

UNITED NATIONS ECONOMIC COMMISSION FOR EUROPE
CONFERENCE OF EUROPEAN STATISTICIANS

Expert meeting on Statistical Data Confidentiality

26–28 September 2023, Wiesbaden

SACRO: Semi-Automated Checking Of Research Outputs

Jim Smith¹, Richard Preen¹, Maha Albashir¹, Felix Ritchie¹, Elizabeth Green¹, Simon Davy², Pete Stokes², Sebastian Bacon² 1: University of the West of England, UK, 2: Bennett Institute, University of Oxford

james.smith@uwe.ac.uk

Abstract

Output checking can require significant resources, acting as a barrier to scaling up the research use of confidential data. We report on a project, SACRO, that is developing a general-purpose, semi-automatic output checking systems that works across the range of restricted research environments. SACRO is designed to

- Automate checking of most common statistics, using best-practice principles-based modelling.
- Support researchers using the major analytical languages (R, Python and Stata), with minimal changes, by exploiting the ‘wrapper’ approach successfully trialled already.
- Support secure environments with different operating models and output checking workflows, through a process of co-design to maximise useability.

SACRO builds on previous work: (ACRO, funded by Eurostat and reported in in the 2021 Workshop) to establish the proof-of-concept; and Py-ACRO which showed how a software-independent tool might be developed. It differs from those earlier projects in terms of a wider range of statistics covered, and a requirement to achieve general applicability. To do this, the project draws on our extensive networks of practitioners. A series of workshops and ‘hands-on’ evaluations ensure the design frameworks support buy-in from a wide range of prospective users across health and social sciences, and from the public and private sectors.

1 Introduction

Statistical agencies and other custodians of secure facilities such as Trusted Research Environments (TREs) [Hubbard et al. \(2020\)](#) provide researchers with access to confidential data under the ‘Five-Safes’ framework [Ritchie \(2017\)](#). This enforces five orthogonal layers of safety procedures, and the last requires explicit checking of research outputs for disclosure risk. This can be a time-consuming and costly task, requiring skilled staff. This paper discusses the development of an open source tool for automating the statistical disclosure control (SDC) of routine research outputs. The goal is to make the clearance process more efficient and timely, and to allow the skilled checkers to focus their attention on the less straightforward cases.

The purpose of the tool, (SACRO, for Semi-Automated Checking of Research Outputs) is to assist researchers and output checkers by distinguishing between research output that is safe to publish, output that requires further analysis, and output that cannot be published because of substantial disclosure risk.

This work builds upon a previous Eurostat-funded project [Green et al. \(2020, 2021\)](#) in which Green, Ritchie and Smith developed a proof-of-concept prototype for the proprietary Stata software. The primary new contributions reported in this paper are:

- The implementation of a Python toolkit.
- An extensible multi-language platform with interfaces familiar to users of popular statistical tools.
- ‘Skins’ in Stata and the language R, demonstrating cross-language support.
- An open source repository with examples, help, documentation, etc.

2 Background

The Five Safes framework [Ritchie \(2017\)](#) is a set of principles that enable services to provide safe research access to their data and has been adopted by a range of TREs, including the Office for National Statistics (ONS), Health Data Research-UK (HDR-UK), and the National Institute for Health Research Design Service (NIHR), as well as many others worldwide.

Ensuring the last of these, ‘safe outputs’ is a complex and often costly human labour-intensive process. Automated output checking aims to improve the rigour and consistency of the output disclosure control process and reduce human workload by automatically identifying, reporting, and (optionally) suppressing disclosive outputs where possible and categorising outputs as ‘safe’ or ‘unsafe’. ‘Safe’ outputs requiring no or minimal further changes can be expedited through the clearing process whereas ‘unsafe’ outputs can be prioritised for human review [Ritchie \(2008\)](#).

A small number of SDC tools have been produced to assist in the process of achieving ‘safe outputs’, such as tauArgus and sdcTable¹, however these are primarily designed for users such as National Statistic Institutes as they require expert knowledge of SDC to use effectively. Moreover, they are designed for tabular outputs, and do not cover the range of statistics produced by researchers

With the aim of improving the efficiency of the process, and (where applicable) reducing the amount of user training required, a recent Eurostat project [Green et al. \(2021\)](#) developed a proof-of-concept prototype in Stata where primary disclosure is regulated by a set of simple rules. For example, a minimum threshold rule applied to the number of observations used by a statistic ensures that there is sufficient uncertainty with respect to any individual respondent. Dominance rules protect large respondent values from being approximated where the contribution to a statistic is dominated by only a few individuals. For example, the $p\%$ -rule sorts the N observations by magnitude and checks whether the sum of the smallest $N - 3$ observations is at least $p\%$ of the largest observation. The NK rule checks that the largest N observations contribute less than $K\%$ of the total. Also, not all aggregation statistics are permitted: reporting minima or maxima values of a subgroup are prohibited, and regressions are protected by checking that the Residual degrees-of-freedom exceeds a minimum threshold.

¹Respectively, <https://github.com/sdcTools/tauargus> and <https://github.com/sdcTools/sdcTable>

Building on the experience of the initial proof-of concept, funding was secured from the UK Research Council’s DARE initiative² for the project: Semi Automated Checking of Researcher Outputs (SACRO) which involves:

- Computer scientists with backgrounds ranging from AI research to commercial software development.
- A range of TREs as co-designers of a toolset.
- SDC theorists and statisticians to provide a conceptual framework for handling different types of output and providing guidance to researchers and output checkers.
- Public Involvement and Engagement specialists and groups to develop a consensus statement around the use of (semi)-automation in disclosure control
- Researchers from a previous DARE project examining the output checking of machine learning models trained on sensitive data within a TRE [Jefferson et al. \(2022\)](#).

In this paper we report on the principal tools developed within the SACRO project, specifically:

1. A toolkit for researchers to use within TREs that produces automated reports on disclosure risk with minimal changes to their practice - simply prefixing common commands with the word ‘acro’.
2. Explicit support for researchers to reduce the number of disclosive outputs they request.
3. Cross-language support: with exemplar interfaces provided for Stata and R.
4. Support for the output types that our TRE partners tell us form the majority of requested releases.
5. A stand-alone viewer for TRE output staff to facilitate rapid, informed, and audited, decision making.
6. A revised guide incorporating theoretical developments, directly linked to its implementation in SACRO.

3 The SACRO toolkit

SACRO is composed of three parts which may be deployed independently: the main ‘ACRO-engine’, a stand-alone viewer, and ‘AI-SDC’ - support for disclosure control of machine learning models (described elsewhere).

3.1 Design Philosophy

The operational design philosophy is extensively documented in [Green et al. \(2020\)](#), who studied the characteristics that an automated solution needs to have to be feasible, effective, and a positive choice for users. Essential criteria are that it should be:

- Acceptable to users, output checkers and TRE managers;
- Able to implement an organisation’s business rules for primary and secondary disclosure, which may vary across datasets or users;
- Comprehensive, even if the automated tool’s response is “I don’t know so this needs manual checking”;
- Consistent, providing the same results across different studies within a TRE, and across TREs;
- Able to support exceptions under principles-based regimes;
- Scalable over users and outputs.

Key operational requirements were for the tool to work in different technical environments, and to be easily updated through well understood mechanisms. This meant separating the software itself (distributed through a recognised channel³, from the specification of a given TRE’s risk appetite (held in a human and machine readable and editable file).

Acceptability to users was identified as the most crucial element. If researchers and output checkers see the tool as something that makes their life better and easier, then they are more likely to use it effectively. Hence, designing the user interface was identified as a separate workstream in SACRO, and given the same resources as the design and implementation of the output-checking component. This is also one reason why SACRO set up a large network of potential users and tests (see Sec. 6 below).

²<https://dareuk.org.uk/>

³for example, PyPi (<https://pypi.org>) or CRAN (<https://cran.r-project.org>)

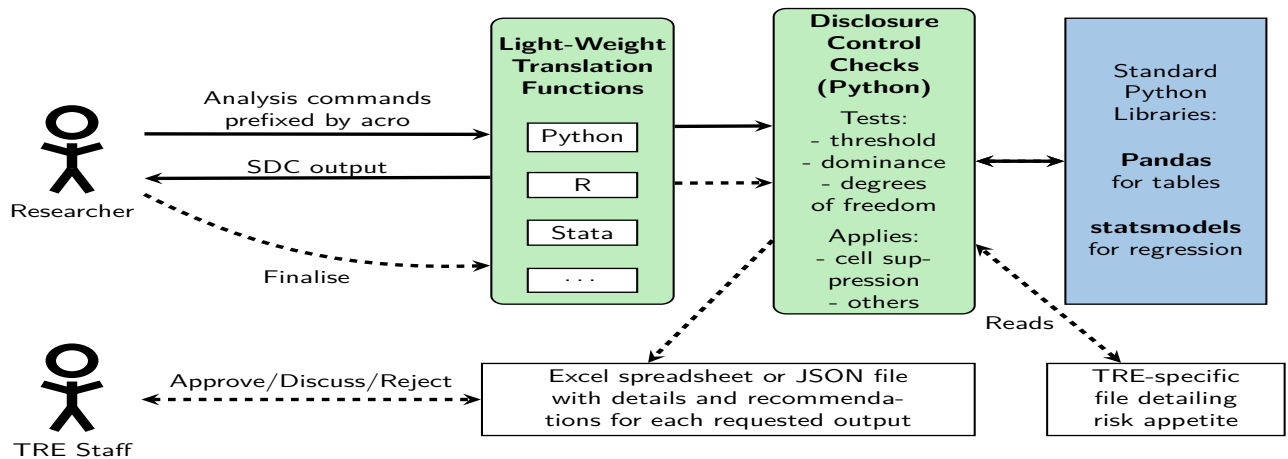


FIGURE 1. Schematic illustration of ACRO.

The ‘proof-of-concept’ version of ACRO did not address secondary disclosure (such as checking for differencing across tables), for two reasons. First, business rules for secondary checking are often not clear or comprehensive. Second, ACRO/SACRO works by intercepting commands and assessing disclosure risk at the time the output is being produced. Analysing results post-hoc is a considerably harder problem, requiring the researcher to produce a lot more information and also locate the other outputs to be compared. Although SACRO does not currently (as of July 2023) carry out secondary disclosure review, we are investigating how to at least flag potential differencing risks across the set of outputs from a research ‘session’, and in future, create a library of outputs which might allow secondary disclosure to be assessed, even if only partially.

3.2 Workflow

ACRO Preen et al. (2023) is an open source toolkit (MIT License) that provides a light-weight ‘skin’ that sits over well-known analysis tools, in a variety of languages researchers might use. The process is illustrated in Fig. 1. This adds functionality to identify potentially disclosive outputs against a range of commonly used disclosure tests and report to researchers and TREs reasons why outputs should not be released ‘as-is’. It creates simple summary documents TRE staff can use to streamline their workflow.

ACRO has been designed with the following aims:

- Reducing barriers to adoption via a front-end application programming interface (API) that is similar to those already commonly used by researchers in their favoured language.
- Providing researchers with: immediate feedback on the results of disclosure checks (on-screen alongside their query results); facilities to add comments or exception requests, and control over what is submitted for review, e.g., removing disclosive outputs if they use feedback to design non-disclosive ones.
- Having a single back-end code base constituting a single source of truth for performing checks, with extensibility for different languages and ongoing support and consistency.
- Providing easy to understand help and documentation.

In practice, researchers prepare their data and statistical queries in the usual way, in their preferred language, using common commands prefixed by ‘acro’. The lightweight ACRO translation functions then call the Python back-end, which executes the queries and performs the requisite output checks. The results of the checks, and the queries are immediately displayed to the researcher, and full details are stored in a list. When the user calls `acro.finalise()` to end their session, outputs and all SDC details are saved to file for review by a TRE output checker. A schematic illustration of the ACRO workflow is shown in Figure 1 and some notebooks demonstrating example code usage and output are available via the ACRO project wiki⁴.

⁴<https://github.com/AI-SDC/ACRO/wiki>

3.3 Checks Implemented

For tabular data (e.g., cross tabulation and pivot tables), we prohibit the reporting of the maximum or minimum value in any cell that represents a sub-group of one or more contributors. Moreover, we suppress, and report the reason, the value of the aggregation statistic (mean, median, variance, etc.) for any cell deemed to be *sensitive*. ACRO currently supports the three most common tests for sensitivity: ensuring the number of contributors is above a frequency threshold, and testing for dominance via $p\%$ and NK rules. ACRO builds a series of suppression masks, which indicate which cells are to be suppressed for each check. A summary outcome table indicating which suppression rule was applied to each cell is presented to the researcher (the grey box in Fig. 2, alongside the query results. For regressions, e.g., linear, probit and logit regression, the tests verify the number of degrees of freedom exceeds a threshold. Immediate feedback on all these checks is designed to support researchers to improve their practice and so reduce the SDC bottleneck by making fewer disclosive requests. The checking of graphical plots is not currently implemented, as this is a complex problem with many different methods for producing visualisations. However, we expect to have some support by Autumn 2023. As noted above, all of these tests and checks are configurable according to the TRE’s risk appetite. The data custodian, e.g., TRE staff member, specifies the parameter values used for the output checks in a YAML⁵ configuration file, which is loaded upon ACRO initialisation. The default ACRO parameters are shown in Table 1. Future releases will offer the option to over-ride these on a dataset, or even attribute level.

TABLE 1. ACRO Default Parameters for sensitivity tests

Description	Parameter	Value
Min frequency threshold for tabular data	<code>safe_threshold</code>	10.0
Min degrees-of-freedom for analytical stats	<code>safe_dof_threshold</code>	10.0
N parameter in NK test	<code>safe_nk_n</code>	2.0
K parameter in NK test	<code>safe_nk_k</code>	0.9
Min ratio for $p\%$ test	<code>safe_ratio_p</code>	0.1

3.4 The SACRO Python ‘Engine’

Python is a popular multi-platform language widely used for data analysis and machine learning. PyPI provides a simple package management system for distributing open source Python libraries. Pandas and Statsmodels⁶ are industry-standard, mature, popular, and well-supported python packages for data analysis, statistical testing, and statistical data exploration. Pandas is currently used by more than 55% of all Python users [Python Software Foundation \(2021\)](#) and there are many web-sites and user groups providing help with formulating queries.

The use of Python as the primary implementation therefore enables the leveraging of existing expertise and community support with these packages so that the ACRO front-end can be as similar to the API researchers already know and trust, and further facilitates the rapid development of disclosure checking functionality on the back-end. As the PyPI distribution system is simple and allows the use of semantic versioning, it supports a rapid and iterative develop-and-deploy strategy to provide continuing functionality and improvements.

For example, the current version of ACRO may be installed [or updated] as simply as:

```
pip install [--upgrade] acro
```

The currently implemented methods are listed below, split into analysis commands, and sessions management commands. For more details see the ACRO project documentation⁷.

⁵<https://yaml.org>

⁶<https://github.com/pandas-dev/pandas> and <https://www.statsmodels.org/stable/index.html> respectively

⁷<https://ai-sdc.github.io/ACRO/>

3.4.1 Analysis commands for Researchers. These are implemented via the use of multiple inheritance from Pandas and Statsmodels. For making tables, the relevant methods are:

- : **crosstab**(index, columns[, values, rownames, ...])
Compute a simple cross tabulation of two (or more) factors,
with options for hierarchies in rows/columns and multiple aggregation functions.
Same API as pandas.crosstab.
- : **pivot_table**(data[, values, index, columns, ...])
Create a spreadsheet-style pivot table as a DataFrame.
Same API as pandas.pivot_table.

and for regression analysis:

- : **logit**(endog, exog[, missing, check_rank])
Fits Logit model.
Same API as statsmodels.discrete.discrete_model.Logit.
- : **logitr**(formula, data[, subset, drop_cols])
Fits Logit model from an R-style formula and DataFrame.
Same API as statsmodels.formula.api.logit.
- : **ols**(endog[, exog, missing, hasconst])
Fits Ordinary Least Squares Regression.
Same API as statsmodels.regression.linear_model.OLS.
- : **olsr**(formula, data[, subset, drop_cols])
Fits Ordinary Least Squares Regression from an R-style formula and DataFrame.
Same API as statsmodels.formula.api.ols.
- : **probit**(endog, exog[, missing, check_rank])
Fits Probit model.
Same API as statsmodels.discrete.discrete_model.Probit.
- : **probitr**(formula, data[, subset, drop_cols])
Fits Probit model from an R-style formula and DataFrame.
Same API as statsmodels.formula.api.probit.

3.4.2 Session Management Commands.

- : **ACRO** () (config,suppress)
Creates an ACRO session object with optional parameters for a config (risk appetite) filename
and whether disclosive tables should have suppression applied (default False).
- : **print_outputs**()
Prints the current results dictionary - i.e., the outputs that would be sent for checking.
- : **remove_output**(key)
Removes an output from the results dictionary.
- : **rename_output**(key, newname=)
Assigns a new (ideally more self-explanatory) name to an output from the results dictionary.
- : **add_comments**(key,text)
Allows researcher to add a description for an output
- : **add_exception**(key,text)
Allows a user to request and justify an exception to strict rules-based checking.
- : **custom_output**(filename,description)
Adds a file containing output from unsupported analysis to an ACRO session for inclusion in outputs
shown in viewer.
- : **finalise**(directory_name, format)
Creates a results file for checking in the desired format(*json* or *xlsx*).

```

» safe_table = acro.crosstab(
  df.recommend, df.parents, values=df.children, aggfunc="mean")
» print(safe_table)

```

```

INFO:get_summary:fail;
threshold: 4 cells may need suppressing

INFO:outcome_df:
parents      great_pret  pretentious  usual
recommend
not_recom    ok           ok           ok
priority     ok           ok           ok
recommend    threshold  threshold  threshold
spec_prior   ok           ok           ok
very_recom   threshold  ok           ok

INFO:acro:add():  output_1

```

grant_type	great_pret	pretentious	usual
recommend			
not_recom	1440	1440	1440
priority	858	1484	1924
recommend	0	0	0
spec_prior	2022	1264	758
very_recom	0	132	196

FIGURE 2. Example ACRO query for the ‘nursery’ data(top), with immediate disclosure control reporting (middle, grey background - pink onscreen) followed output (bottom). This ‘researcher-view’ corresponds to the top image in the viewer screenshots

An example ACRO query run on the nursery admission dataset⁸ and its output is shown in Fig. 2. This is the ‘researchers-view’ of the output at run-time. The corresponding ‘TRE-view’ is shown in the top screenshot in Fig. 3. This example does not have an aggregation function so dominance rules are not applied, otherwise they would also show in the ‘INFO’ section of the report in any relevant cells. Note that if the user starts their session with `acro= ACRO(suppress=True)` then any disclosive cells would have their values set to NaN

3.5 The R interface to ACRO

The R front-end is an example of cross-language support. It provides a set of wrapper functions that execute Python back-end checking via the reticulate⁹ package, which provides automatic conversions for many types, e.g., R data frame to Pandas DataFrame.

A session is created when the acro package is called `source("../acro.R")` and thereafter the acro methods work as callable functions with the prefix `acro_` e.g., `acro_rename_output(output5, "xy-plot")` etc., and to end a session the user calls `acro_finalise(results_dir, "json")`

For regressions, the common R `lm()` and `glm()` functions were shadowed with equivalent versions implemented as `acro_lm()` and `acro_glm()`, respectively. For tabular data, the dplyr¹⁰ package is commonly used within R, however no simple cross tabulation or pivot table functions are provided; instead various combinations of `groupby()` and `summarize()` etc. are used. Therefore, at this stage of development, the Python cross tabulation and pivot table functions were directly interfaced with `acro_crosstab()` and `acro_pivot_table()`.

⁸<https://www.openml.org/search?type=data&sort=runs&id=1568&status=active>

⁹<https://github.com/rstudio/reticulate>

¹⁰<https://github.com/tidyverse/dplyr>

3.6 Stata Interface

This makes extensive use of Stata's `SFIToolkit` library to manage a python session, transfer data in memory from stata to a Pandas dataframe in the python session, and results back to the Stata window. A simple `acro.ado` file defines a new function `acro` which takes as parameters either one of the ACRO session management methods (adding `init()` to start a session) or the name of a standard Stata function such as `table`, `regress` etc. Stata's inbuilt parsing functions are used to separate out the parts of command and pass them as lists to a python function `parse_and_run()` which handles the rest of the translation between the two languages.

4 SACRO Viewer for Output Checking

We have also created an open-source platform-independent stand-alone viewer for output checkers to use to: view outputs and their risks; make decisions with reasons (all recorded for auditing purposes); and produce zipped packages of files for release [Open-Safely \(2023\)](#). Figure 3 illustrates two screenshots from the version currently (July 2023) being evaluated by TREs. The viewer supports and renders a range of different file types for results from unsupported queries. A separate script lets TRE staff create an ACRO session from a set of output files in a directory, and hence use the viewer for making and recording decisions, even if the researcher has not used ACRO during their analysis. Automated disclosure risk analysis is not provided in those cases.

5 Linking theory and implementation

As part of the project, the SACRO team committed to review and re-develop the theory and operational guidelines for output SDC. The aim was threefold; first, to bring together key points from the OSDC literature (and fill in some of the theoretical gaps) to provide an integrated guide to both theory and practice of output checking; second, to develop a new approach to OSDC based on classifications into groups (see [Derrick et al. \(2023\)](#), for details); third, to explicitly link theory to operational rules and their implementation in manual and automatic checking regimes. The third aim is essential to demonstrating that SACRO is not seen as a 'black box' implementing its own rules, but is fully integrated into core theory. It is also important for showing how manual and automatic output checking necessarily differs. For example, dominance checks are almost impossible for a human, but straightforward for computers; on the other hand, computers cannot easily identify whether zero cells in tables are structural or disclosive, but humans can. The purpose of the guide is to show precisely what checks have been made, where differences occur between humans and computers, and why they are necessary.

6 Engagement with TREs

One of the lessons learned from the original Stata version of ACRO [Green et al. \(2021\)](#) was the importance of user buy-in. Although that version met its design goals (and has subsequently been adopted by Eurostat in its TRE), reaction to it was a mixture of "*this looks useful, I'll give a go*", "*this looks useful, I'll wait to see it installed before I commit myself*", and "*I've read the installation manual and have no idea what's going on, so it's a no*". As a result, that version of ACRO has remained largely within the project remit: a demonstration of possibilities. The SACRO project was intended to involve co-design from the outset to take ACRO to the next stage, of general utility and application. This involved three tests:

1. Would a new tool be acceptable to users?
2. Would a new tool be acceptable to output checkers?
3. Could a new tool be installed in secure research environments?

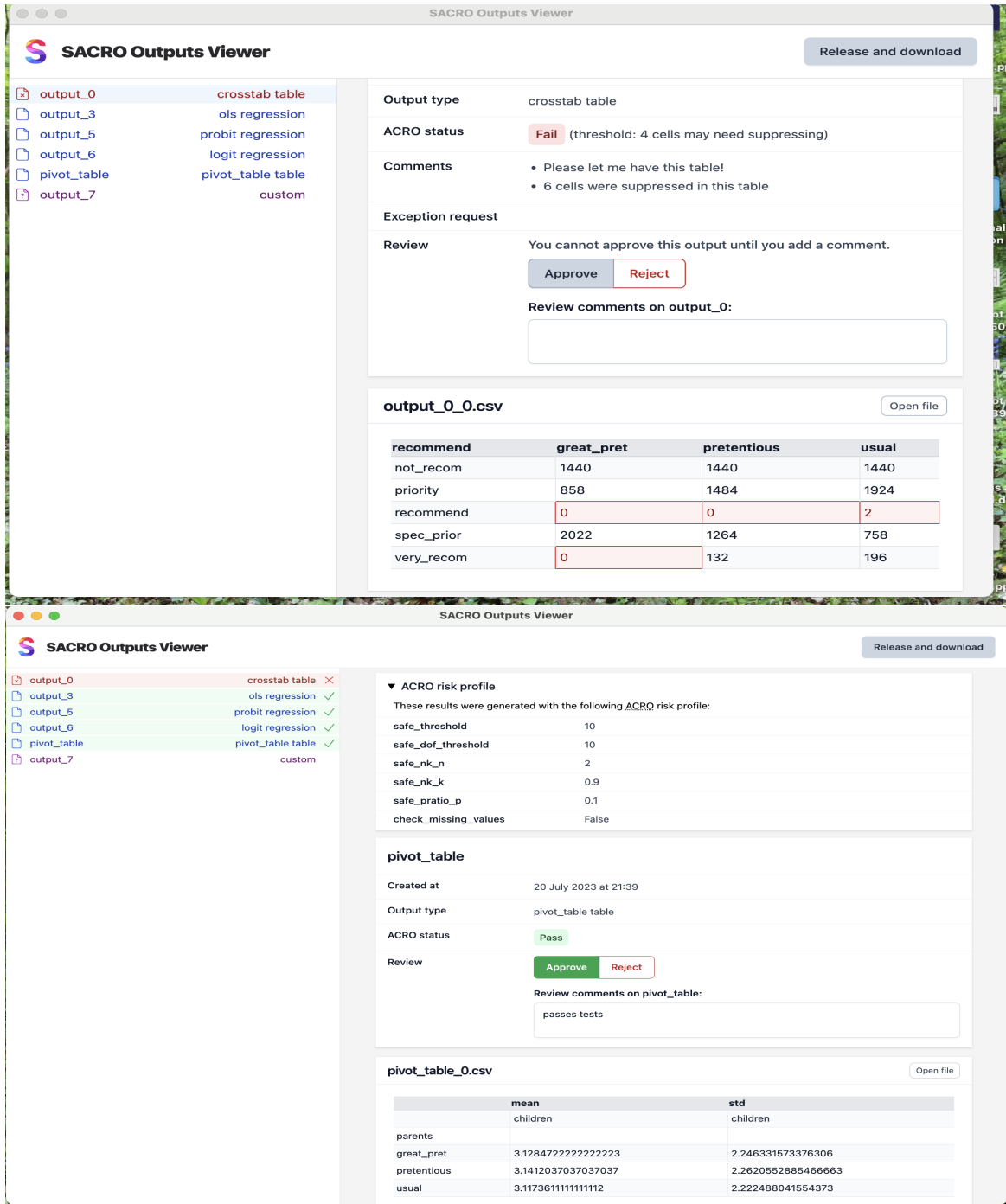


FIGURE 3. Two screenshots of viewer. The left hand column shows list of files requested. In top image, colouring of file names suggests which files require special attention. In lower image background colour-coding and tick/cross symbols show decisions made by output checker. Top image shows checker viewing table that fails disclosure tests, with problematic cells highlighted in red. Bottom shows acceptable table. Also in this image the top right hand panel shows option to view TRE 'risk appetite' expanded.

The SACRO project took two approaches. First, six TREs (OpenSafely at the University of Oxford, and the five Scottish Safe Havens) were funded as co-investigators on the project to provide detailed feedback on user and output checker perspectives (OpenSafely also took the lead in the design of the user interface). This group also directly tested the feasibility of installing and allowing the Python code to run on their systems as TREs differ in their perceptions of python’s ‘riskiness’. Second, the SACRO team contacted a large number of TREs in the UK and abroad, and set up a network of interested parties potentially willing to be testers. Several engagement events with this group identified how they worked and what they would expect from an automatic solution. At the time of writing (July 2023), the first ‘external’ TRE’s are starting to install and run the tool with genuine users. SACRO has a workpackage dedicated to helping TREs set up their systems, and then collecting evaluation feedback. This aims to make sure that the tool is tested in as wide a variety of environments as possible, given the time constrain. A secondary aim is to involve TREs in the development, to build a sense of ownership and lay the foundations for widespread adoption. This helps to address the concerns of ‘wait-and-see’ TREs.

7 Future Plans

By the current project end in October 2023 we aim to have added support for: more common types of analyses (including simple plots); different versions of Stata; and more ways of creating tables within R. Additional features and improved user experience will be facilitated by the involvement of end-users and output checkers. Beyond then, UWE has committed to web hosting various resources for the indefinite future, and partners have agreed to continue support and development of the toolkits. We are keen to engage with any interested parties to enrich and build an on-going community of support for SACRO.

References

- Derrick, B., E. Green, F. Ritchie, J. Smith, and P. White (2023). Towards a comprehensive theory and practice of output SDC. In *UNECE/Eurostat Workshop on Statistical Data Confidentiality*.
- Green, E., F. Ritchie, and J. Smith (2020). Understanding output checking. Technical report, European Commission (Eurostat - Methodology Directorate).
- Green, E., F. Ritchie, and J. Smith (2021, October). Automatic checking of research outputs (ACRO): A tool for dynamic disclosure checks. *ESS Statistical Working Papers 2021*, 1–27. doi: 10.2785/75954.
- Hubbard, T., G. Reilly, S. Varma, and D. Seymour (2020, July). Trusted research environments (TRE) green paper. *ZENODO 2020*, 1–31. doi: 10.5281/zenodo.4594704.
- Jefferson, E., J. Liley, M. Malone, S. Reel, A. Crespi-Boixader, X. Kerasidou, F. Tava, A. McCarthy, R. Preen, A. Blanco-Justicia, E. Mansouri-Benssassi, J. Domingo-Ferrer, J. Beggs, A. Chuter, C. Cole, F. Ritchie, A. Daly, S. Rogers, and J. Smith (2022, September). GRAIMATTER Green Paper: Recommendations for disclosure control of trained Machine Learning (ML) models from Trusted Research Environments (TREs).
- Open-Safely (2023). Sacro:a tool for fast, secure and effective output checking, which can work in any TRE. <https://github.com/opensafely-core/sacro>.
- Preen, R. J., J. Smith, M. Albashir, and S. Davy (2023). ACRO. <https://github.com/AI-SDC/ACRO>.
- Python Software Foundation (2021). Python developers survey 2021 results. <https://lp.jetbrains.com/python-developers-survey-2021/>. Accessed: 24/07/2023.
- Ritchie, F. (2008). Disclosure detection in research environments in practice. In *Joint UNECE/Eurostat work session on statistical data confidentiality*, Volume WP. 73. United Nations Statistical Commission and Economic Commission for Europe Conference of Europe Statisticians, European Commission Statistical Office of the European Communities (Eurostat).
- Ritchie, F. (2017, September). The ‘five safes’: A framework for planning, designing and evaluating data access solutions. *Zenodo 2017*, 1–5. doi: 10.5281/zenodo.897821.