/>
- Bond S., Brandt M. and de Wolf P-P (2015) *Guidelines for Output Checking*. Data Without Boundaries Project, Workpackage 11 deliverable.  
[https://ec.europa.eu/eurostat/cros/system/files/dwb\\_standalone-document\\_output-checking-guidelines.pdf](https://ec.europa.eu/eurostat/cros/system/files/dwb_standalone-document_output-checking-guidelines.pdf).
- Greci C., Griffiths E., Scott. J., Welpton R., Wolters A. and Woods C. (2019) *Handbook on Statistical Disclosure Control for Outputs*.  
<https://securedatagroup.org/sdc-handbook/>
- Green E., Ritchie F., Newman J. and Parker T. (2017) "Lessons learned in training 'safe users' of confidential data". *UNECE worksession on Statistical Data*

*Confidentiality 2017*. Eurostat.

[https://www.unece.org/fileadmin/DAM/stats/documents/ece/ces/ge.46/2017/7\\_lessons\\_learned\\_training.pdf](https://www.unece.org/fileadmin/DAM/stats/documents/ece/ces/ge.46/2017/7_lessons_learned_training.pdf)

- Hafner H-P., Lenz R., Ritchie F., and Welpton R. (2015) "Evidence-based, context-sensitive, user-centred, risk-managed SDC planning: designing data access solutions for scientific use", in UNECE/Eurostat Worksession on Statistical Data Confidentiality 2015, Helsinki. [https://www.unece.org/fileadmin/DAM/stats/documents/ece/ces/ge.46/20150/Paper\\_9\\_Session\\_4\\_-\\_Various\\_Hafner\\_et\\_al..pdf](https://www.unece.org/fileadmin/DAM/stats/documents/ece/ces/ge.46/20150/Paper_9_Session_4_-_Various_Hafner_et_al..pdf)
- Lowthian P. and Ritchie F. (2017) *Ensuring the confidentiality of statistical outputs from the ADRN*. Administrative Data Research Network Technical Report. <https://uwe-repository.worktribe.com/output/888435/ensuring-the-confidentiality-of-statistical-outputsfrom-the-adrn>
- ONS (2019) *Safe Researcher Training*. <http://www.saferesearchertraining.org/>.
- Ritchie, F., & Green, E. (2016). *Data Access Project Final Report*. Australian Department of Social Services. <https://uwe-repository.worktribe.com/output/908255/departement-of-social-services-data-access-project-final-report>