

# **An informational infrastructure for the E-Science Age – On the way to remote data access for business data**

Maurice Brandt<sup>1</sup>, Markus Zwick<sup>2</sup>

<sup>1</sup>Federal Statistical Office Germany, e-mail: [maurice.brandt@destatis.de](mailto:maurice.brandt@destatis.de)

<sup>2</sup>Federal Statistical Office Germany, e-mail: [markus.zwick@destatis.de](mailto:markus.zwick@destatis.de)

## **Abstract**

For a couple of years, access to microdata in Germany has been possible through different ways of access. The researcher can use so called scientific use files, even for business microdata, in his or her own institution, he or she can visit the safe centre or can use remote data execution.

There are still reservations in the scientific community concerning the data perturbation methods for business microdata. The user needs are developing towards on-site data access, because most often researchers would like to work with original data.

On the one hand this leads to a higher burden for the employees of the research data centre (RDC) because they have to apply the analysis programs and have to deal with the manual output checking of the results. On the other hand the researchers themselves have to wait longer to get their results.

The project "An informational infrastructure for the E-Science Age" deals with the improvement of remote access in the National Statistical Institutes (NSIs). The project aims to find solutions for better remote access in Germany through so called data structure files and automatic output checking procedures.

**Keywords:** [Remote Access, Output Checking, Business Data]

## **1. Introduction**

Since the research data centres (RDCs) of the statistical offices of the German Federation and Länder were set up in 2001, they have become firmly established and today empirical science is unimaginable without them.

The demand for data of official statistics has reached a dimension which is very difficult to deal with by the research data centres within a reasonable period of time. What is particularly time-consuming is on-site use of the data, i.e. the data are processed at the safe centre of the research data centre or through remote data execution; at the end of either of those processes, the tables containing the results are checked for confidentiality.

Another option is to use the scientific use files (SUFs), which can be analysed on the researcher's own workstation in the relevant institution. There is much demand for

SUFs with regard to person-related surveys such as the microcensus. For business data, more intensive anonymisation measures are performed because data perturbation methods such as microaggregation or stochastic noise are inevitable to protect the enterprises. Due to reservations regarding data perturbation methods of anonymisation and due to longer waiting times for providing a SUF, demand for data from economic statistics is increasingly shifting towards original microdata through on-site use, and in particular towards remote data execution. That form of access is highly comfortable for researchers as it does not require travelling to, and staying at the RDCs. Upon request, the researcher will get a so-called data structure file consisting of a fully anonymised data set whose structure is identical to that of the original data and which can be used to write the program for data analysis. For the RDC staff, that form of access is time-consuming for two reasons. First, the programs must be adjusted several times because the current data structure files do not reflect very well the original data as they are strongly anonymised (sampling and exchanging). Second, manual checking of results and checking for confidentiality takes much of the overall time needed for a remote data execution order.

Both points are to be treated by the project "An informational infrastructure for the E-Science Age - On the way to remote data access for business data" described in this paper.

## **2. Current situation at the research data centres**

The research data centres (RDCs) of the statistical offices observe a fundamentally changing demand for their products. It turns out that remote data execution and safe centres have become the most frequently used forms of access to microdata of economic statistics in Germany. Demand thus focuses on on-site data use, which is highly time-consuming for both data producers and data users (Zwick 2006 and Zühlke et al. 2004).

With the new projects *Amtliche Firmendaten für Deutschland (AFiD)* (Official company data for Germany) and *Kombinierte Firmendaten für Deutschland (KombiFiD)* (Combined company data for Germany) the statistical offices, together with other partners, deal with highly complex data sets (cf. Bender et al. 2009). Ensuring de facto anonymity and, at the same time, maintaining a maximum of the analysis potential is a complex job even for cross-section and panel data of economic statistics. Considerable impairment of the analysis potential is inevitable. It is foreseeable that the new data sets, which are currently produced in AFiD and KombiFiD and will be linked in cross-section terms and over time, would – as a result of the anonymisation measures required to produce a SUF for off-site use – lose their information potential that will have been redesigned by such linking. As a result, demand for time-consuming on-site use will further increase. However, the RDCs of the various data producers have to cope with heavy workload even now as their capacities are widely used especially in remote data execution.

## **3. Project activities**

A real remote access application, which is fully automated and does not require any manual handling, is a vision for the future which has not been achieved in countries with a legal frame like Germany either. A reasonable interim goal for Germany in the

medium term is a remote access solution like those implemented in the Netherlands, Denmark or at NORC<sup>1</sup>. The project “An informational infrastructure for the E-Science Age - On the way to remote data access for business data” is necessary to take first steps towards that goal. The next (or parallel) step for Germany can be an “RDC in RDC” solution where data of an RDC can be processed in another RDC, using remote access. This can be used as a test implementation for real remote access to be established later and permits shifting activities towards more data exploration, data documentation, internationalisation and less visitor care.

Before a pure remote access application can completely be implemented, however, many technical, legal and – the focus of the project – methodical problems must be solved. Although first applications are available – Lissy in Luxembourg<sup>2</sup> and the methods of the Dutch<sup>3</sup> and Danish<sup>4</sup> national statistical institutes –, none of them provides fully automated access routines and part of them, e.g. Lissy, are limited to specific applications. In Sweden, MONA<sup>5</sup> is a feasible remote data access solution which, however, can be applied only because of the particular legal situation of data use in Sweden. In Germany, SAM<sup>6</sup> is a first technical solution. JoSua<sup>7</sup>, too, might be extended to become such an application.

The research project is to set the bases for the following three methodical steps:

- (1) Developing anonymous data structure files which can be used to specify analysis models and must therefore be suitable for semantic analysis and which allow developing analysis programs that are error-free in terms of syntax.
- (2) Developing and assessing methods of standardised and completely automated checking of results.
- (3) Simultaneous consideration of microdata anonymisation and checking of results.

The project’s purpose is to develop basic strategies for producing anonymised data structure files which allow checking a program run for syntactic and semantic errors. Current data structure files – which is what users of remote data execution get – allow only syntactic checking. Methods that might be applied to produce such data structure files are in particular the data perturbation methods of multiplicative stochastic noise, multidimensional microaggregation and multiple imputation.

Checking the results is always time-consuming and labour-intensive. Results of remote data execution and of activities performed at safe centres are checked for confidentiality before being released. Such checks are extremely difficult for complex tables and large estimation output. Automated procedures have been developed only

---

<sup>1</sup> NORC: National Opinion Research Center at the University of Chicago

<sup>2</sup> See e.g.: Coder, John and Cigrang, Marc (2003): LISSY Remote Access System

<sup>3</sup> Hundepool, Anco and de Wolf, Peter-Paul (2005): OnSite@Home: Remote Access at Statistics Netherlands

<sup>4</sup> Borchsenius, Lars (2005): New Developments in the Danish system for access to microdata.

<sup>5</sup> MONA: Microdata ON-line Access, Statistics Sweden

<sup>6</sup> Heitzig, Jobst (2006): *Wissenschaftsserver zur Auswertung von Mikrodaten* (Science servers for microdata analysis)

<sup>7</sup> The data centre of the Institute for the Study of Labor in Bonn (IIZA) has developed an application allowing external researchers to start microdata analyses via the internet. That application (JoSuA) is, first, user-friendly because researchers can monitor the status of their orders from their workstations and, second, it facilitates IIZA activities because it is no longer necessary to start the programs manually.

in some cases for standardised results becoming available regularly. For the flexible analyses performed in the research data centres, the methods developed so far – even at the international level – are far from sufficient. The project's function here is to extend the issue and to perform a systematic comparison between data-based and result-based safeguarding of the protection of the carriers of variables with regard to the analysis potential. It should be a goal to develop methods allowing the user to decide, before running the analysis, whether it should be performed with anonymised data and without restrictions for results or on the basis of the original data set and with restricted release of results. At this point, mixed forms could be envisaged.

### **3.1 Producing data structure files**

A first goal of the project is the standardisation of the data in the form of so-called data structure files. Such anonymised data sets, which have the same structure as the original data sets, are sent to researchers who made a request for use, so that they can develop their program codes for analysis and send them to the relevant research data centre. That program code is then applied by the RDC staff to the original data and the output is returned to the researchers after checking for data security and confidentiality.

So far, many of the data structure files consist of a sample of the original material, which has been subjected to additional anonymisation measures, or of values generated at random within the value range of the data set. Although the variables are maintained in both approaches, their attributes and the dependence structures (filter, variance-covariance matrix) with regard to other variables are completely destroyed. Although researchers can check whether their programs are executable, they do not get any information on whether the actual issue has adequately been implemented. Therefore, in many cases, the researchers' analysis programs cannot be taken in an identical form for the subsequent application to the original data. Often adjustments must be performed by the researchers and the RDC staff.

For more complex data such as the Linked Employer-Employee Data Sets (structure of earnings survey (1995, 2001, 2006), Linked Employer-Employee Data Set of the Institute for Employment Research – IAB), the data structure files in their current form are not very helpful because, in practical work, they cause enormous adjustment problems and require much co-ordination between external researchers and the research data centres. For example, such data structure files do not allow performing consistency checks of whether the analysis program developed by the researcher is correct and can be applied without errors to the original data. To allow suitable adjustment of the programs to the remote data execution procedure, empirically working scientists must be able to test not only univariate calculations but also, and to an increasing extent, programs referring to multivariate issues. Evaluating multivariate analyses with the existing data structure files is difficult because the covariances are not maintained here and because researchers, too, must repeatedly adjust the models with regard to the original data. This causes problems especially where the full number of observations is needed (e.g. in analyses with several waves). In addition, not all logical restrictions are always correctly represented. This means that the developed program codes must be sent several times between the researchers and the staff of the relevant institution before the desired result is available. But it also

means that all analyses performed must be checked by the staff to ensure that data protection is not violated when providing the results.

Another advantage of data structure files that are error-free in terms of semantics and syntax is that the researchers can determine more exactly the number and extent of their tables containing the results and that they can adjust the evaluations on their own workstation until the desired tables are produced. This reduces the efforts required to check tables for confidentiality which possibly are not included in the publication.

Data perturbation methods of anonymisation have already been developed or been adjusted to the requirements of economic statistics of the German statistical offices or of the Federal Employment Agency (cf. Ronning et al. 2005 and Bender et al. 2008). The question of whether the developed data sets are fully anonymised public use files or de facto anonymised must be answered in the course of the project. A public use file has the advantage that it can be sent to the researcher already before, or during the process of making the request. For de facto anonymised data structure files, the request for use must first of all be decided upon and a contract on data use must be signed before the data material can be sent to the researchers. Where a public use file sufficiently reflects the data structure, it should be preferred to a de facto anonymised data set because the analysis program will be applied to the original data anyway, irrespective of whether the data structure file has been fully or de facto anonymised. The question of whether the data structure can sufficiently be represented by a fully anonymous data set or only with a de facto anonymous data set can be answered only within the scope of the project activities.

Before a true remote access solution can be implemented in the form of a functioning technical infrastructure, it is indispensable to use data structure files because this provides the researchers with some flexibility when working on the analysis programs and because those programs can be developed independent of time and location.

For a functioning remote access, too, the data structure files – once developed – continue to be necessary. As viewing the original data on the user's screen would represent a transmission of original data (which is problematic in terms of data protection), it might be useful for remote access to present on the researcher's screen the view of the microdata from the data structure files. When calculating the analyses, however, the original data are used. The data themselves, i.e. both the data structure files and the original data, remain on the servers in the protected rooms of the statistical offices. In this respect, the data structure files are both an interim solution until real remote access has been developed and a final product which can be implemented as an important component of a hardware-based access solution.

### **3.2 Result-based confidentiality**

Remote access is an optimal solution to provide the scientific community with access to confidential data. It allows external researchers to perform analyses on their own computer via a remote server and the results are displayed in real time. A problem here is the question of how a researcher can be controlled in case of data protection being violated or how such violation of data protection can definitely be prevented.

The safest way is to check the results for potential data protection risks before they are transmitted in real time.

So far, the results of remote data execution and of activities performed in safe centres are manually checked for confidentiality before they are released. Such checking is highly difficult for complex tables and large estimation output and it is very time consuming and labour-intensive.

To ensure confidentiality of tables, part of the cells are generally suppressed in case of business data. As the tabular data disseminated by the statistical offices are linearly linked with each other through subtotals and marginals, i.e. they are additive, additional cells must be suppressed (“secondary suppression”), to protect the primarily secret cells against disclosure through subtraction. Establishing suitable secondary suppression – which would minimise the loss of information caused by the suppression – is a complex linear problem of optimisation. An overview of the common standard methods is given, for example, in Giessing, 1999.

Frequently, part of the cells of a specific table are identical to cells of another table. In such cases, the selection of secondary suppressions must be co-ordinated across tables to avoid a situation where users can disclose suppressed cells by performing comparisons between such overlapping tables.

Checking the results of tables produced in safe centres or through remote data execution is necessary, among other reasons, to ensure that the results finally published by the users cannot disclose the secondary suppressions from the statistical offices’ own publications. This requires co-ordinating secondary suppression between standard publications and users’ tabular data. There are no suitable automated methods available for that purpose. As handling of user tables is subordinate to handling standard publications here, considerable impairment of the quality of results due to massive secondary suppression is expected.

For the above, and other reasons, data perturbation methods have increasingly been proposed in the literature over the last few years as an alternative or in addition to the suppression methods. Most of the methods proposed perturbate the data at the aggregate level. However, there are also methods perturbing the microdata. As maintaining the quality of results in specific tables is the main focus also for those methods to be applied to microdata, such methods are considered as confidentiality methods for tables rather than anonymisation methods for microdata.

Data protection regarding table output has been a topical issue for a long time already at the statistical offices of the German Federation and Länder, whereas no systematic study has been performed yet on the problem of data protection for estimation output and non-linear analyses. The studies by Gomatam, Karr, Reiter and Sanil (2005) may be used as a basis here.

Heitzig (2005) provided an approach to result-based confidentiality. It transfers the idea of the jackknife approach to estimating standard errors to the confidentiality of microdata. The method’s basic idea is to replace a value of the underlying original data by a random value (from a distribution with sufficient variance). The analyses will then be performed successively with all modified data sets and the interval of the values of the results is published.

There are other approaches to result-based confidentiality. The US Census Bureau, for example, follows an approach in which automated checking of results is not obtained by presenting those results in less detail but by applying restrictions of use which are controlled by the system (Zayatz, 2007). This refers, first, to the data made available to users for analysis (e.g. showing microdata in less detail, such as by combining categories of specific variables) and, second, to a limitation of the studies that can be done with the system to a limited catalogue of analyses.

Other approaches known from the literature deal specifically with the problem of confidentiality of regression residuals (Reiter, 2004 and Sparks et al. (2008)). While Reiter proposes, among other things, to provide users with synthetically generated residuals instead of the original residuals, the approach by Sparks et al. is based on the production of box plots for the residuals.

#### **4. Summary and outlook**

The project described here forms an important bridge between the developments in improving data access channels for the scientific community over the last few years and the concepts planned already today for the future by the research data centres. It is a major milestone on the way towards real remote access.

With the current development of demand for microdata in Germany, and especially the development of access for on-site use, manual remote data execution is getting more and more difficult anyway for the capacity reasons mentioned above. Due to the increasing demand for various statistics, timely provision of scientific use files from highly different surveys is nearly impossible, too.

In the long term, real remote access seems to be the only feasible solution both nationally and internationally; all the more so as a method, once developed, can rapidly be transferred to other surveys and could allow “just in time” delivery of data. The technical developments have reached a phase where online access is possible from anywhere or will be possible soon with the relevant range.

Real remote access allows researchers to process the data independent of time and location and has the advantage that the data remain in the protected rooms (and on the protected servers) of the statistical offices. Also, that kind of data access increases the networking among researchers and the scientific transparency because any researcher may access the data and replicate results any time.

What is more, care should be taken especially in the e-science age that the development of the informational infrastructure is not left behind by technological development. The possibilities offered by the technical infrastructure are far from exhausted and have further potential for development in the future. Nevertheless legal issues, too, must be settled at this point.

The purpose of the project is, first, to provide the methodical bases for fully automated remote access. Second, it will reduce the burden on the staff of the research data centres of the statistical offices of the Federation and the Länder already in the course of the project. This can be achieved by producing data structure files and the

tools needed to produce data structure files regarding any statistics for controlled remote access as well as the guidelines and tools for categorising and automated confidentiality checks. The project will be able to take account of the methods developed in other countries and to benefit in international working groups from the experience acquired. It also benefits from the methodical projects performed in the last few years in the area of anonymisation of business microdata.

## References

- Bender, S., Himmelreicher, R., Zühlke, S., and Zwick, M. (2009) *Improvement of Access to Data from Official Statistics – the case of Germany* (in this publication)
- Bender, S., Rosemann, M., Zühlke, S., and Zwick, M. (ed.) (2008) *Betriebs- und Unternehmensdaten im Längsschnitt – Neue Datenangebote und ihre Forschungspotenziale* (Longitudinal data on local units and enterprises – New data offers and their research potential). In: *AStA – Wirtschafts- und sozialstatistisches Archiv*, Volume 2, Number 3, Springer
- Borchsenius, Lars (2005) New developments in the Danish system for access to micro data, Monographs of Official Statistics, Luxembourg
- Coder, John and Cigrang, Marc (2003) LISSY Remote Access System, working paper No.7 of the Joint ECE/Eurostat work session on statistical data confidentiality, Luxembourg
- Gießing, S. (1999) *Statistische Geheimhaltung in Tabellen* (Statistical confidentiality in tables). In: Statistisches Bundesamt (ed.): *Methoden zur Sicherung der statistischen Geheimhaltung* (Methods of ensuring statistical confidentiality), Wiesbaden, pp. 6-26
- Gomatam, S., Karr, A.F., Reiter, J.P., Sanil, A. (2005) Data dissemination and disclosure limitation in a world without microdata: A risk-utility framework for remote access servers. *Statistical Science* 20, pp. 163-177
- Heitzig, J. (2005) The “Jackknife” Method: Confidentiality Protection for complex statistical Analyses, UNECE/Eurostat Work Session on Statistical Data Confidentiality, Geneva
- Hundepool, Anco and de Wolf, Peter-Paul (2005) OnSite@Home: Remote Access at Statistics Netherlands, UNECE/Eurostat Work Session on Statistical Data Confidentiality, Geneva
- Reiter, J.P., (2004) New Approaches to Data Dissemination: A Glimpse into the Future (?). *Change* 17, pp. 12-16
- Ronning, G., Sturm, R., Höhne, J., Lenz, R., Rosemann, M., Scheffler, M. and Vorgrimler, D. (2005) *Handbuch zur Anonymisierung wirtschaftsstatischer Mikrodaten* (Manual on anonymisation of microdata of economic statistics), *Statistik und Wissenschaft*, Vol. 4/2005. Statistisches Bundesamt
- Sparks, R., Carter, C., Donnelly, J., O’Keefe, C.M., Duncan, J., Keighley, T., McAullay, D. (2008) Remote Access Methods for Exploratory Data Analysis and Statistical Modelling: Privacy-Preserving Analytics. *Comput. Methods Programs Biomed*
- Zayatz, L., (2007) New Implementations of Noise for Tabular Magnitude Data, Synthetic Tabular Frequency and Microdata, and a Remote Microdata Analysis System. Proceedings of the Joint UNECE/Eurostat work session on statistical data confidentiality, Manchester, United Kingdom, 17-19 December 2007

- Zühlke, S., Zwick, M., Scharnhorst, S. and Wende, T. (2004) The research data centres of the Federal Statistical Office and the statistical offices of the Länder, *Journal of Applied Social Science Studies* 124 (4), pp. 567 – 578
- Zwick, M. (2006), *Forschungsdatenzentren – Nutzen und Kosten einer informationellen Infrastruktur für Wissenschaft, Politik und Datenproduzenten* (Research data centres – Benefits and costs of an informational infrastructure for science, politics and data producers), *Wirtschaft und Statistik* 12, pp. 1233 – 1240