

FDZ-LIfBi

## Data Manual

NEPS Starting Cohort 5—First-Year Students  
*From Higher Education to the Labor Market*

Scientific Use File Version 14.0.0

Copyrighted Material  
Leibniz Institute for Educational Trajectories (LIfBi)  
Wilhelmsplatz 3, 96047 Bamberg  
Director: Prof. Dr. Cordula Artelt  
Executive Director of Research: Dr. Jutta von Maurice  
Executive Director of Administration: Dr. Robert Polgar  
Bamberg; May 26, 2020

## Research Data Documentation

The *NEPS Research Data Documentation Series* presents resources prepared to support the work with data from the National Educational Panel Study (NEPS).

Full citation of this document:

FDZ-LIfBi. (2020). *Data Manual NEPS Starting Cohort 5—First-Year Students, From Higher Education to the Labor Market, Scientific Use File Version 14.0.0*. Bamberg, Germany: Leibniz Institute for Educational Trajectories, National Educational Panel Study

This release of Scientific Use Data from Starting Cohort 5—First-Year Students “From Higher Education to the Labor Market” was prepared by the staff of the Research Data Center at Leibniz Institute for Educational Trajectories (Forschungsdatenzentrum, FDZ-LIfBi). It represents a major collaborative effort. *The contribution of the following persons is gratefully acknowledged:*

Eva Akins  
Dietmar Angerer  
Nadine Bachbauer  
Pia Bechtloff  
Daniel Bela  
Hannes Götz  
Daniel Fuß  
Lydia Kleine  
Tobias Koberg  
Gregor Lampel  
Sven Pelz  
Benno Schönberger  
Mihaela Tudose  
Katja Vogel  
Clara Wolf

*For their support in writing this manual, special thanks go to:*

Isabelle Fiedler, Annika Grieb, Marie Kühn, Uta Liebeskind (DZHW Hannover)

*We also appreciate the work of the former colleagues at the Research Data Center:*

Simon Dickopf, Thomas Leopold, Manuel Munz, Sebastian Pink, Marcel Raab, Jan Skopek, Knut Wenzig, Markus Zielonka

Leibniz Institute for Educational Trajectories (LIfBi)  
Research Data Center (FDZ)  
Wilhelmsplatz 3  
96047 Bamberg, Germany

E-mail: [fdz@lifbi.de](mailto:fdz@lifbi.de)  
Web: <https://www.neps-data.de/datacenter>  
Phone: +49 951 863 3511



# Contents

<b>1</b>	<b>Introduction</b>	<b>1</b>
1.1	About this manual . . . . .	1
1.2	Further documentation . . . . .	1
1.3	Data release strategy . . . . .	3
1.4	Data access . . . . .	5
1.5	Publications with NEPS data . . . . .	6
1.6	Rules and recommendations . . . . .	7
1.7	On using the Federal State label ( <i>Bundeslandkennung</i> ) . . . . .	9
1.8	User services . . . . .	9
1.9	Contacting the Research Data Center . . . . .	11
<b>2</b>	<b>Sampling and Survey Overview</b>	<b>12</b>
2.1	From higher education to the labor market . . . . .	12
2.2	Sampling strategy . . . . .	13
2.3	Competence measures . . . . .	14
2.4	Survey overview and sample development . . . . .	17
2.4.1	Wave 1: 2010/2011 (CATI+competencies) . . . . .	19
2.4.2	Wave 2: 2011 (CAWI) . . . . .	20
2.4.3	Wave 3: 2012 (CATI) . . . . .	21
2.4.4	Wave 4: 2012 (CAWI) . . . . .	22
2.4.5	Wave 5: 2013 (CATI+competencies) . . . . .	23
2.4.6	Wave 6: 2013 (CAWI) . . . . .	24
2.4.7	Wave 7: 2014 (CATI+competences) . . . . .	25
2.4.8	Wave 8: 2014 (CAWI) . . . . .	26
2.4.9	Wave 9: 2015 (CATI) . . . . .	27
2.4.10	Wave 10: 2016 (CATI) . . . . .	28
2.4.11	Wave 11: 2016 (CAWI) . . . . .	29
2.4.12	Wave 12: 2017 (CATI) . . . . .	30
2.4.13	Wave 13: 2018 (CATI) . . . . .	31
2.4.14	Wave 14: 2018 (CAWI) . . . . .	32
<b>3</b>	<b>General Conventions</b>	<b>33</b>
3.1	File names . . . . .	33
3.2	Variables . . . . .	35
3.2.1	Conventions for general variable naming . . . . .	35
3.2.2	Conventions for competence variable naming . . . . .	38
3.2.3	Labels . . . . .	41
3.3	Missing values . . . . .	42
3.4	Generated variables . . . . .	44

<b>4</b>	<b>Data Structure</b>	<b>47</b>
4.1	Overview . . . . .	47
4.1.1	Identifiers . . . . .	48
4.1.2	Panel data . . . . .	48
4.1.3	Episode or spell data . . . . .	49
4.1.4	Revoked episodes . . . . .	51
4.2	Data files . . . . .	52
4.2.1	Basics . . . . .	54
4.2.2	Biography . . . . .	56
4.2.3	CohortProfile . . . . .	58
4.2.4	EditionBackups . . . . .	60
4.2.5	Education . . . . .	62
4.2.6	MethodsCATI . . . . .	64
4.2.7	MethodsCAWI . . . . .	66
4.2.8	MethodsCompetencies . . . . .	68
4.2.9	pTargetCATI . . . . .	70
4.2.10	pTargetCAWI . . . . .	72
4.2.11	pTargetMicrom . . . . .	74
4.2.12	spChild . . . . .	76
4.2.13	spChildCohab . . . . .	78
4.2.14	spCourses . . . . .	80
4.2.15	spEmp . . . . .	82
4.2.16	spFurtherEdu1 . . . . .	84
4.2.17	spFurtherEdu2 . . . . .	86
4.2.18	spGap . . . . .	88
4.2.19	spInternship . . . . .	90
4.2.20	spMilitary . . . . .	92
4.2.21	spParLeave . . . . .	94
4.2.22	spPartner . . . . .	96
4.2.23	spSchool . . . . .	98
4.2.24	spSchoolExtExam . . . . .	100
4.2.25	spSibling . . . . .	102
4.2.26	spUnemp . . . . .	104
4.2.27	spVocExtExam . . . . .	106
4.2.28	spVocPrep . . . . .	108
4.2.29	spVocTrain . . . . .	110
4.2.30	Weights . . . . .	112
4.2.31	xEcoCAPI . . . . .	114
4.2.32	xInstitution . . . . .	116
4.2.33	xPlausibleValues . . . . .	118
4.2.34	xTargetCompetencies . . . . .	120
<b>5</b>	<b>Special Issues</b>	<b>122</b>
5.1	Service Variables (Area of studies, ISCED-97 subject) . . . . .	122

5.2	Coding subject of study . . . . .	123
5.2.1	Recruitment . . . . .	123
5.2.2	Panel Waves . . . . .	123
<b>A</b>	<b>Appendix</b>	<b>127</b>
A.1	R examples . . . . .	127
A.2	Release notes . . . . .	156
A.3	Comparison of _v1 variables . . . . .	163

# 1 Introduction

## 1.1 About this manual

This manual is intended to facilitate your work with data of NEPS Starting Cohort 5—First-Year Students (NEPS SC5). It serves both as a first guide for getting started with the complex data and as a reference book. The primary emphasis is on practical aspects such as sample development, data structure, and variable merging. The manual is neither complete nor exhaustive, but several links to other resources are provided in the respective paragraphs.

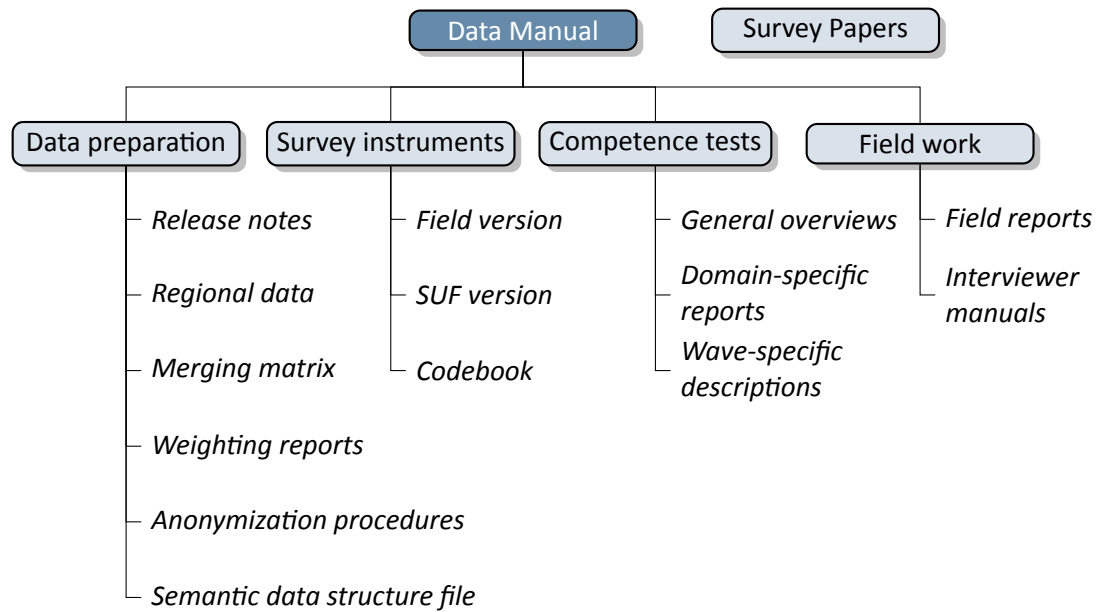
The first chapter refers to further documentation material, requirements for data access, instructions for data citation, some general rules and recommendations, and selected user services. In the second chapter, the fundamental objectives of Starting Cohort 5 and its sampling strategy are briefly introduced. The main part of this chapter is devoted to the sample development across the waves including field times, realized case numbers, survey modes, and the measurement of competency domains. The principles of Scientific Use File data-editing processes as well as conventions for naming the data files and variables are explained in the third chapter, supplemented by missing value definitions and an overview of additionally generated variables. The fourth chapter focuses on the data structure with information about data types, identifiers, and short portraits of all available datasets in the Scientific Use File. These short portraits include recommendations on how to use the dataset as well as syntax examples for merging variables of this dataset with variables from other files. The last chapter addresses some specific issues that should be noted when working with data of Starting Cohort 5.

According to the cumulative release strategy—each new Scientific Use File contains the data of all previous survey waves plus the data of the currently prepared wave(s)—this manual will be regularly updated and revised. While the given information remain valid over time, at least the sample development has to be continuously complemented. In other words, the latest published manual replaces the previous ones. All relevant adjustments and extensions in future releases of this manual will be listed in a separate appendix.

## 1.2 Further documentation

The data manual cannot cover all issues in detail. Hence, a bunch of supplementary reports and other materials with background information on data preparation, survey instruments, competence tests, and field work is offered (see Figure 1). This frequently updated and enhanced data documentation can be downloaded from our website at:

→ [www.neps-data.de](http://www.neps-data.de) > Data Center > Data and Documentation  
    > Starting Cohort First-Year Students > Documentation



**Figure 1:** NEPS supplementary data documentation

**Release notes** All Scientific Use Files are accompanied by release notes that log changes in the data compared to prior versions and list bugs eliminated or at least known. For the latter, short syntax corrections are usually given. Please consult these notes when working with the data. See also section A.2 for a depiction of the current notes.

**Regional data** Fine-grained regional indicators from a commercial provider (microm) are available in our On-site environment. The report describes the regional levels covered by these indicators, their content, and how to merge them to the survey data.

**Merging matrix** This matrix provides an overview of how to link information from different datasets, taking into account the relevant identifier variables.

**Weighting reports** These reports entail information regarding the design principles of the sampling process and the creation of weights.

**Anonymization procedures** The document describes the anonymization measures carried out and provides an overview regarding the opportunity to access sensitive data.

**Semantic data structure file** This data package corresponds to the Scientific Use File but does not contain any observations (*purged datasets*). It provides all metadata including variable names, labels and answering scheme options to be used for exploring the data structure and for preparing analyses.

**Survey instruments** For each wave, the survey instruments are offered in the form of Scientific Use File (SUF) and field versions. While the field versions consist of the originally deployed instruments (in German only), the SUF versions are enriched by additional information



such as variable names and value labels used in the Scientific Use File. *Please note, that the competence test booklets are not publicly available.*

**Codebook** The codebook lists all variables and their corresponding labels plus the basic frequencies by waves in concordance with the datasets in the Scientific Use File.

**Competence tests** Information about competence testing is provided in various documentations, including general overviews and wave-specific descriptions; also, for each domain there is usually a brief description of the construct with sample items, a description of the data, and of the psychometric properties of the test.

**Field reports** The field reports document the overall data-collection process conducted by the survey institute. This information about survey preparation, interviewer deployment, respondent tracking, initial contacts, incentives, and sample realization is available in German only.

**Interviewer manuals** The interviewer manuals are a collection of instructions for the interviewers. In particular, they exemplify the interview process as well as the content of each of the questionnaire modules. They are available in German only.

**NEPS Survey Papers** Finally, there is a series of NEPS Survey Papers that address several topics of more general interest. These papers are listed for download on our website at:

→ [www.neps-data.de](http://www.neps-data.de) > Data Center > Publications > NEPS Survey Papers

Additional documentation material might be available for specific cohorts and/or waves. Please visit the website above for further details.

### 1.3 Data release strategy

NEPS data are published in the form of Scientific Use Files. They are provided free of charge to the scientific community. Each Scientific Use File consists of multiple datasets, forming a complex data structure with cross-sectional, panel and episode or spell information (see section 4). The release of NEPS Scientific Use Files follows a cumulative strategy, i. e., the latest data release replaces all former data releases. Hence, we recommend to use the most current release of a Scientific Use File.

#### File Format

All Scientific Use Files are disseminated in Stata and SPSS format with bilingual variable labels and value labels in German and English. In the SPSS format, there are separate data files for both languages. Data stored in Stata format contain both languages within one file; the switch is induced by the following Stata command:

```
label language [de/en]
```

Due to the change of encoding to “Unicode” in Stata14 and the fact that older Stata versions are not able to open such data files, the NEPS Scientific Use Files contain two Stata formats, namely Stata14 and Stata12.

### Versioning and Digital Object Identifier

Every time a new Scientific Use File is released, the data files existing up to now are either extended, usually by information from a new survey wave, or updated with changes due to larger or smaller corrections. The three digits of the version number inform about the number of waves integrated in the specific Scientific Use File, the frequency of major updates, and the frequency of minor updates. The version number is part of all relevant designations: that of the Scientific Use File, its data files (see Table 3), and the respective Digital Object Identifier.

Every release of a NEPS Scientific Use File is registered at [data.neps.gesis.org](https://data.neps.gesis.org) and clearly labeled with a unique Digital Object Identifier (DOI, cf. Wenzig, 2012). This DOI has two main functions. On the one hand, it enables researchers to cite the utilized NEPS data in an easy and precise way (see section 1.5). This in turn is a basic precondition for any replication analysis. On the other hand, the DOI directs to a landing page with further information about the Scientific Use File and the data access options. The DOI of the current release is `doi:10.5157/NEPS:SC5:14.0.0`. Other releases of Scientific Use Files for Starting Cohort 5 can be accessed by substituting the version number at the end of the DOI and the URL respectively (see Table 1).

**Table 1:** Release history of SUF in Starting Cohort 5

SUF Version	DOI	Date of release
<b>14.0.0</b> (current)	<code>doi:10.5157/NEPS:SC5:14.0.0</code>	<b>2020-05-27</b>
13.0.0	<code>doi:10.5157/NEPS:SC5:13.0.0</code>	2020-02-14
12.0.0	<code>doi:10.5157/NEPS:SC5:12.0.0</code>	2019-07-26
11.0.0	<code>doi:10.5157/NEPS:SC5:11.0.0</code>	2018-09-06
10.0.0	<code>doi:10.5157/NEPS:SC5:10.0.0</code>	2018-04-19
9.0.0	<code>doi:10.5157/NEPS:SC5:9.0.0</code>	2017-06-23
8.0.0	<code>doi:10.5157/NEPS:SC5:8.0.0</code>	2016-12-23
6.0.0	<code>doi:10.5157/NEPS:SC5:6.0.0</code>	2016-03-31
4.0.0	<code>doi:10.5157/NEPS:SC5:4.0.0</code>	2014-09-30
3.1.0	<code>doi:10.5157/NEPS:SC5:3.1.0</code>	2014-05-16
3.0.0	<code>doi:10.5157/NEPS:SC5:3.0.0</code>	2013-07-05

### 1.4 Data access

Access to the NEPS data is free of charge but limited to the purpose of research and members of the scientific community. Granting the right to obtain the data requires the conclusion of a Data Use Agreement. The existence of a valid Data Use Agreement entitles to work with all NEPS Scientific Use Files, i. e., the full data portfolio is at the disposal of the data recipient and all further persons involved in the agreement.

#### Application for data access

- Fill in the online form for a NEPS Data Use Agreement either in German or in English. Enter a title, the duration, and a short description of the intended research project. Make sure that all project participants with NEPS data access are specified in the form and that these persons have signed the agreement. Submit one copy of the complete agreement by e-mail, fax, or post. Further instructions and the relevant forms are provided on our website at:

→ [www.neps-data.de](http://www.neps-data.de) > Data Center > Data Access > Data Use Agreements

- After approval by the Research Data Center, the registered NEPS data user receives a user name and a password to log in to our website.
- The basic Data Use Agreement permits the download of all available Scientific Use Files from our website at:  
→ [www.neps-data.de](http://www.neps-data.de) > Data Center > Data and Documentation > NEPS Data Portfolio
- There are two other modes of access to the NEPS data (see below); each demanding a supplemental agreement in addition to the basic Data Use Agreement.
- Another form is provided to state changes of the Data Use Agreement regarding further project participants or a prolonged project duration.

#### Modes of data access

Three modes of accessing the NEPS Scientific Use Files are available. They are designed to support the full range of researchers' interests and maximize data utility while complying with national and international standards of confidentiality protection. Each modus corresponds to a data version that is different with regard to the accessibility of sensitive information as the three versions of a Scientific Use File vary according to their level of data anonymization.

- *Download* from the website = highest level of anonymization
- *RemoteNEPS* as browser-based remote desktop access = medium level of anonymization
- *On-site* access at secure working stations at LfBi = lowest level of anonymization

While working with RemoteNEPS requires a biometrical authentication and Internet access, the On-site use of NEPS data necessitates a guest stay at LIfBi in Bamberg. More details about the three access modes and their implications for application and utilization are given on our website at:

→ [www.neps-data.de](http://www.neps-data.de) > Data Center > Data Access

### Sensitive information

The download version of a Scientific Use File contains the least amount of information. For instance, institutional context data and the Federal State label (*Bundeslandkennung*, see section 1.7) are only available in the controlled environments of RemoteNEPS and our On-site data security rooms. Indicators of a certain sensitivity are modified in the download data, such as aggregated categories for countries of citizenship or languages of origin. A few datasets and variables are exclusively accessible in the On-site version, e.g. the fine-grained regional indicators or open text entries. For a full picture of the availability of sensitive information, please refer to the overview on our website at:

→ [www.neps-data.de](http://www.neps-data.de) > Data Center > Data Access > Sensitive Information

The hierarchical concept of data dissemination translates into an onion-shaped model of datasets. The most sensitive on-site level represents the outer layer with the remote and download levels being subsets of these data. That is, any data contained within a less sensitive level are also included in the higher level(s). A detailed list of variables offered at the different levels together with notes on the generation of the three data versions can be found for each release of a Scientific Use File in the respective report on "Anonymization Procedures".

## 1.5 Publications with NEPS data

Referencing the use of data from the National Educational Panel Study (NEPS) is essential for a good scientific practice as well as for revealing the scientific value of this study. The following citation rules apply to all publications based on NEPS data of Starting Cohort 5.

It is obligatory to acknowledge the NEPS study in general and to indicate the utilized data version by including a phrase like this in your publication:

This paper uses data from the National Educational Panel Study (NEPS): Starting Cohort 5—First-Year Students, doi:10.5157/NEPS:SC5:14.0.0. From 2008 to 2013, NEPS data was collected as part of the Framework Program for the Promotion of Empirical Educational Research funded by the German Federal Ministry of Education and Research (BMBF). As of 2014, NEPS is carried out by the Leibniz Institute for Educational Trajectories (LIfBi) at the University of Bamberg in cooperation with a nationwide network.

Please also add these bibliographic details to your list of references:

Blossfeld, H.-P., Roßbach, H. G., & von Maurice, J. (Eds.). (2011). Education as a Lifelong Process: The German National Educational Panel Study (NEPS). *[Special Issue] Zeitschrift für Erziehungswissenschaft: 14.*

Authors of any kind of publications based on NEPS data are requested to notify the Research Data Center about their articles and to provide an electronic version or a special print or a copy. All reported publications are listed in the NEPS Bibliography on our website at:

→ [www.neps-data.de](http://www.neps-data.de) > Data Center > Publications

### Citing documentation

To refer to any of the documentation material published in the *NEPS Research Data Documentation Series* (e. g. this manual), please make use of the following citation templates:

FDZ-LifBi. (2020). *Data Manual NEPS Starting Cohort 5– First-Year Students, From Higher Education to the Labor Market, Scientific Use File Version 14.0.0.* Bamberg, Germany: Leibniz Institute for Educational Trajectories, National Educational Panel Study

Or another example:

Schönberger, K., & Koberg, T. (2017). *Regional Data: Microm.* Bamberg, Germany: Leibniz Institute for Educational Trajectories, National Educational Panel Study

If no author is given, please take a universal *NEPS* instead:

NEPS (Ed.). (2020). *Starting Cohort 5: First-Year Students (SC5), Wave 14, Questionnaires (SUF Version 14.0.0).* Bamberg, Germany: Leibniz Institute for Educational Trajectories, National Educational Panel Study

If a document has not been published in this series, please refer to the author and the title as in the following citation of a field report by one of our survey institutes:

Steinwede, J., & Aust, F. (2012). *Methodenbericht, NEPS Startkohorte 5 – CATI-Haupterhebung Herbst 2010, B52.* Bonn, Germany: infas

## 1.6 Rules and recommendations

Working with NEPS data is bound to a couple of rules that are codified in the Data Use Agreement. Each data user has to confirm these rules by his or her signature. The already mentioned obligation to cite the study and to indicate any kind of publication resulting from the use of NEPS data (see section 1.5) are just two examples. The major part of rules refers to issues of data privacy and the requirements of a careful data handling.

### Rules

- *Avoidance of re-identification:* Any action aimed at and suitable for re-identifying persons, households, or institutions is strictly forbidden. This also includes the combination of NEPS data with other data that allow for a re-identification of persons. In case of any accidental re-identification, the Research Data Center has to be informed immediately and all individual data gained therefrom have to be kept secret.
- *Avoidance of data disclosure:* NEPS data are exclusively provided on the basis of a valid Data Use Agreement—for a defined purpose (research project) and to a defined group of persons (data recipient and further project members that are involved in the contract). Any use for commercial or other economic purposes is not permitted just as any transfer of the data to third parties. Please handle the provided NEPS data with strict confidentiality!
- *Regulations on using the Federal State label:* For NEPS data collected in connection with schools or higher education institutions it is not allowed to use Federal-State-related information directly or indirectly contained in the data for analyses aiming at direct comparisons of the German Federal States (*Bundesländer*), or aiming at direct conclusions to be drawn about a Federal State, or aiming at a reconstruction of the concrete Federal State affiliation of persons, households, and institutions. Any kind of ranking between the Federal States based on NEPS data is prohibited (see section 1.7).

Please note that violation of these rules may lead to severe penalties as stated in the NEPS Data Use Agreement. If there is any doubt or question regarding the given regulations, please contact the Research Data Center (see section 1.9). The same applies in case of encountering any deficiencies in data quality or any security leaks with regard to NEPS data protection and data security.

### Recommendations

In addition to the binding rules, there are some recommendations for the use of NEPS data:

- *As a matter of course:* Always be critical when working with empirical data! Although a big effort is being made to ensure the integrity of the provided data we cannot guarantee absolute correctness. Notices on problems or errors in the data are welcome at any time at the Research Data Center.
- *Enhanced understanding of the data:* Consult the documentation and survey instruments! The analysis of complex data necessitates a precise idea of how the information were collected and edited. All relevant material is available online (see section 1.2).
- *Facilitated handling of the data:* Utilize the tools that are offered! Several user services are provided to support NEPS data analyses—reaching from specific Stata commands (e. g., for an easy and adequate recoding of missing values) to a meta search engine (e. g., for an interactive exploration of all instruments) to a discussion forum (e. g., for the clarification of questions). These tools are also available online, see section 1.8 for more details.

### 1.7 On using the Federal State label (*Bundeslandkennung*)

In concurrence with the regulations of the Research Data Center at the Institute for Educational Quality Improvement (Institut zur Qualitätsentwicklung im Bildungswesen, IQB), using the Federal State label in conjunction with NEPS data collected in connection with schools or higher education institutions is permitted in the context of exploring scientific research questions, if it is exclusively used for:

- control purposes in order to incorporate it as a covariate in the planned analysis; the identification of individual Federal States in the displayed results is not permitted
- incorporating contextual characteristics or other third-party variables; the identification of individual Federal States in the displayed results is not permitted
- comparing aggregated groups of Federal States where at least two states are combined to form a single meaningful group with regard to substantive issues; the identification of individual Federal States in the displayed results is not permitted
- for sample descriptions (e.g., the distribution of participants by state and by different types of schools within states)

When using data collected in connection with schools or higher education institutions, it is not allowed to use Federal-State-related information directly or indirectly contained in the data for analyses aiming at a direct Federal State comparison, direct conclusions to be drawn about a Federal State, or a reconstruction of the concrete Federal State affiliation of persons, households, and institutions.

The Federal State label in the starting cohorts of schools and higher education institutions is provided by LIfBi to the scientific community only via remote access (RemoteNEPS) and—depending on availability—via guest working stations in Bamberg (On-site). The respective analysis results are reviewed by LIfBi to ensure that this agreement has been observed before being passed on electronically to the researcher in a password-protected environment. The abovementioned restrictions concerning the use of the Federal State label do not apply to data collected in a nonschool context and/or in Federal-State-specific educational reform studies.

### 1.8 User services

In addition to a comprehensive data documentation there are several user services to support researchers working with NEPS data. First and foremost, the Research Data Center maintains a regularly updated and enhanced website with detailed information on all available Scientific Use Files, a complete list of NEPS projects, a NEPS bibliography, a reference to NEPS events, and a NEPS newsletter. All subsequently introduced services and tools can also be reached via this website:

→ [www.neps-data.de](http://www.neps-data.de) > NEPS

### NEPSforum

The *NEPSforum* is an open online discussion platform for experienced users as well as for persons who are searching for NEPS related information. It offers the opportunity to exchange with NEPS staff members and with other researchers in a transparent dialogue. That way, the forum will become a rich archive of knowledge with practical solutions for numerous problems and questions. We highly encourage you to browse the forum first when struggling with NEPS issues or when help is needed with specific data matters. If there is no available solution, please take the opportunity to share your question by posting it to the forum. Active participation requires no more than a one-time registration. The entire NEPS user community will benefit from a broad participation. You can find the *NEPSforum* at:

→ [www.neps-data.de > Data Center > NEPSforum](http://www.neps-data.de/Data_Center/NEPSforum)

### NEPSplorer

The *NEPSplorer* facilitates an interactive and quick full text search through all instruments of released NEPS surveys, with the exception of competence tests. The tool is particularly suitable for getting a first idea of the availability of constructs, items, and variables in the datasets. It is currently based on keyword search with several filtering options, but a hierarchical construct search will be added soon. The *NEPSplorer* offers some helpful functions such as displaying univariate statistics, listing relevant metadata, and enabling registered users to create their own personal watch list of interesting items. As a web application—a mobile version aligned for smartphone usage is also available—the *NEPSplorer* relies on the most up-to-date information; any correction in the metadata is thus instantly visible. Start the tool here:

→ [www.neps-data.de > Data Center > Overview and Assistance > NEPSplorer](http://www.neps-data.de/Data_Center/Overview_and_Assistance/NEPSplorer)

### NEPStools

*NEPStools* is a free to use collection of Stata commands that is created and supplied by the Research Data Center. The package includes some programs (“ado files”) that make NEPS data handling easier. As an example, the `nepsmiss` command automatically recodes all of the numeric missing values (-97, -98, etc.) into Stata’s “Extended Missings” (.a, .b, etc.) with correctly recoded value labels. Another example is the `infoquery` command that displays additional attributes of the variable such as the question text and the initial variable name in the instrument. The *NEPStools* set can be easily installed from our repository through Stata’s built-in installation mechanism:

```
net install nepstools, from(http://nocrypt.neps-data.de/stata)
```

A description of the programs and further information are given on the website at:

→ [www.neps-data.de > Data Center > Overview and Assistance > Stata Tools](http://www.neps-data.de/Data_Center/Overview_and_Assistance/Stata_Tools)



### User trainings

The Research Data Center offers a series of regular user training courses at the Leibniz Institute for Educational Trajectories in Bamberg. The standard 2-day courses are free of charge. On the first day, there is a general introduction to the design of the NEPS study, the structure of NEPS Scientific Use Files, the terms and conditions of data access and data usage, and the handling of documentation materials. The second day is more focused on data of a certain starting cohort and on selected methodological and/or theoretical concepts. Both parts come along with guided hands-on sessions. A crucial aspect of all user trainings is the sensitization of participants to issues of privacy and data protection. In this context, participation is obligatory for those who want to enroll in the biometric authentication system in order to gain access to the NEPS remote or On-site environment. A schedule of all training dates together with information on how to register for a course can be retrieved from our website at:

→ [www.neps-data.de](http://www.neps-data.de) > Data Center > User Training

### 1.9 Contacting the Research Data Center

The Research Data Center at the Leibniz Institute for Educational Trajectories (Forschungsdatenzentrum, FDZ-LIfBi) accounts for large parts of the NEPS data preparation and documentation process, for the data dissemination, and for the user support including individual advice. We welcome your feedback at any time to further improve our products and services. This particularly applies to this manual as the guiding document to facilitate your work with NEPS data of Starting Cohort 5.

Please contact us with your questions, comments, requests, and suggestions:

E-mail: [fdz@lifbi.de](mailto:fdz@lifbi.de)

Web: → [www.neps-data.de](http://www.neps-data.de) > Data Center > Contact Data Center

Phone: +49 951 863 3511

## 2 Sampling and Survey Overview

### 2.1 From higher education to the labor market

German higher education system has been facing a number of challenges and developments since the early 2000ies, that raised new issues for research. To name but a few, there is the introduction of a two-stage structure in higher education according to the Bologna Process, a growing demand for outcome orientation, the evolution of higher education towards lifelong learning, an increase of (international) competitiveness, and the emerging shortage of highly qualified professionals. At the same time, key issues remained core challenges for the higher education system, such as student dropouts, social selectivity in university entrance, and the relationship between higher education and working life. In order to answer research questions associated with these issues, a cohort of first-year students was followed through their years of study since winter term 2010/11, including their entrance into working life. Central issues to be studied are educational choices, the outcomes of university education, and the entry into the job market.

The main focus is on

- Educational choices during the course of studies and success in studies: What are the determinants of educational decisions and success in studies while studying at a higher education institution – such as dropping out, changing subjects, studying abroad, and pursuing a Master's degree? What is the importance of competencies and social factors, such as social background, gender or migration experiences in this process? Which consequences do decisions have for subsequent education and working life?
- Entrance into working life and professional success: When thinking about students' transition into the job market and their professional success (e.g., occupational position, income, employment security), how important are acquired competencies, on the one hand based on formal qualifications (diplomas), social background, gender, and on the other hand based on social and cultural capital? What role do general competencies play in comparison to subject-specific ones?
- Students' competencies: Which general competencies do students possess to crucial points of time in their students' and young adults' lifecourse (beginning of studies, end of studies/labour market entry)? How does the competence level influence transitions during studies and beyond (change of subject, higher education drop out, transition to the labour market)? How do competencies correlate with learning environments provided by higher education institutions?

### 2.2 Sampling strategy

The target population of Starting Cohort 5 is defined as all first-year students of the academic year 2010/2011, independent of their nationality and their knowledge of the German language, who are:

- enrolled for the first time in a public or state-approved institution of higher education in Germany
- aiming at a Bachelor's degree or a state examination (Staatsexamen) in medicine, law, pharmacy, and teaching, or a diploma or Master's degree in Roman Catholic or Protestant theology or specific art and design degrees
- not attending higher education institutions run by Federal Ministries or Federal States for members of their public services (e. g., University of Applied Labour Studies/Hochschule der Bundesagentur für Arbeit)

The sampling process was designed to incorporate an oversampling of teacher education students and students at private higher education institutions. For that reason, a stratified cluster approach has been applied. Administrative data provided by the Federal Statistical Office of Germany constituted the corresponding sampling frame. Each cluster referred to the total of students enrolled in a certain subject at a particular higher education institution (e. g., social sciences at the University of Bamberg). On the primary level, the stratification differentiated between the following four strata; on the secondary level these strata were combined with groups of related subjects:

- clusters linked to teacher education at public universities
- clusters linked to all other fields of study at public universities
- clusters linked to all fields of studies at public universities of applied sciences (Fachhochschulen)
- clusters linked to all degree programs at private higher education institutions

In a second step, all institutions of selected clusters were contacted by the survey agency in order to gain access to the students. The administration of 261 institutions declared their co-operativeness, thereof 104 public universities, 108 public universities of applied sciences, and 49 private university institutions.

In the subsequent recruitment process, two different modes of contact were employed to approach the students and to receive their consent to participate in the panel study:

- conventional mail via higher education institutions administration
- personal information in lectures for freshmen students in the selected fields of studies via interviewers

The former strategy has been applied at all sampled institutions. Recruiting questionnaires in prepared envelopes were transferred to the university administrations together with detailed instructions on how to select the targeted student population. Part of this instruction was the request to include all non-traditional first-year students, i. e., all students with a higher education admission other than the general higher education certificate (Abitur or Fachabitur). It was the task of the higher education institution to compile the respective postal addresses and to send the letters plus reminder letters. Altogether 16,887 filled questionnaires were sent back to the survey agency. The latter strategy presupposed the explicit agreement by the higher education institution and the lecturer to recruit students in appropriate freshmen courses by professional interviewers. In the course of 299 visits at 99 higher education institutions, another 17,229 filled questionnaires could be collected. While the two strategies were conducted parallel during the winter semester 2010/2011, a simplified procedure was applied in the summer semester 2011. Based on postal distribution and display of reduced questionnaires, so-called NEPS address cards, additional 4,169 contact information were gathered.

The returned information of all 38,285 persons were then checked with regard to the belonging to the target population, the existence of double recruitments, and the quality of provided contact details. Finally, 21,438 cases were administrated in the first CATI survey wave of Starting Cohort 5. This first CATI was the prerequisite for staying in the panel.

The sampling design and its consequences for the derivation of sampling weights are fully described in Zinn, Steinhauer, and Aßmann, 2017. Further remarks on the recruiting process are given in the CATI field report of the first survey wave (in German only). Both documents are available on our website at:

→ [www.neps-data.de](http://www.neps-data.de) > Data Center > Data and Documentation  
    > Starting Cohort First-Year Students > Documentation

### 2.3 Competence measures

The collection and provision of data on the development of competencies and skills throughout the life course is a key element of the German National Educational Study (NEPS). Competence measurements are carried out across different waves in all NEPS starting cohorts covering domain-general and domain-specific cognitive competencies as well as metacompetencies and stage-specific competencies.

Data from the competence tests pass through an editing process before they get integrated into the Scientific Use File. This data preparation enables users to work with scored items and test scores such as the sum or mean of correct answers. Detailed descriptions on how these scores were estimated can be found in separate reports for the respective competence domains (see section 1.2). The scores are compiled in a dataset named `xTargetCompetencies`. This dataset is structured in the so-called wide format, that is, all responses of a single respondent are represented in one row of the data matrix. As a consequence, variable names for competence scores follow a specific nomenclature. It not only allows for the identification of the

respective domain, the target group, the testing modus, and the kind of scoring, but also informs about the repeated administration of a test item in a different wave or starting cohort (see section 3.2.2).

The next table shows the schedule of competence measures in Starting Cohort 5 with domains by waves including test modus. The overview contains released data as well as data that is not yet published.

**Table 2:** Schedule of competence measures. P = Paper-Based Test (proctored), C = Computer-Based Test (proctored), W = Web-Based Test (unproctored)

		2011 <b>Wave 1</b> (2nd Sem.)	2013 <b>Wave 5</b> (6th Sem.)	2014 <b>Wave 7</b> (7th Sem.)	2017 <b>Wave 12</b> (13th Sem.) <sup>3</sup>
<b>Domain-General Competencies</b>					
DGCF: Cognitive Basic Skills	dg	—	P, C, W	—	—
<b>Domain-Specific Competencies</b>					
Reading Competence <sup>1</sup>	re	P	—	—	C, W
Reading Speed	rs	P	—	—	—
Mathematical Competence <sup>1</sup>	ma	P	—	—	C, W
Scientific Competence <sup>1</sup>	sc	—	P, C, W	—	—
<b>Metacompetencies</b>					
ICT Literacy <sup>1</sup>	ic	—	P, C, W	—	—
<b>Stage-Specific Competencies</b>					
Business Administration and Economics <sup>2</sup>	ba	—	—	P	—
English Reading Competence <sup>1</sup>	ef	—	—	—	C, W

<sup>1</sup> Subsequent to the respective competence test the target persons had to assess their own test performance (Procedural Metacognition, mp).

<sup>2</sup> Reduced testing: In wave 7, the stage-specific competence test (ba) was realized in a subsample of students and graduates of business sciences only.

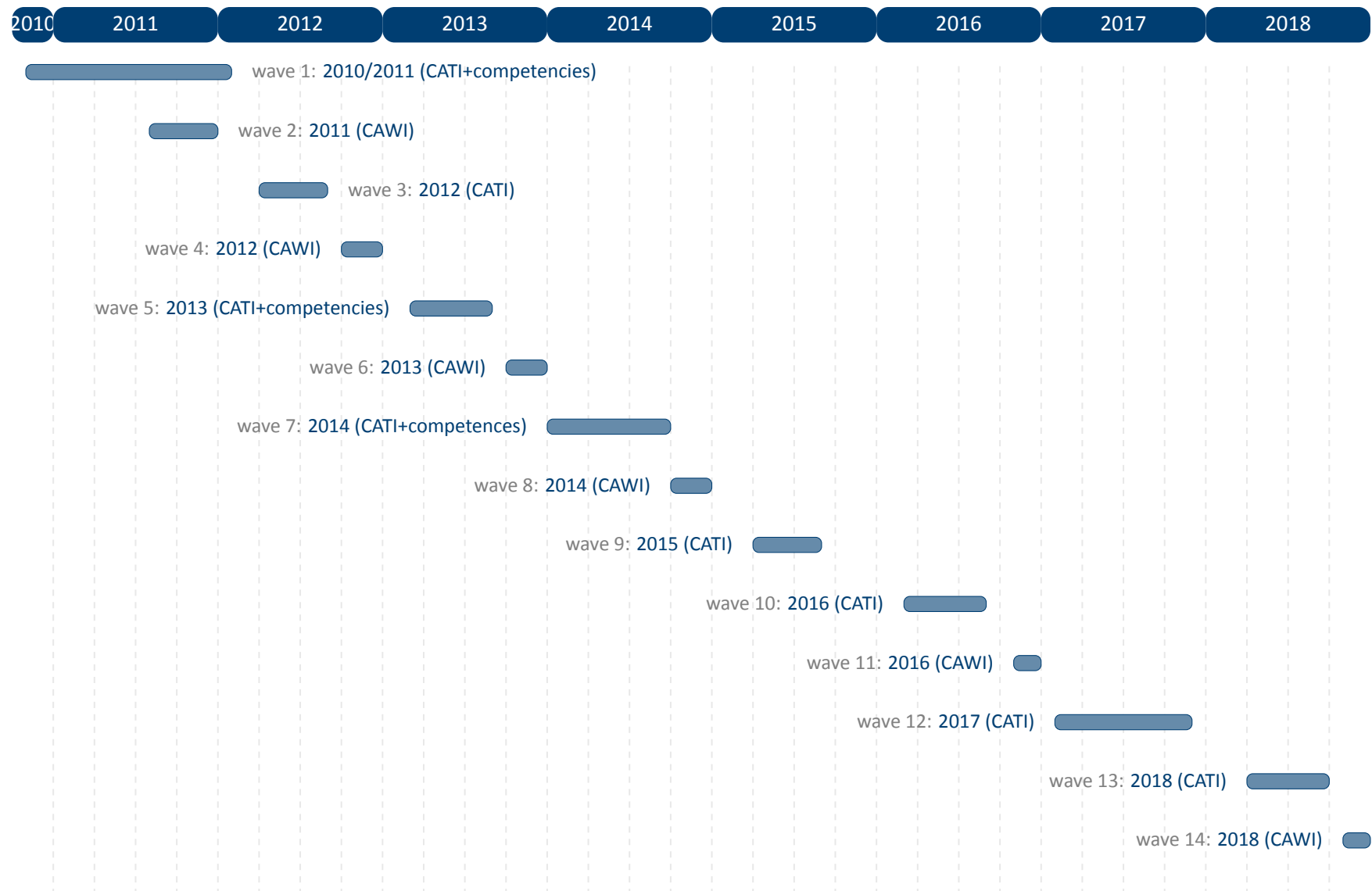
<sup>3</sup> Reduced testing: In wave 12, a randomized allocation of competence tests with two out of the three domains (re, ma or re, ef or ma, ef) has been applied.

### 2.4 Survey overview and sample development

This section informs about the progress of the Starting Cohort 5 sample. For each survey wave included in the current Scientific Use File there is a short characterization in terms of field time, number of realized cases, relevant subsamples and domains of competence testing (if appropriate), survey modus, and the institution(s) responsible for collecting the data. A more detailed insight into all relevant field work issues is provided by the *Field Reports* of the survey institutes, which are available on the website (in German only) as part of the data documentation for each (sub-)study:

→ [www.neps-data.de](http://www.neps-data.de) > Data Center > Data and Documentation  
    > Starting Cohort First-Year Students > Documentation

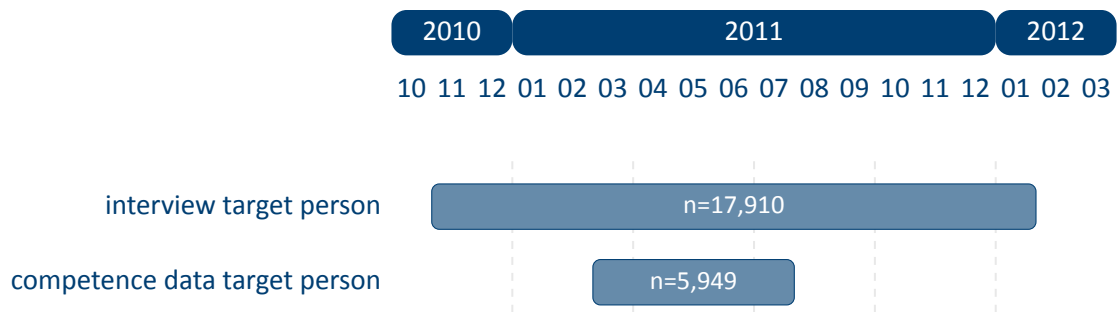
Figure 2 starts with an overview illustrating the panel progress of Starting Cohort 5 in terms of field times and survey modes from wave 1 to 14.



**Figure 2:** Survey progress of Starting Cohort 5 (waves 1 to 14)



### 2.4.1 Wave 1: 2010/2011 (CATI+competencies)



**Figure 3:** Field times and realized case numbers in wave 1

- Target persons

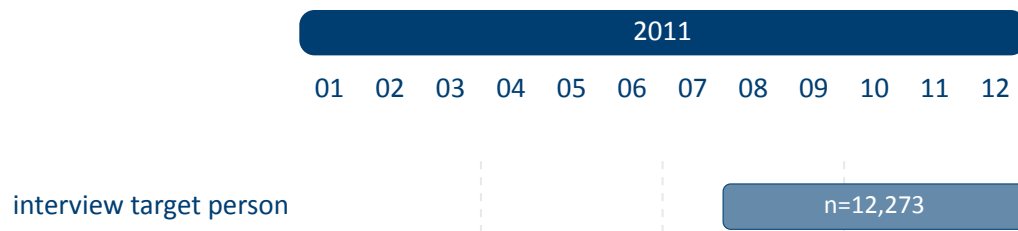
**Sample** First-year students in winter semester 2010/11 (for details about the sampling strategy, see section 2.2)

**Competence tests** Reading Competence, Reading Speed, Mathematical Competencies

**Data collection** infas – Institute for Applied Social Sciences, Bonn

**Mode of survey** Written questionnaires (in each case for recruiting and competence test, PAPI) and computer-assisted telephone interview (CATI)

### 2.4.2 Wave 2: 2011 (CAWI)



**Figure 4:** Field times and realized case numbers in wave 2

- Target persons

**Sample** Survey with the participants of the main survey 2010/2011 additional to CATI-survey

**Data collection** DZHW - German Centre for Higher Education Research and Science Studies,  
Hannover

**Mode of survey** Online survey (CAWI)

### 2.4.3 Wave 3: 2012 (CATI)



**Figure 5:** Field times and realized case numbers in wave 3

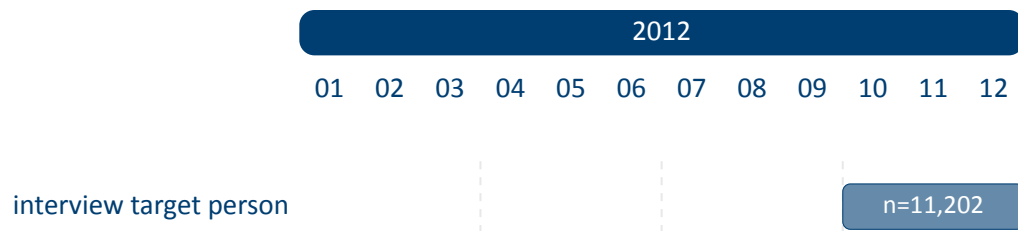
- Target persons

**Sample** Panel sample. Follow-up survey with interviewees willing to participate in the panel.

**Data collection** infas – Institute for Applied Social Sciences, Bonn

**Mode of survey** Computer-assisted telephone interview (CATI)

### 2.4.4 Wave 4: 2012 (CAWI)



**Figure 6:** Field times and realized case numbers in wave 4

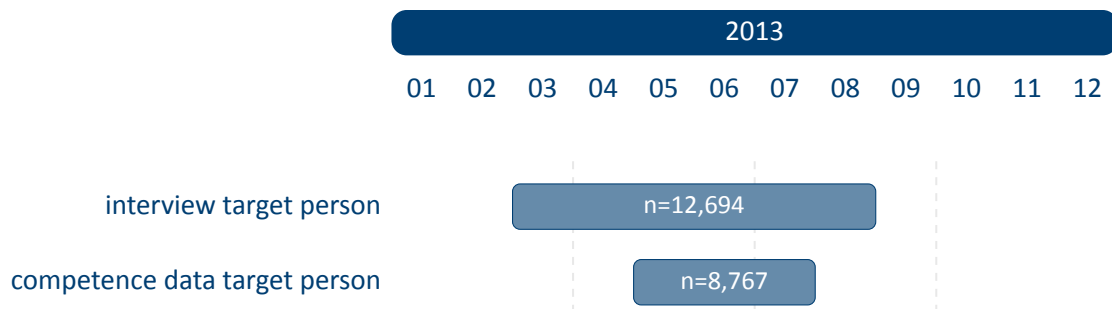
- Target persons

**Sample** Panel sample. Follow-up survey with interviewees willing to participate in the panel.

**Data collection** DZHW - German Centre for Higher Education Research and Science Studies, Hannover

**Mode of survey** Online survey (CAWI)

### 2.4.5 Wave 5: 2013 (CATI+competencies)



**Figure 7:** Field times and realized case numbers in wave 5

- Target persons

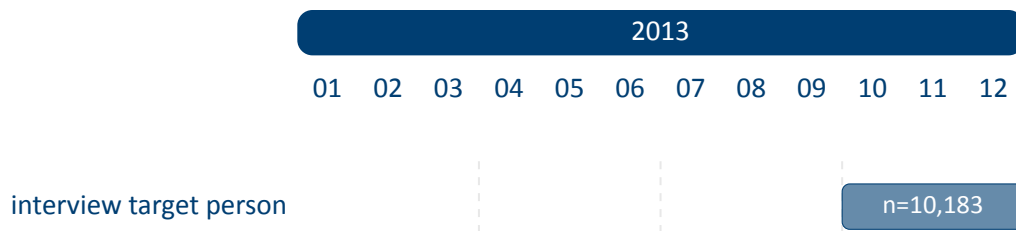
**Sample** Panel sample. Follow-up survey with interviewees willing to participate in the panel.

**Competence tests** DGCF (Cognitive Basic Skills), Scientific Competence, ICT Literacy

**Data collection** infas – Institute for Applied Social Sciences, Bonn

**Mode of survey** Computer-assisted telephone interview (CATI) and group testing (conventional paper-based testing (PAPI), paper-based testing with electronic pens (E-Pen) or computer-based testing with notebooks (CBA)) or individual testing (computer-based online testing, CBWA)

### 2.4.6 Wave 6: 2013 (CAWI)



**Figure 8:** Field times and realized case numbers in wave 6

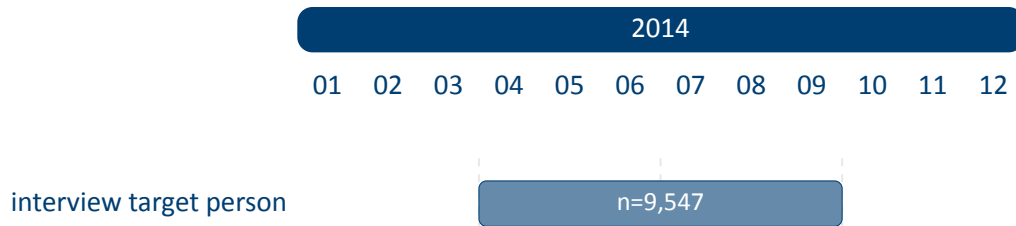
- Target persons

**Sample** Panel sample. Follow-up survey with interviewees willing to participate in the panel.

**Data collection** DZHW - German Centre for Higher Education Research and Science Studies, Hannover

**Mode of survey** Online survey (CAWI)

### 2.4.7 Wave 7: 2014 (CATI+competences)



**Figure 9:** Field times and realized case numbers in wave 7

- Target persons (Subsample A)

**Current wave** All students excluding the teaching-oversampling. (see section 2.2 for more information about this subpopulation).

**Sample** Panel sample. Follow-up survey with interviewees willing to participate in the panel.

**Data collection** DZHW - German Centre for Higher Education Research and Science Studies, Hannover

**Mode of survey** Computer-assisted telephone interview (CATI)

- Target persons (Subsample B)

**Current wave** Students who study an economic subject or have graduated from such studies. (identifiable via tx80921 in CohortProfile).

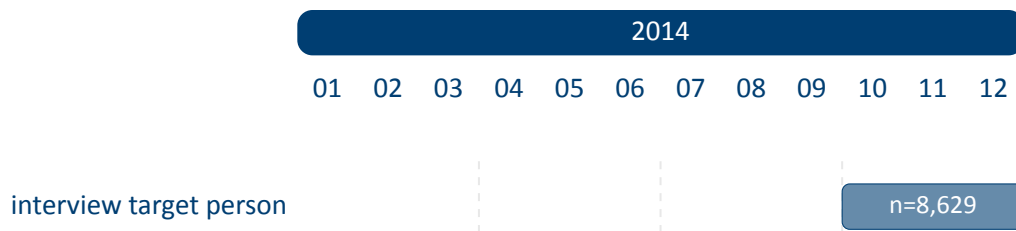
**Sample** Panel sample. Follow-up survey with interviewees willing to participate in the panel.

**Competence tests** Business Administration and Economics

**Data collection** DZHW - German Centre for Higher Education Research and Science Studies, Hannover

**Mode of survey** Paper-based competence testing within a personal-verbal interview (CAPI)

### 2.4.8 Wave 8: 2014 (CAWI)



**Figure 10:** Field times and realized case numbers in wave 8

- Target persons

**Sample** Panel sample. Follow-up survey with interviewees willing to participate in the panel.

**Data collection** DZHW - German Centre for Higher Education Research and Science Studies, Hannover

**Mode of survey** Online survey (CAWI)



### 2.4.9 Wave 9: 2015 (CATI)



**Figure 11:** Field times and realized case numbers in wave 9

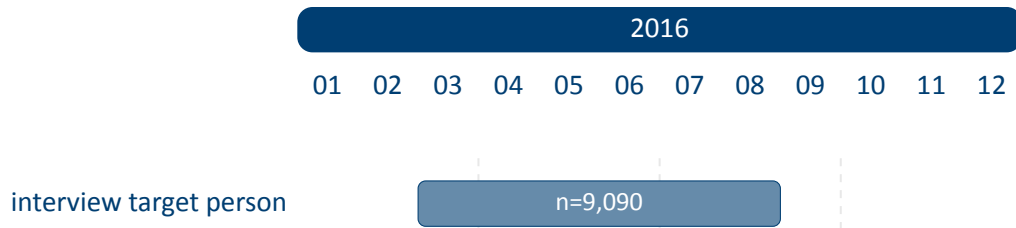
- Target persons

**Sample** Panel sample. Follow-up survey with interviewees willing to participate in the panel.

**Data collection** infas – Institute for Applied Social Sciences, Bonn

**Mode of survey** Computer-assisted telephone interview (CATI)

### 2.4.10 Wave 10: 2016 (CATI)



**Figure 12:** Field times and realized case numbers in wave 10

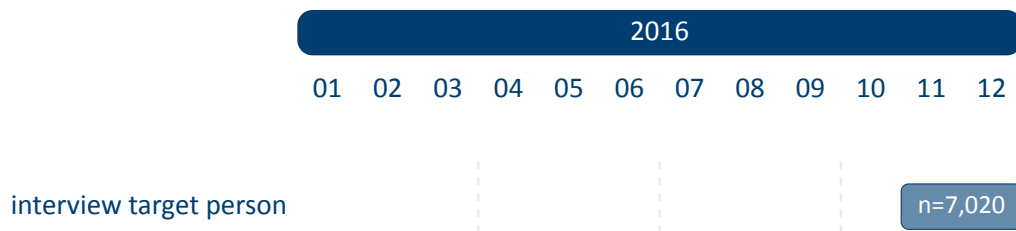
- Target persons

**Sample** Panel sample. Follow-up survey with interviewees willing to participate in the panel.

**Data collection** infas – Institute for Applied Social Sciences, Bonn

**Mode of survey** Computer-assisted telephone interview (CATI)

### 2.4.11 Wave 11: 2016 (CAWI)



**Figure 13:** Field times and realized case numbers in wave 11

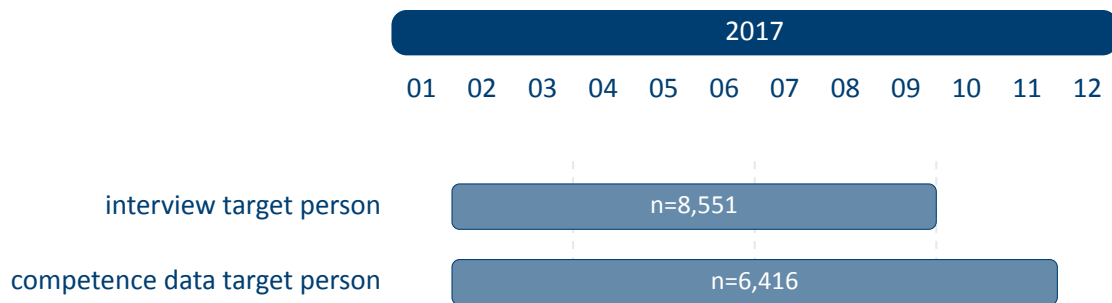
- Target persons

**Sample** Panel sample. Follow-up survey with interviewees willing to participate in the panel.

**Data collection** DZHW - German Centre for Higher Education Research and Science Studies, Hannover

**Mode of survey** Online survey (CAWI)

### 2.4.12 Wave 12: 2017 (CATI)



**Figure 14:** Field times and realized case numbers in wave 12

- Target persons

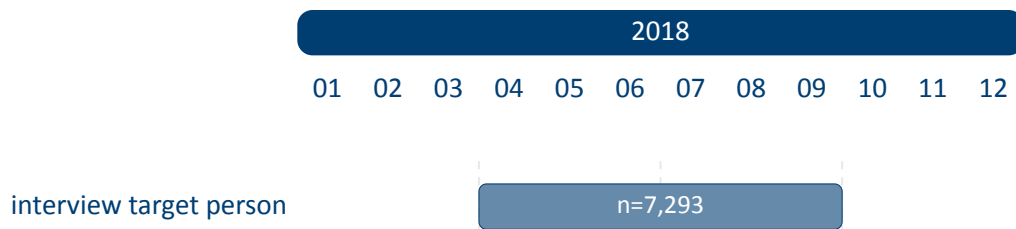
**Sample** Panel sample. Follow-up survey with interviewees willing to participate in the panel.

**Competence tests** Reading Competence, Mathematical Competence, English Reading Competence

**Data collection** infas - Institute for Applied Social Sciences, Bonn

**Mode of survey** Computer-assisted telephone interview (CATI) and group testing (computer-based testing with notebooks (CBA)) or individual testing (computer-based online testing, CBWA)

### 2.4.13 Wave 13: 2018 (CATI)



**Figure 15:** Field times and realized case numbers in wave 13

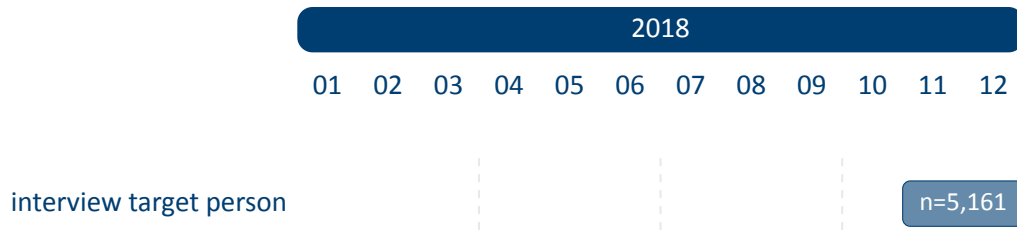
- Target persons

**Sample** Panel sample. Follow-up survey with interviewees willing to participate in the panel.

**Data collection** infas - Institute for Applied Social Sciences, Bonn

**Mode of survey** Computer-assisted telephone interview (CATI)

### 2.4.14 Wave 14: 2018 (CAWI)



**Figure 16:** Field times and realized case numbers in wave 14

- Target persons

**Sample** Panel sample. Follow-up survey with interviewees willing to participate in the panel.

**Data collection** DZHW - German Centre for Higher Education Research and Science Studies, Hannover

**Mode of survey** Online survey (CAWI)

## 3 General Conventions

The compilation of NEPS Scientific Use Files follows two general paradigms on how to edit the source data (i. e., the data that is delivered to the LfBi Research Data Center by the survey agencies). There may be exceptions to these principles that are explicitly noted in the respective documentation material.

The first and foremost paradigm in creating NEPS Scientific Use Files is the one of unaltered data. Wherever possible, the data editing procedures do neither change nor destruct the content of the original data. We consider this to be the basis for preserving the full research potential of the collected data. For this reason, no corrections are made during the entire data editing process to ensure the content validity of the source data. As a consequence, this means that the data in the Scientific Use File may contain implausible values, unless corresponding controls were already provided in the survey instrument. Only in rare cases, in which the responsible developers of a variable require the removal of clearly implausible information, these values are replaced by the special missing code *implausible value removed* (–52, see Table 6). The most prominent (and only systematic) exception to this general paradigm concerns the recoding of open responses that could originally have been recorded directly as closed responses (see section 3.4 for details). NEPS Scientific Use Files are equipped with a dataset `EditionBack-ups` that contains backup information for all content that has been modified by such recoding procedures (see section 4.2.4 for details).

The second paradigm is to integrate the data as much as possible without compromising the usability of the Scientific Use File. The underlying assumption is that for a vast majority of data users it is far more comfortable to reduce already integrated data for a specific analysis as opposed to correctly compile the relevant information from scattered source data themselves. In the end, each Scientific Use File contains only a few dozen integrated panel and spell datasets according to a general structure (see section 4.1.2 and section 4.1.3 for details), even if the compilation is based on several hundred separate source dataset files.

In addition to these two basic principles of data editing, there are several conventions for the data structure of all NEPS Scientific Use Files. The aim of this structuring is to ensure a maximum of consistency between the data of the different starting cohorts. In other words, a researcher who is familiar with the data logic of a particular NEPS cohort should be able to immediately recognize this structure when starting to work with data from another NEPS cohort. These conventions are explained in more detail in the following sections.

### 3.1 File names

The naming of the data files in NEPS Scientific Use Files follows a series of rules that are summarized in Table 3. The different elements are concatenated with an underscore ( `_` ) to generate the complete file name.

**Table 3:** Naming conventions for NEPS file names

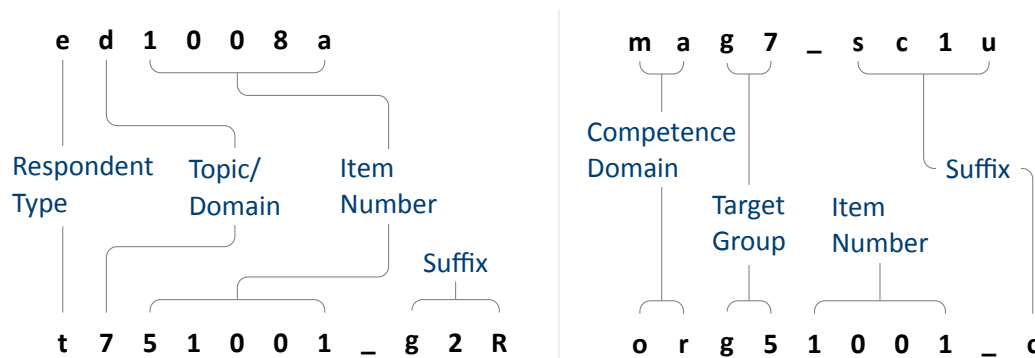
Element	Definition
SC[1–6]	<p><b>Indicator for the starting cohort</b></p> <p>1 = Newborns  2 = Kindergarten  3 = Fifth-grade students  4 = Ninth-grade students  5 = First-year university students  6 = Adults</p>
[filename]	<p><b>Meaning of the file name</b></p> <p><i>Prefix:</i> x = cross-sectional file; sp = spell file; p = panel file</p> <p><i>Keyword:</i> indicates the content of the corresponding file (e. g., data file xTarget contains cross-sectional data from the target questionnaire; spSchool contains spell data from the school history)</p> <p>File names of generated datasets do not have a prefix and always start with a capital letter (e. g., CohortProfile, Weights)</p>
[D,R,O]	<p><b>Indicator for the confidentiality level</b></p> <p>D = Download version  R = Remote access version  O = On-site access version</p>
[#]–[#]–[#](_beta)	<p><b>Indicator for the release version</b></p> <p><i>First digit:</i> the main release number is incremented with every further wave in the Scientific Use File; e. g., the first digit 5 implies that data of the first five survey waves are included in the release</p> <p><i>Second digit:</i> the major update number is incremented with every bigger change to the Scientific Use File; major updates affect the data structure, so updating the syntax files may be necessary</p> <p><i>Third digit:</i> the minor update number is incremented with every smaller change to the Scientific Use File; minor updates affect the content of cells, so updating the syntax files is not necessary</p> <p>_beta: this suffix indicates a preliminary Scientific Use File release which allows users to test the data before the main release; the beta release is no longer available after the main release</p>

For instance, the file SC5\_CohortProfile\_D\_14.0.0.dta refers to the *CohortProfile* data of *Starting Cohort 5* in its *Download* version of the Scientific Use File release 14.0.0.



### 3.2 Variables

The naming conventions for variables in NEPS Scientific Use Files aim to ensure maximum consistency both between the panel waves and between the starting cohorts. The names also refer to different characteristics and thus allow the data user an orientation regarding the contents of the variables. The principles of these naming conventions are exemplified in Figure 17. It has to be noted that a separate nomenclature is used for variables from competence measurements. Section 3.2.1 offers a detailed description of the general naming conventions for NEPS variables; the logic of naming competence variables is introduced in section 3.2.2.



**Figure 17:** General variable naming (left) and competence variable naming (right)

#### 3.2.1 Conventions for general variable naming

A variable name consists of up to four elements: the respondent type, the domain of information, an item number, and an optional suffix providing further information.

**Table 4:** Conventions for variable names

Digit	Description
1	<b>Respondent type</b>
	Indicator to which group of respondents the variable refers; note that variables related to the target person start with t even if the target person was not the actual informant (e. g., generated variables, list data from schools/kindergartens)
t	= Target person
p	= Parent of target person
e	= Educator/childminder
h	= Head/manager of institution (information about school/kindergarten)

(...)

**Table 4:** (continued)

Digit	Description
2	<p><b>Topic/domain</b></p> <p>Indicator to which theoretical dimension or educational stage the variable refers</p> <ul style="list-style-type: none"> <li>1 = Competence development</li> <li>2 = Learning environments</li> <li>3 = Educational decisions</li> <li>4 = Migration background</li> <li>5 = Returns to education</li> <li>6 = Interest, self-concept and motivation</li> <li>7 = Socio-demographic information</li> <li>a = Newborns and early childhood education</li> <li>b = From kindergarten to elementary school</li> <li>c = From elementary school to lower secondary school</li> <li>d = From lower to upper secondary school</li> <li>e = From upper secondary school to higher ed./occ. training/labor market</li> <li>f = From vocational training to the labor market</li> <li>g = From higher education to the labor market</li> <li>h = Adult education and lifelong learning</li> <li>s = Basic program</li> <li>x = Generated variables</li> </ul>
3–7	<p><b>Item number</b></p> <p>Indicator for the item number which typically consists of four numeric characters plus one alphanumeric character</p>
8–11	<p><b>Suffixes</b> (optional, see below)</p> <p>Indicator for several types of variables; separated from the previous characters by an underscore</p>

### Suffixes

- **Generated variables:** The \_g# suffix indicates a generated variable; the running number after \_g is in most cases a simple enumerator (e. g., \_g1). Since scale indices are generated by a set of other variables, they are also identified by a \_g# suffix. Note that scale indices are named after the first of the set of variables from which they were generated. In this case, numbering is only relevant if the first variable is identical for several scale indices. The number after \_g is in most cases a simple enumerator. However, there are two types of generated variables that assign specific meanings to digits, namely regional and occupational variables. The former are based on the Nomenclature of Territorial Units for Statistics (NUTS):

- g1: Indicator for East or West Germany
- g2: NUTS level 1 (federal state/Bundesland)
- g3: NUTS level 2 (government region/Regierungsbezirk)
- g4: NUTS level 3 (district/Kreis)

Generated variables for occupational classifications and prestige indices (see also section 3.4):

- g1: KldB 1988 (German Classification of Occupations 1988)
- g2: KldB 2010 (German Classification of Occupations 2010)
- g3: ISCO-88 (International Standard Classification of Occupations 1988)
- g4: ISCO-08 (International Standard Classification of Occupations 2008)
- g5: ISEI-88 (International Socio-Economic Index of Occupational Status 1988)
- g6: SIOPS-88 (Standard International Occupational Prestige Scale 1988)
- g7: MPS (Magnitude Prestige Scale)
- g8: EGP (Erikson, Goldthorpe, and Portocarero's class categories)
- g9: BLK (Blossfeld's Occupational Classification)
- g14: ISEI-08 (International Socio-Economic Index of Occupational Status 2008)
- g15: CAMSIS (Social Interaction and Stratification Scale)
- g16: SIOPS-08 (Standard International Occupational Prestige Scale 2008)
- *Versions of variables:* If question formulations, interviewer instructions, etc. change between panel waves to such an extent that sufficient meaning equivalence is no longer guaranteed, the answers to these questions are stored in different versions of a variable. The data for the latest and most current version of a question are provided under the variable name without any version suffix. Previous item versions are identified by `var_v1` for the data before the question was modified for the first time, `var_v2` for the data before the question was modified for a second time, and so on.
- *Harmonized variables:* The suffix `var_ha` indicates a harmonized variable in which common information from different versions of a variable is integrated. This is often done by aggregating detailed value characteristics into common superordinate categories. In other words, a harmonized variable reflects the lowest common denominator of information from a variable and its version(s).

- **Wide format variables:** The `_w#` suffix indicates variables that are stored in wide format. Note that this suffix does not necessarily imply a wave logic. The presence of a set of variables `var_w1`, `var_w2`, ..., `var_w10` may mean that there are up to 10 values for this variable per person or episode. This is the case, for example, if the corresponding item in the survey instrument was repeatedly measured in a loop. Another example concerns the date of the competence measurement within a survey wave if it took place on two different days.
- **Confidentiality level:** The `_D`, `_R`, or `_O` suffix indicates variables that have been modified during the anonymization process (see section 1.4). The suffix `_O` signals that data in this variable is only available via on-site access; `_R` refers to variables where access to detailed information is only possible via RemoteNEPS and on-site stay; and `_D` means that data in this variable has been extracted from the corresponding `_O` or `_R` variable to make at least some information available in the download version of the Scientific Use File. The confidentiality suffixes stand either alone (e. g., country of birth: `t405010_R`) or in combination with other suffixes (e. g., district of place of birth: `t700101_g3R`).

### Teaching specific variables

Certain parts of the survey in Starting Cohort 5 refer to teaching. The corresponding information in the datasets can be identified by variable names: Variables with the first three characters `tg6` or `tg8` indicate questions specifically addressed to (prospective) teachers.

### 3.2.2 Conventions for competence variable naming

The naming of variables from competence measurements and direct measures follows an alternative logic. In contrast to other data files, the competence datasets (`xTargetCompetencies` and `xDirectMeasures` in Starting Cohort 1) are structured in wide format; that is, all values for a single respondent are represented in one row of the data matrix. Thus, the integration of information from several competence domains collected across several survey waves requires specific conventions for variable naming. Competence variables are characterized by three name components and supplementing suffixes. The first component indicates the competence domain of the measurement (two characters, e. g., `vo` for vocabulary). The second part identifies the target group and the survey wave or class level in which the measurement was first used (two or three characters, e. g., `k1` for kindergarten children during the first wave). The target group identification does not necessarily indicate the cohort or testing wave of the measurement. Please refer to the explanations in the next section for the special features of repeatedly used test items. Some competence measurement are not designed for specific age groups, but are implemented unmodified in different cohorts and testing waves. In these cases the target group is defined as `ci` (cohort invariant). The third component denotes the item number. Table 5 contains a list of all possible specifications of the three parts of a competence variable name.

The additional suffixes inform about the mode of test execution if more than one survey modus has been applied for a measurement and about the sort of item score and overall competence score. There is a distinction between scored items named [varname]\_c and scored partial credit-items named [varname]s\_c. The latter is relevant if more than one correct solution is possible (e.g., value 0 = 0 out of two points, value 1 = 1 out of two points, value 2 = 2 out of two points), whereas the former is applied for dichotomous solutions (value 0 = not solved, value 1 = solved). In addition to the item scores, several aggregated scores are provided for competence measurements. They are indicated by \_sc[number] and a few special suffixes for Starting Cohort 1. A letter appended to the suffix indicates that more than one aggregated score for a competence measurement is available (e.g., \_sc3a, \_sc3b for different sum scores of any test). Detailed descriptions on how the aggregated competence scores were estimated can be found in the domain-specific documentation reports. The last part of Table 5 shows all possible suffixes and their meanings.

**Table 5:** Conventions for competence variable names

### Part I: Competence Domain (2 chars)

ba	Business administration and economics
bd	Backwards digit span: Phonological working memory
ca	Categorization: SON-R subtest
cd	Cognitive development: Sensorimotor development
de	Delayed gratification: Executive control
dg	Domain-general cognitive functions (DGCF): Cognitive basic skills
ds	Digit span: Phonological working memory
ec	Flanker task: Executive control
ef	English foreign language: English reading competence
fa	FAIR: Concentration abilities
gr	Grammar: Listening comprehension at sentence level
hd	Habituation-dishabituation paradigm
ic	Information and communication technology literacy (ICT)
ih	Interaction at home: Parent-child interaction
ip	Identification of phonemes: Phonological awareness
li	Listening: Listening comprehension at text/ciscourse level
lk	Early knowledge of letters
ma	Mathematical competence
md	Declarative metacognition
mp	Procedural metacognition
nr/nt	Native language Russian/Turkish: Listening comprehension
on	Blending of onset and rimes: Phonological awareness
or	Orthography
re	Reading competence
ri	Rimes: Phonological awareness

(...)

**Table 5:** (continued)

rs	Reading speed
rx	Early reading competence
sc	Scientific competence
st	Scientific thinking: Science propaedeutics
vo	Vocabulary: Listening comprehension at word level

**Part II: Target Group (1 char), followed by wave or grade (1-2 digits)**

n#	Newborns in wave #
k#	Kindergarten children in wave #
g#	Students at school in grade #
s#	University students in wave #
a#	Adults in wave #
ci	Cohort invariant (for instruments administered unchanged in all cohorts)

**Part III: Item number (3-4 chars)**

For some competence domains, these item numbers follow a certain scheme, but for most competence domains they only indicate the different items

**Part IV: Suffixes (starting with an underscore)**

_pb	Paper-based test modus (proctored)
_cb	Computer-based test modus (proctored)
_wb	Web/Internet-based test modus (unproctored)
_c	Scored item variable (s_c for partial credit-items)
_sc1	Weighted likelihood estimate (WLE) <sup>12</sup>
_sc2	Standard error for the WLE <sup>2</sup>
_sc3	Sum score
_sc4	Mean score
_sc5	Difference score (for procedural metacognition)
_sc6	Proportion correct score (for procedural metacognition)
_p	Maximum value for an item (only in Starting Cohort 1)
_b	Minimum value for an item (only in Starting Cohort 1)
_m	Mean value for an item (only in Starting Cohort 1)
_s	Sum value for an item (only in Starting Cohort 1)
_n	Number value for an item (only in Starting Cohort 1)

**1** WLEs and their standard errors are estimated in tests that are scaled based on models of Item Response Theory (cf. Pohl and Carstensen, 2012).

**2** WLEs and their standard errors are corrected for test position; uncorrected WLEs and standard errors are indicated by an additional u in the suffix (\_sc1u, \_sc2u).

### Identification of repeated test items

In some competence measurements identical items are implemented in different testing waves (e. g., mathematics). Identifying repeatedly measured test items in NEPS data can be easily done by looking for competence variables with an identical word stem. If the same test item is surveyed in different survey waves or starting cohorts, the variable name is equipped with an additional suffix. It is important to know that the two or three characters for the target group (second part of the variable name) always indicate the wave or cohort in which the item was initially used. The word stem is then fixed and does not change when the item is used again in later waves or other cohorts. If the variable name does not contain a suffix for repeated use, then the second part of the word stem refers to the target group of the realized measurement. However, if the variable name includes a suffix for repeated use, then the values of the variable do not refer to the target group according to the word stem, but to the target group according to the suffix. The suffix that points to the repeated use consists of two parts: The first element indicates the starting cohort of current item administration and the second element indicates the cohort or testing wave of current item administration.

The following example illustrates this logic: The competence variable `vok10067_sc2g1_c` is a vocabulary item (vo) that was initially measured during the first kindergarten survey wave (k1). However, the values in this variable reflect the scored measurements of this item's repeated use among the target persons of Starting Cohort 2 in the course of the survey wave in grade 1 (`_sc2g1`), and thus two years after the first measurement.

### 3.2.3 Labels

As a rule, the seven-digit variable names are not sufficient to uniquely identify the respective contents of the variables and to differentiate sufficiently between items. All variables therefore have *variable labels* for more detailed description. In addition, most variables contain *value labels* for the respective value characteristics. All information is available in German and English and is typically displayed directly in the editor of the statistics program, e.g. for frequency calculation or when searching the data (applies to SPSS and Stata, see also section 1.3). For users of R, see section A.1 for hints on this.

In addition to the variable and value labels, the datasets also contain extended characteristics for variables. These include the question text from the survey instrument, any associated interviewer instructions and filter conditions, as well as other meta information. All extended features can be accessed directly within data files. Stata users apply the `infoquery` command for this, which is part of the *NEPS tools* package (see section 1.8). SPSS users will find the additional meta information in the "Variable View" at the end of each variable line.

As explained in more detail in section 4, NEPS data from different waves are integrated as much as possible. For panel data, this primarily means that many variables contain information from multiple waves. In most cases of such a data integration, the meta information between the

waves does not change. However, if there are changes to the meta information of a repeatedly measured item, and if these changes are not significant enough to store the information in separate variables, the assignment of meta information follows a general rule: **All meta information available in a dataset always corresponds to the most recently instrument in which the respective item was used.**

A concrete example is the adaptation of interviewer instructions or question texts from the informal salutation (“Du”) to the formal salutation (“Sie”). Since these changes are not expected to have any effect on how a question is answered, the corresponding values across multiple waves get integrated into one variable. If you request the meta information of such a variable in the dataset, the wording of the latest item formulation will be displayed (in the given example with the formal salutation “Sie”). In case of uncertainties regarding the continuity of meta information of a variable across different waves, we recommend to consult the respective *survey instruments* for the individual waves.

### 3.3 Missing values

The NEPS data contain various missing codes to differentiate between various types of missing values. All missing codes have negative values or are defined as system missing. Depending on the statistics program used, you must ensure that these codes are processed correctly. In the offered SPSS datasets, the missing codes are already defined as missing values. When using Stata, the missing codes must first be excluded from the analyses by the user as missing values. For this purpose the command `nepsmiss` is available in the *NEPS tools* package (see section 1.8). The general recommendation is to always carefully check the frequency distributions of the relevant variables before running an analysis.

We distinguish between three types of missing codes, which are summarized in Table 6 and described in more detail below.

**Item nonresponse:** The first type of missing codes occurs when a person has not (validly) replied to a question.

- The most common cases of item nonresponse are *refused* (–97) answers and *don’t know* (–98) answers.
- Missing values specified by the survey agency due to an incorrect use of the instrument are coded as *implausible value* (–95).
- Within the competence data, there is a special missing code indicating that a question or test item was *not reached* (–94) due to time constraints or other test setting restrictions. It usually signals that the respondent had to quit the test somewhere before this point.
- Other missing codes refer to various categories of *item-specific nonresponse* (–20, ..., –29) such as –20 for “*stateless*” in the citizenship variable `p407050_D`.



**Table 6:** Overview of missing codes

Code	Meaning	Note
<b>Item nonresponse</b>		
–94	not reached	only relevant for instruments with time restrictions (e. g., competency test measures)
–95	implausible value	assigned by the survey agency (e. g., multiple answers to a one-answer question in PAPI mode)
–97	refused	as default answer option to the question
–98	don't know	as default answer option to the question
–20,...,–29	various	item-specific missing with informative value label (e. g., “no grade received” for question about school grades)
<b>Not applicable</b>		
–54	missing by design	question not included in (sub)sample-specific instrument (e. g., not asked in all waves)
–90	unspecific missing	in PAPI mode (e. g., question not answered, empty field)
–93	does not apply	as default answer option to the question
–99	filtered	filtered out question, in other than CATI/CAPI mode
.	system	filtered out question, in CATI/CAPI mode
<b>Edition missings (recoded into missing)</b>		
–52	implausible value removed	only at the request of the responsible item developers
–53	anonymized	sensitive information removed (e. g., country of birth of parents in the download version)
–55	not determinable	not sufficient information to generate the variable value (e. g., net household income t510010_g1)
–56	not participated	in case of unit nonresponse, only used in certain datasets

**Not applicable:** The second type of missing codes occurs when an item does not apply to a respondent.

- The code *missing by design* (–54) is assigned when respondents in a (sub)sample have not been asked the respective questions. This is usually the case if the administered survey instrument contains (sub)sample-specific questionnaire modules. The code is also used for the more general case where values of a variable are not available due to the design of the survey (e. g., measurement rotation with either easier or heavier test tasks).
- If the respondent him-/herself or the interviewer indicates that a particular question is not applicable to the person, the missing value is coded as *does not apply* (–93). If, on the other hand, filtering takes places automatically via the survey instrument, the coding of the filtered out questions depends on the survey mode: in CATI and CAPI interviews, a system missing value (.) is assigned for this; in all other modes the respective code is *filtered* (–99).

- Missing values that cannot be assigned to any of the above categories are coded as *unspecific missing* (–90). This missing code usually occurs in PAPI questionnaires when a respondent has not answered a question for unknown reasons.

**Edition missings:** The third type of missing codes is defined in the process of data preparation for the Scientific Use File.

- If in the data edition process certain values which are not considered to be meaningful are requested to be removed, the missing code *implausible value removed* (–52) is assigned in their place. As a rule, however, all values from the field instruments are included in the Scientific Use File without further plausibility checks (see section 3). Only in exceptional cases, when the responsible item developers explicitly recommend a removal of implausible answers, this missing coding is done.
- Sensitive information that is only available via remote and/or on-site access is encoded in the more anonymized data access option as *anonymized* (–53).
- In general, coding schemes are used to generate variables (e. g., occupational coding; see section 3.4). However, if the information from the original data is not sufficient to generate a suitable value, the missing code *not determinable* (–55) is used instead.
- If a person was not present during the interview or did not complete a questionnaire at all, even though it was administered to the person, the concerning variables receive the code *not participated* (–56). This missing code is special in that target persons without survey data for a certain wave (e. g., due to illness) are usually not included in the corresponding datasets. This missing code is only used in the special cases of datasets that integrate several waves in wide format (e. g., xTargetCompetencies) or that also contain observations for non-participating persons in a wave (e. g., CohortProfile).

### 3.4 Generated variables

#### Coding and recoding of open responses

At various points in the NEPS survey instruments there are so-called open questions where respondents can or should enter their answers as text. A typical example is information about occupation.

The open text format allows respondents to specify anything they want. A practical way to deal with the resulting string information is to code and recode the information for further processing and later analyses. In general, coding describes the process of assigning one or more codes from selected category schemes to the string information, e. g. the classification of occupational data according to DKZ (database of documentation codes, *Datenbank der Dokumentationskennziffern*) or WZ (classification of economy branches, *Klassifikation der Wirtschaftszweige*).

The term “recoding” is used here to describe the process of assigning a code from an already presented closed answer scheme. This usually applies to semi-open question formats where

respondents enter a text under the category “other”, but which can be assigned ad hoc to one of the given closed answer categories. The recoding therefore does not define any new codes; the presented answer scheme of the respective question is not extended.

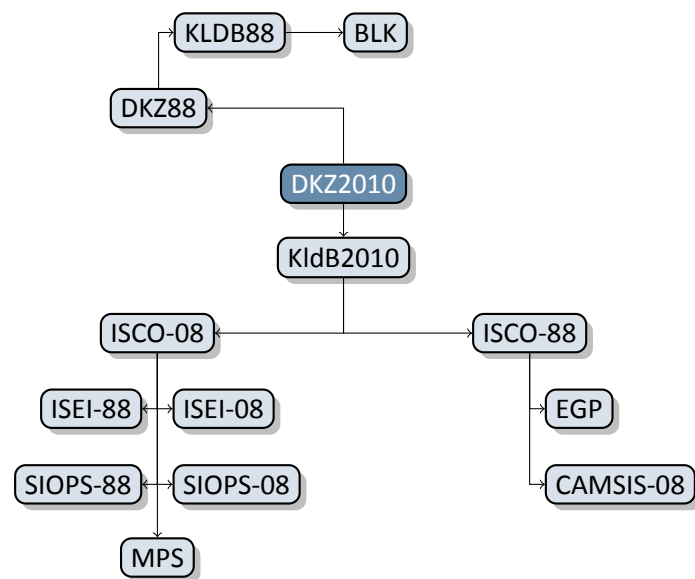
The most common and comprehensive coding scenarios in the fields of occupation, education, branches, courses, and regional information are processed by the Research Data Center (FDZ-LIfBi) itself. Other coding tasks are distributed among the responsible departments at the LIfBi in Bamberg and the partners in the NEPS consortium.

### Derived scales and classifications

The (re-)coding of open answers or string entries into primary classifications (such as DKZ2010 or WZ08) is a first and essential step towards making this information available within the NEPS Scientific Use Files in a user-friendly and analyzable way. The standardized derivation of further classifications or scales, especially in the area of educational qualifications and occupational titles, is a second and no less important step. At least three types and objectives of derivations can be distinguished:

- Derivations from primary classifications (and originated from string entries/open answers) into other classifications that function as a standard schema in other studies or international comparisons, e. g. ISCO instead of KldB in the field of occupations
- Derivations from primarily closed response schemes into general classifications and schemes using auxiliary information, e. g. ISCED or CASMIN from school certificate and training data plus additional information on the type of school/training
- Combination of the two types, e. g. EGP class scheme via derived ISCO classification plus information on self-employment and supervisory status

Figure 18 shows the derivation paths for several occupational scales and schemes provided in the NEPS. A detailed description of the standard derivations for educational attainment (ISCED, CASMIN and Years of Education) can be found in the corresponding documentation report by Zielonka and Pelz, 2015.



**Figure 18:** Derivation paths for several occupational scales and schemes provided in the NEPS

## 4 Data Structure

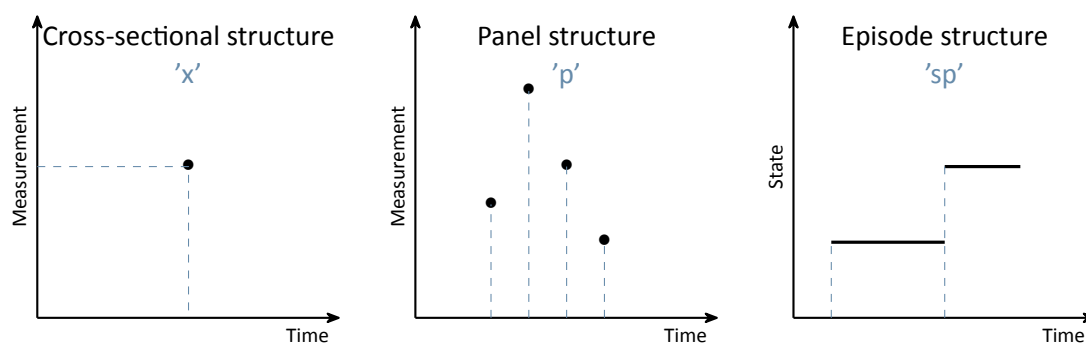
### 4.1 Overview

The broad objectives and the large size of the longitudinal NEPS surveys inevitably lead to a complex database. The crucial task is to organize this data in a well-structured, reproducible and user-friendly way, while at the same time preserving a maximum level of detail in the data. To meet this challenge, a number of additionally generated variables and datasets is included in the Scientific Use File to facilitate the preparation and analysis of the data.

In principle, all information collected in the course of a panel wave is appended to the information from previous waves in the corresponding data file, together with the required identifiers. Data files containing longitudinal information from several waves are denoted with a *p* in the file name. For example, the *pTarget* file(s) contain(s) information from the target persons' interviews with one row in the dataset representing the information of one target in one wave.

This convention does not apply to all longitudinal data. For example, there are competence measurements that were repeatedly carried out with the same target persons. However, since the instruments, i.e. the content of competence tests, vary over time, the corresponding information is structured in wide format (for more details, see section 3.2.2 or section 4.2.34). Such cross-sectionally structured data files with one line representing information of a respondent from all waves are marked with a *x*.

Another type of data structuring refers to episode data. For the information collected prospectively and retrospectively using iterative question sets, the Scientific Use File provides life area-specific spell datasets. These datasets are marked by a preceding *sp*. An example is the file *spEmp*, which informs about current and former episodes of employment.



**Figure 19:** Different types of data structures

In addition to interview and test data provided by the respondents as well as episode data, there are also so-called paradata or derived information. These data files can be identified by

the leading capital letter in the name (e.g. `Weights` or `CohortProfile`). In most cases, these datasets correspond to the panel structure.

### 4.1.1 Identifiers

The multi-level and multi-informant design of the NEPS and the distribution of survey information across different datasets requires the use of multiple identifiers. The following identifier variables are relevant in this Starting Cohort for linking data:

**ID\_t** identifies a target person. The variable `ID_t` is unique across waves and samples (and also starting cohorts).

**wave** indicates the sample wave in which the data was collected.

**ID\_i** identifies the respective educational institutions such as kindergartens or day care centers, schools, universities, etc. The variable `ID_i` is unique across waves and starting cohorts.

In addition, there are other identifier variables to indicate a target person's membership in a particular test group (`ID_tg` in `CohortProfile`, not applicable to all starting cohorts) and to indicate the interviewer who conducted the respective interview (`ID_int` in `Methods` datasets). However, these identifiers are not relevant for the merging of information from different datasets and are negligible for most empirical applications.

### 4.1.2 Panel data

As mentioned above, all information from subsequent survey waves are appended to the already existing data files (as far as possible). This method of data processing generates *integrated panel data* files in a long format as opposed to providing one separate file per wave (where each file contains only the information from a single wave). When working with the integrated panel data in the NEPS Scientific Use Files, the following points should be considered:

- A row in the dataset contains the information of one respondent from one survey wave.
- This means that more than one identifier variable is needed to identify a single row for uniquely selecting and merging information from different datasets. These are usually `ID_t` and `wave`.
- It also means that although not all variables were administered in each survey wave, the integrated structure of the dataset contains cells for all variables of all waves. If no data is available, e.g. because a variable was not queried in a particular wave, the corresponding cells are filled with a missing code (see section 3.3).
- Once information about a variable has been surveyed from one individual across multiple waves, the corresponding data is distributed across multiple rows in the dataset.

This long format is usually the preferred data structure for the analysis of panel items with information from several waves. However, cross-sectional information is often also required, e. g. because it depicts time-invariant characteristics or was collected only once for other reasons. In most analysis scenarios, the combined set of relevant variables is not measured in a single wave. Therefore, the corresponding data cannot be analyzed together straightaway because it is stored in *different rows* of the dataset. Cross-tabulating these variables in their current state results in an L-shaped table in which all observations of one variable fall into the missing category of the other variable and vice versa. The best way to deal with this issue depends very much on the intended analysis and the methods used. Two typical procedures are:

- First, the integrated panel data file is split into wave-specific subfiles so that each dataset contains only information from one wave. The relevant information from these subfiles is then merged together by using only the respondent's identifier (ID\_\*) as key variable. The wave variable is not needed here and remains neglected. Before this step, variables may need to be renamed to make them wave-specific identifiable. The result is a dataset with a cross-sectional structure in which the information of a respondent is summarized in one single row (wide format). Stata's *reshape* command (and similar tools in other software packages) basically follow this strategy.
- Alternatively, the panel structure is retained and the values from observed cells are copied into the unobserved cells. If, for example, the place of birth was only surveyed in the first wave, the corresponding value can be transferred to the respective cells of the other waves of the respondent. This method is particularly useful for time-invariant variables (e. g. country of birth, language of origin), which are usually collected only once in a panel study.

### 4.1.3 Episode or spell data

Handling cross-sectional data is usually not a problem. Most data users also know how to work with and analyze panel data. Episode or spell data, on the other hand, present a particular challenge for understanding data processing. The following explanations should help to deal with this data format in a meaningful and appropriate way.

In the episode (or spell) data, there is one row for each episode that was captured. Note that the number of episodes per se is independent of the survey wave. This means that several episodes (=several rows) can be recorded in a single wave. Usually, a start date and an end date describe the duration of an episode. The remaining variables in such spell datasets contain additional information about that episode. These characteristics are chronologically linked to the episode. In other words: Especially for time-variant variables (e. g. ISEI, CASMIN) it is important to know that the respective values indicate the status of the respondent **at the time of the episode** and not necessarily the current status.

To give an example: In the spell dataset spEmp there is a period of time for a certain respondent in which he or she worked without interruption in a particular job. If this person changes to a new job, this marks a new episode which is stored in a new data row. Further changes in

this context may also lead to new episodes, e. g. a change of employer or the conclusion of a new employment contract (but not if the salary, working hours or other characteristics of the respective job change). Episodes can therefore be understood as the smallest possible units of one's life history, in this case the employment biography. As soon as there are several relevant changes in such a biography between two consecutive interview dates, this is reflected in several data rows per survey wave.

In addition to these (time) episode data, which we call *duration spells*, there are two other types of episode data: Occurring events or the transition from one state to another (e. g., change of marital status, change of educational level) are recorded in so-called *event spells*; the existence of children, partners, etc. is recorded in so-called *entity spells* with one row per entity. Regardless of the type of episode, two variables are usually necessary to identify a single row in the data file, namely the respondents' identifier `ID_t` and an episode, event or entity numerator such as `spell` that identifies a duration spell. More detailed information on the required identifier variables can be found on the respective data file pages in section 4.2.

There is one important circumstance to consider when working with NEPS spell data. This concerns *subspells*. Biographical episode data are collected retrospectively. During an interview the respondents are asked about all episodes that have occurred since the last interview (in the first interview it is since birth or a certain age). If an episode is finished at the time of the interview, the respondent reports a corresponding end date and the spell is completed. Difficulties arise when the episode is not yet finished at the time of the interview, i.e. it is ongoing. Such an episode appears as right-censored in the dataset. In the next interview, this episode is then queried using preloads in the course of "dependent interviewing" in such a way that the respondent can report whether it has been finished in the meantime or whether it continues. Technically this leads to several rows in the data structure, which can be distinguished by the variable `subspell`:

- original (right-censored) episode reported in initial wave (`subspell=1`)
- continued episode reported in next wave(s) (`subspell=2`, `subspell=3`, etc.)

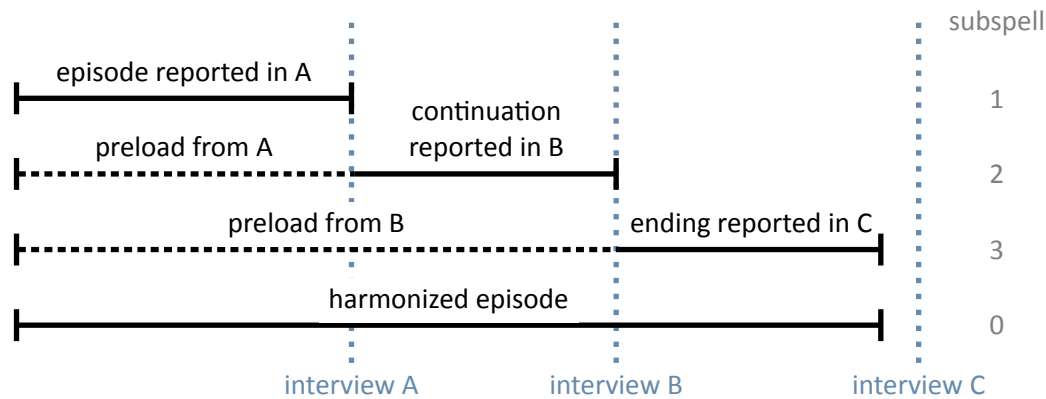
Normally, attention is paid to the last `subspell`, as it contains the most up-to-date information about an episode. However, the most recently captured information for an episode may contain missing values, or the value from the last-mentioned sub-episode may have been transferred to it. To facilitate the handling of spell datasets, *harmonized rows* with the (presumed) most relevant information for all episodes are generated from the summary of the corresponding `subspells`. This information often corresponds to the last (non-missing) reported entry, but sometimes also to the first (non-missing) reported entry of an episode. The harmonized rows are selectable by the filter condition `subspell=0`. The same applies to all episodes reported as completed without any `subspells`.<sup>1</sup> If there is no particular interest in `subspell` information, it is recommended to use only the harmonized data rows for analysis:

```
keep if subspell==0
```

<sup>1</sup> The variable `spgen` indicates whether an episode was originally reported as finished (`spgen=0`) or whether it is a harmonized (generated) episode (`spgen=1`).



Note that the information cumulated in the harmonized spells (last valid available) may originate from different waves. Please be also aware that the selection of harmonized spells should **not** be used when working with information stored in wide format (e. g., interruption episodes of vocational training spells in `spVocTrain`).



**Figure 20:** Logic of subspells

### 4.1.4 Revoked episodes

In order to reduce seem bias, spell data are preloaded by prior wave information. This information from previous waves can be revoked by the respondent during the current interview. Spell datasets therefore also contain information about revocations (variables `disagint`, `disagwave`). The reasons for a revocation or contradiction are manifold; they depend mainly on the information that is presented to the respondent to remember the episode (see the questionnaires for the exact wording of the episode data collection).

If an episode is later revoked by the respondent, this episode is marked accordingly in the dataset. The respective information is collected again in the current interview and saved as a new episode in the actual data collection wave. The updated spell is not flagged as a corrected spell. The identification of related spells (=previously given information plus their correction in the following wave) is up to the data user. Please note: Since it is technically impossible to specify a start date for an episode prior to the last interview date, virtually all corrected spell episodes are left-censored. The only exception are episodes that started on the interview date of the last wave.

In addition to the possibility of revoking an episode in the course of the subsequent survey wave, there is also the possibility of revoking an episode during the interview. For this purpose, a *check module* is used after the biographical information has been recorded. It ensures that the life course is captured as completely as possible. The biographical episodes asked in the thematically structured questionnaire modules are already examined in the interview for their chronological plausibility. To verify the temporal consistency of the events across the questionnaire modules, a complete overview of all types of events is created. For this purpose, all

recorded biographical episodes are displayed in tabular form in the check module. If gaps or overlaps are indicated, the respondent will be asked again. He or she can then make corrections, add new episodes, or revoke already recorded episodes. The identification of episodes revoked in the check module is possible in the spell datasets by the variable “Biography: Type of event (edited)” (spms=20); the addition of new episodes in the check module is indicated in the variable “Episode mode” (ts23550=4 in spEmp). A detailed description of the functionality of the check module for reported life courses (in German language) can be found on the website in the section “Data Documentation”:

→ [www.neps-data.de](http://www.neps-data.de) > Data Center > Data and Documentation  
    > Starting Cohort First-Year Students > Documentation

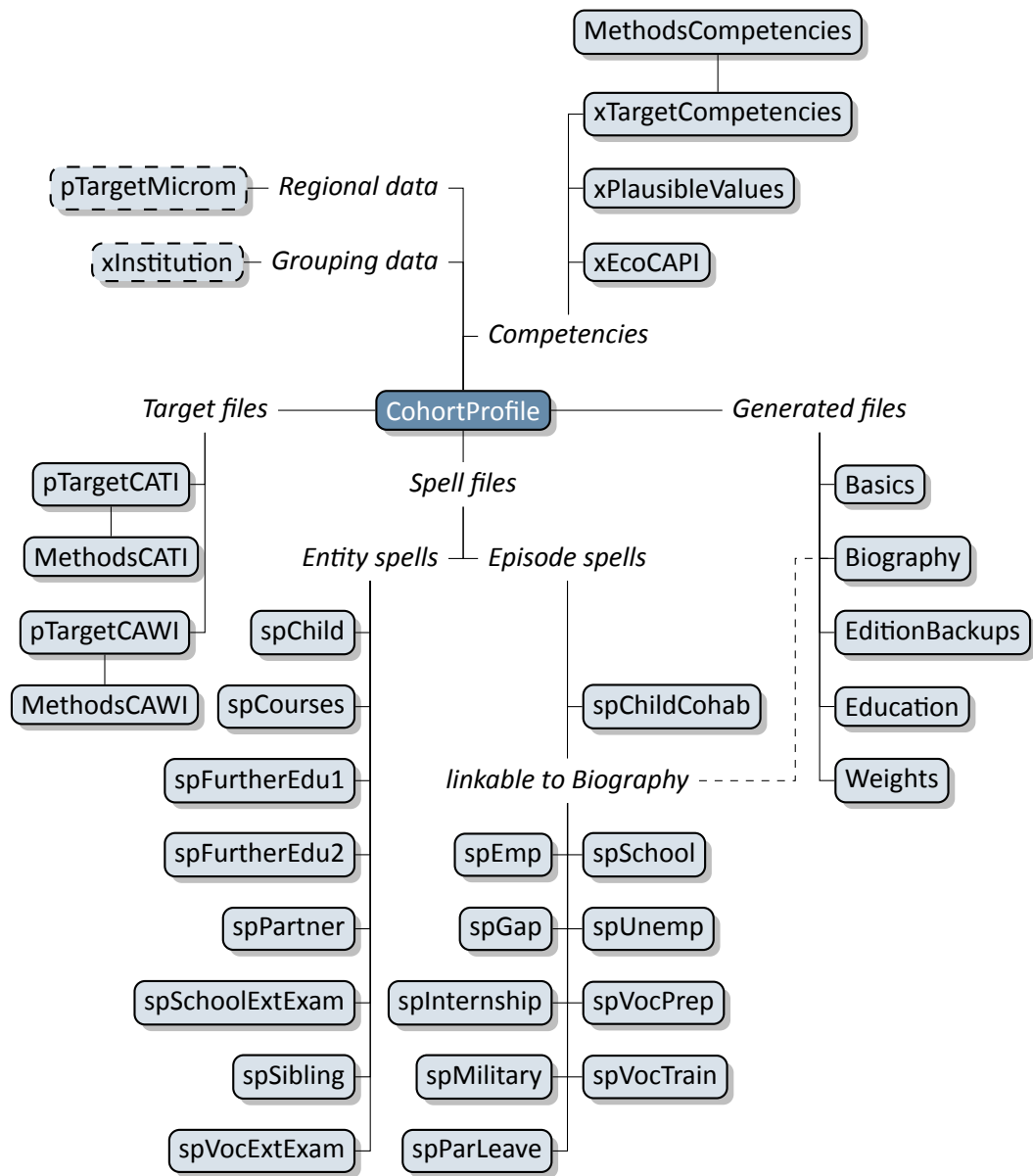
### 4.2 Data files

In the following section, every data file of this Starting Cohort is explained in a subsection, including a data snapshot and an example of data usage (in Stata). The examples are written so that everyone knowing Stata should easily understand it. You also do not need additional ado files installed, although you are highly advised to use the `nepstools` (see section 1.6).

To ease your understanding of the relationship of those files, Figure 21 provides an overview. The edges in this graph symbolize how a data file may be linked to other files. This is not meant to document every possible data link you could do but rather tries to give you an idea which data files relate most. By clicking on a node, you get directed to this data file’s explanatory page.

You need to set the following globals for the Stata examples to work. Just adapt and copy the lines below to the top of the syntax files or execute them in your Stata command line before running the syntax:

```
** Starting Cohort
global cohort SC5
** version of this Scientific Use File
global version 14-0-0
** path where the data can be found on your local machine
global datapath Z:/Data/${cohort}/${version}
```



**Figure 21:** Graphical overview of all data files. Each node represents one data file. Relations are indicated by connection lines. Files with a dashed border are not available in the Download version of the Scientific Use File. Click on a data file to get more information.

## 4.2.1 Basics

[« go back to overview](#)

## Description

Simplified information about respondents in a plain format

## File structure

wide format: 1 row = 1 respondent

## ID variables needed to identify a single row

ID\_t

## Other ID variables useful for linkage

none

## Number of variables / number of rows in file

82 / 17,910

## Contains data from waves



## Exemplary variables

ID_t	ID target
tx29000	Age at interview month (years)
t70000m	Date of birth: month
t70000y	Date of birth: year
t700001	Gender
tx29003	Mother tongue: German
tx29004	Nationality: German
tx29005	Born in Germany
t741001	Size of Household (persons)
tx29060	currently employed
tx29904	Main spells of type 'Emp' (number)
tx29007	Age at migration to Germany

## Exemplary data snapshot

ID_t	tx29000	t700001	tx29005	t741001	tx29060	tx29904
7005513	24.67	[w] female	yes	1	yes	4
7002273	22.75	[w] female	yes	1	yes	2
7004930	30.08	[m] male	yes	1	yes	3
7010083	28.67	[w] female	yes	1	yes	9
7011856	27.92	[w] female	yes	2	yes	2

This file contains the latest reported basic information on each respondent, e. g., sociodemographic variables like age in month (tx29000), born in Germany (tx29005), gender (t700001), currently employed (tx29060), but also household characteristics, etc. It also contains meta information about some episodes like the number of main employment spells (tx29904). This data is generated from the pTarget files and a number of spell files. The Basics file is updated prospectively. That is, the file is cross-sectional (i. e., one row per person) and always includes updated information from the latest panel wave a respondent has participated. This simplified data structure can help to gain a first insight in the data. However, it should be handled with care, as it may not feature the *best* information about the respondent. **Please use this file only to get a first overview of the data. Use the original panel or episode files for analyses!**

**Example 1 (Stata):** Working with Basics ([find R example here](#))

```
** open the data file
use ${datapath}/SC5_CohortProfile_D_${version}.dta, clear

** merge the data from Basics, enhancing every entry in CohortProfile
** (i.e. every wave, this is why m:1 merge is needed)
** with information from Basics
merge m:1 ID_t using ${datapath}/SC5_Basics_D_${version}.dta

** change language to english (defaults to german)
label language en

** tabulate gender by wave
tab wave t700001

** please note that now, you have the most recent information known about respondents
** in every wave. This does not have to be equal to the information actually surveyed
** in that wave!
** Proceed at your own risk!
```

## 4.2.2 Biography

[« go back to overview](#)

## Description

Integrated and edited life course data

## File structure

spell format: 1 row = 1 episode of 1 respondent

ID variables needed to identify a single row

ID\_t splink

Other ID variables useful for linkage

wave sptype

Number of variables / number of rows in file

10 / 221,618

Contains data from waves



## Exemplary variables

ID_t	ID target
splink	Link for spell merging
wave	Wave
sptype	Spell type
startm	Episode start (month)
starty	Episode start (year)
endm	Episode end (month)
endy	Episode end (year)
spms	Type of event
splast	Episode is ongoing

## Exemplary data snapshot

ID_t	splink	wave	sptype	starty	endy
7004953	240002	7	24	2012	2013
7007211	220002	1	22	2001	2010
7010899	220001	1	22	1997	2000
7014876	270002	7	27	2014	2014
7016550	300004	12	30	2016	2016

The file Biography serves to facilitate the analysis of complex life course data collected both retrospectively and prospectively. The dataset puts together harmonized episodes with educational and employment relevance from the following duration spell files: spSchool, spVocPrep, spVocTrain, spMilitary, spEmp, spUnemp, spInternship, spParLeave, and spGap. The variable sptype is provided to identify the source of each episode.

In contrast to the “raw” biographical data from each of the module-specific spell modules, the Biography file provides more consistent life course data that has been additionally checked and edited. In particular, inconsistencies in the individual life course data were identified and corrected during the interview with the help of a “check module”. Corrected times are stored in the duration spell files as \_g1 variables. For example, the variable ts2311y\_g1 in spEmp contains the starting date of an employment episode which was corrected within the check module. Such corrected times form the basis for further adjustments that are implemented in

the data editing process for Biography. Essentially, the following measures are taken to ensure the integrity of the life course data in this file:

- All subspells have been removed, i. e., Biography contains only completed, harmonized, or right-censored episodes (`subspell=0`).
- Episodes revoked by respondents during the interview or in the subsequent survey wave (see section 4.1.4) are deleted, unless the episode was re-recorded in the current wave. Revoked episodes are included in the original spell files and can be identified there with the corresponding marker variables (`spms` or `disagint`).
- Starting and end dates of episodes are smoothed and corrected, i. e., overlaps of one month and more between adjacent episodes have been resolved.
- Gaps between adjacent episodes that do not exceed two months are closed; gaps of more than two months are defined as specific gap episodes (`edition gaps`) within the Biography file.

Due to the additional editing steps and the compilation of spells from different biographical modules, it is recommended to use the Biography dataset as a starting point for life course analyses.

### Example 2 (Stata): Working with Biography (find R example [here](#))

```
** open the data file
use ${datapath}/${cohort}_Biography_D_${version}.dta, clear

** change language to english (defaults to german)
label language en

** check out which spell modules you can merge to this file
tab sptype

** check that you will need splink to merge information
** from other modules to this file
** (command gives no result, which means approval)
isid ID_t splink
```

## 4.2.3 CohortProfile

[« go back to overview](#)

## Description

Paradata on the cohort's panel sample

## File structure

long format: 1 row = 1 respondent in 1 wave

ID variables needed to identify a single row

ID\_t wave

Other ID variables useful for linkage

ID\_i ID\_tg

Number of variables / number of rows in file

18 / 250,740

Contains data from waves



## Exemplary variables

ID_t	ID target
wave	Wave
cohort	NEPS Starting Cohort
tx80220	Participation/drop-out status
tx80521	Data available: interview target person
tx80522	Data available: competence data target person
inty	Interview date (year)
intm	Interview date (month)
tx80524	Data available: Institution
testm	Test: Survey day (month)
testy	Test: Survey day (year)
tx80107	Sample: First participation in wave

## Exemplary data snapshot

ID_t	wave	tx80220	tx80521	tx80522	inty	tx80524
7011366	1	Participation	yes	yes	2011	yes
7011366	2	Temporary drop-out	no	Missing by design	-56	Not determinable
7011366	3	Temporary drop-out	no	Missing by design	-56	Not determinable
7011379	1	Participation	yes	yes	2010	yes
7011379	2	Participation	yes	Missing by design	2011	yes
7011379	3	Participation	yes	Missing by design	2012	yes

The file CohortProfile contains all target persons of the panel sample. These are all targets with an initial agreement to participation. For each respondent in each wave, the CohortProfile contains meta information like the ID of the institution (ID\_i), various variables indicating participation (tx80220), availability of survey (tx80521), or availability of test data (tx80522). In addition, there are variables of the dates when the competence tests (testm/y) and the interview (intm/y/d) took place.

**In general, we strongly recommend using this file as a starting point for any analysis!**

**Example 3 (Stata):** Working with CohortProfile (find R example here)

```
** open the data file
use ${datapath}/SC5_CohortProfile_D_${version}.dta, clear
```



```
** change language to english (defaults to german)
label language en

** how many different respondents are there?
distinct ID_t

** as you can see, in this file there is an entry for every
** respondent in each wave
tab wave

** check participation status by wave
tab wave tx80220
```

## 4.2.4 EditionBackups

[« go back to overview](#)

## Description

Backup of original data that were modified during the data edition process

## File structure

long format: 1 row = 1 changed value of a variable in a datafile

ID variables needed to identify a single row

dataset varname ID\_t wave splink subspell partner child

Other ID variables useful for linkage

mergevars

Number of variables / number of rows in file

14 / 1,848

Contains data from waves



## Exemplary variables

dataset	Dataset name
varname	Variable name
mergevars	ID-Variables for merging
sourcevalue_num	Original value (if numeric)
editvalue_num	New value (if numeric)
sourcevalue_str	Original value (if string)
editvalue_str	New value (if string)
ID_t	ID target
wave	Wave

## Exemplary data snapshot

dataset	varname	mergevars	sourcevalue_num	editvalue_num	ID_t	wave
spVocExtExam	ts15304	ID_t wave exam	28.00	18.00	7004987	9
spVocExtExam	ts15304	ID_t wave exam	28.00	18.00	7006307	9
pTargetCATI	t731306	ID_t wave	5.00	2.00	7006357	1
pTargetCATI	t731306	ID_t wave	5.00	3.00	7006753	1
spVocExtExam	ts15304	ID_t wave exam	28.00	19.00	7018354	10

The dataset EditionBackups consists of single values that have been changed or modified in the data edition process. These single values can potentially originate from all other datasets. EditionBackups contains both the original and the changed value of a particular variable in a particular data file (i. e., one change or edition per row). The following variables are provided for each change:

- varname and dataset specify the name of the variable affected by an edition and the respective data file
- mergevars lists the identifier variables that are required to merge the information back to the respective data file

- `sourcevalue_[num/str]` contains the original, unaltered value; variables with the suffix `_num` refer to values from numeric variables and variables with the suffix `_str` refer to values from string variables (if the variable is numeric, `_str` is used to store the value label for this value instead)
- `editvalue_[num/str]` contains the result of the modification, i. e. the value into which the original value was changed; these values correspond exactly to the values in the respective data file (again, there is a version for both numeric and string variables - or the label).
- `ID_t`, `wave`, ... are the different identifier variables needed to merge the original values to the respective data files

### Example 4 (Stata): Working with EditionBackups

```
** In this example, we want to restore the original
** values in variable tg51410 (Intended degree) in datafile pTarget

** open the datafile
use ${datapath}/${cohort}_EditionBackups_D_${version}.dta, clear

** only keep rows containing data of the aforesaid variable
keep if dataset=="pTargetCAWI" & varname=="tg51410"

** check which variables we need for merging
tab mergevars

** then keep the merging variables and the variable with
** the original values (for cross-checking, we also keep the
** variable editvalue, which contains the values found in pTarget)
keep ID_t wave sourcevalue_num editvalue_num

** rename the variables to emphasize affiliation
rename sourcevalue_num tg51410_source
rename editvalue_num tg51410_edit

** temporary save this data extract
tempfile edition
save `edition'

** open pTargetCAWI
use ID_t wave tg51410 using ${datapath}/${cohort}_pTargetCAWI_D_${version}.dta, clear

** add the above data
merge 1:1 ID_t wave using `edition', keep(master match)

** check all edition made
list ID_t wave tg51410* if _merge==3, nolab

** replace the variable in the datafile with its original value
replace tg51410=tg51410_source if _merge==3
```

## 4.2.5 Education

[« go back to overview](#)

## Description

Generated: upward transitions in educational careers

## File structure

spell format: 1 row = 1 event (episode) of 1 respondent

ID variables needed to identify a single row

ID\_t splink

Other ID variables useful for linkage

tx28100

Number of variables / number of rows in file

12 / 55,313

Contains data from waves



## Exemplary variables

ID_t	ID target
number	Sort number
datem	Valid since (month)
datey	Valid since (year)
tx28101	Recent CASMIN
tx28102	Years of education = f(CASMIN)
tx28103	Recent ISCED-97
tx28109	Change in educational classification
splink	Link for spell merging
exam	Exam number
tx28100	Source of information of educational qualification

## Exemplary data snapshot

ID_t	number	datey	tx28101	tx28102	tx28103	splink	tx28100
7001974	1	2003	0	-20	0	220001	22
7001974	2	2007	3	10	2	220002	22
7001974	3	2010	5	13	3	220003	22
7001974	4	2015	8	18	9	240001	24
7001975	1	1999	0	-20	0	220001	22
7001975	2	2005	3	10	2	220002	22
7001975	3	2006	5	13	3	220003	22
7001975	4	2008	6	15	6	240001	24

This generated file provides longitudinal information on transitions in respondents' educational careers. It contains only persons who have an educational degree at a lower secondary level or higher. We used all information on educational attainment from spSchool (lower, intermediate, and upper secondary school degrees – Hauptschule, Realschule, (Fach-)Abitur), spVocPrep (participation in vocational preparation schemes), and spVocTrain (all successfully completed trainings). Also, data from spVocExtExam and spSchoolExtExam have been integrated. Three measures of educational attainment are available: CASMIN (variable tx28101), ISCED-97 (tx28103), and years of education (tx28102; derived from CASMIN). You can easily

merge data from the original spells to Education using the variable splink. The file stores transitions in a long event time format. That is, each row represents a transition in at least one classification (CASMIN and/or ISCED-97). Variables on month and year of the transition (datem and datey) specify the event time. We considered only upward educational transitions in CASMIN levels and upward as well as lateral transitions in ISCED-97 levels (CASMIN is ordinal, whereas ISCED-97 has some nominal elements). Because ISCED-97 and CASMIN follow different concepts, some educational transitions are captured by only one of these classifications.

### Example 5 (Stata): Working with Education (find R example here)

```
** we want to merge the school type from spSchool to this datafile.
** For this to work, we first have to prepare spSchool and keep only
** harmonized episodes (subspell==0)
use ${datapath}/SC5_spSchool_D_${version}.dta, clear
label language en
keep if subspell==0
tempfile temp
save `temp'

** now, open the Education data file
use ${datapath}/SC5_Education_D_${version}.dta, clear

** change language to english (defaults to german)
label language en

** check out which spell modules you can merge to this file
tab tx28100

** check that you will need splink to merge information
** from other modules to this file
** (command gives no result, which means approval)
isid ID_t splink, miss

** merge the previously generated temporary data file
merge 1:1 ID_t splink using `temp', keep(master match) keepusing(ts11204)

** see that this only added information to the rows corresponding to spSchool
tab tx28100 _merge
```

## 4.2.6 MethodsCATI

[« go back to overview](#)

## Description

Paradata from the targets CATI interview

## File structure

long format: 1 row = 1 target in 1 wave

ID variables needed to identify a single row

ID\_t wave

Other ID variables useful for linkage

ID\_int

Number of variables / number of rows in file

39 / 128,671

Contains data from waves



## Exemplary variables

ID_t	ID target
ID_int	Interviewer: ID
wave	Wave
intm	Interview date: month
inty	Interview date: year
tx80400	Willingness: Panel participation
tx80108	Type of recruitment
tx80301	Interviewer: Gender
tx80302	Interviewer: Age group
tx80209	Interview: Length of interview (minutes)
tx80401	Willingness: Merging data from federal employment agency
tx80304	Interviewer: Working experience as interviewer for infas

## Exemplary data snapshot

ID_t	ID_int	wave	intm	inty	tx80301	tx80302
7001968	1028	1	4	2011	2	50-65 years
7001968	1405	3	-54	-54	2	up to 29 years
7001969	1111	1	2	2011	1	50-65 years
7001969	-54	3	-54	-54	.	.

This dataset offers a variety of information on the data collection, e. g., gender (tx80301) and age (tx80302) of the interviewer; interview date (intm, inty); interview duration (tx80209); incentives (tx80210); and individual survey participation (tx80220).

Importantly, MethodsCATI contains all contacted respondents whether an interview was realized or not. Thus, MethodsCATI includes more cases than pTargetCATI.

**Example 6 (Stata):** Working with MethodsCATI (find R example here)

```

** open the data file
use ${datapath}/SC5_MethodsCATI_D_${version}.dta, clear

** change language to english (defaults to german)
label language en

** check out participation status by wave
tab wave tx80220

```

```
** how many different interviewers did CATI surveys?  
distinct ID_int  
  
** create one single variable containing the interview date  
generate intdate=mdy(intm,intd,inty)  
format intdate %td  
list intd intm inty intdate in 1/10
```

## 4.2.7 MethodsCAWI

[« go back to overview](#)

## Description

Paradata from the targets CAWI interview

## File structure

long format: 1 row = 1 target in 1 wave

ID variables needed to identify a single row

ID\_t wave

Other ID variables useful for linkage

none

Number of variables / number of rows in file

21 / 21,974

Contains data from waves



## Exemplary variables

ID_t	ID target
wave	Wave
inty	Interview date: year
intm	Interview date: month
intd	Interview date: day
tx80208	Interview: Length of questionnaire (minutes)
tx80225	Interview: last delivery status
tx80250	Interview: winners of the lottery
tx80200	Interview: Number of all contact attempts
tx80206	Interview: Number of interruptions
tx80207	Interview: Response Code differentiated (final outcome)

## Exemplary data snapshot

ID_t	wave	inty	intm	tx80208	tx80250
7001982	11	2016	11	14.11	0
7001982	14	-54	-54	-54.00	-54
7002030	11	2016	11	25.96	0
7002030	14	2018	11	18.14	-54
7002115	11	2016	11	26.18	0
7002115	14	2018	11	21.80	-54
7002191	11	2016	11	19.14	0

This dataset offers a variety of information on the data collection, e. g., interview date (intm, inty); interview duration (tx80208); winners of the prize draw (tx80250); and individual survey participation (tx80220).

Importantly, MethodsCAWI contains all contacted respondents whether an interview was realized or not. Thus, MethodsCAWI includes more cases than pTargetCAWI.

## Example 7 (Stata): Working with MethodsCAWI

```

** open the data file
use ${datapath}/SC5_MethodsCAWI_D_${version}.dta, clear

** change language to english (defaults to german)
label language en

```



```
** check out participation status by wave
tab wave tx80220

** how many waves have CAWI method data?
tab wave

** create one single variable containing the interview date
generate intdate=mdy(intm,intd,inty)
format intdate %td
list intd intm inty intdate in 1/10
```

## 4.2.8 MethodsCompetencies

[« go back to overview](#)

## Description

Paradata from the targets competency tests

## File structure

long format: 1 row = 1 target in 1 wave

ID variables needed to identify a single row

ID\_t wave

Other ID variables useful for linkage

ID\_i ID\_tg ID\_int

Number of variables / number of rows in file

73 / 29,419

Contains data from waves



## Exemplary variables

ID_t	ID target
ID_i	Institution ID
wave	Wave
ID_tg	ID test group
testm	Test: Survey day (month)
testy	Test: Survey day (year)
tx80422	Survey mode (realized)
ID_int	Interviewer: ID
tx80301	Interviewer: Gender
tx80302	Interviewer: Age group
tx80303	Interviewer: Highest school leaving qualification
tx80661	Number Participants
tx80628	Questions about test tasks
tx80629	Questions about the calculator
tx80633	Signature

## Exemplary data snapshot

ID_t	wave	ID_int	tx80301	tx80302	tx80303
7002115	12	-54	-54	Missing by design	-54
7002163	1	1385	1	30-49 years	7
7002163	12	-54	-54	Missing by design	-54
7002189	1	1385	1	30-49 years	7
7002189	12	-54	-54	Missing by design	-54
7002204	1	1318	2	30-49 years	2
7002204	5	1466	2	50-65 years	18
7002204	12	-54	-54	Missing by design	-54

Parallel to other Methods files, this dataset contains information about the testing situation, like durations, dates, interviewer IDs (ID\_int), information about the interviewer (e. g., sex (tx80301), age (tx80302), and education (tx80303)), individual survey participation (tx80220), number of participants (tx80661), and disruptions and influences during testing (tx80619).

**Example 8 (Stata):** Working with MethodsCompetencies (find R example here)

```

** open the data file
use ${datapath}/SC5_MethodsCompetencies_D_${version}.dta, clear

** how many respondents have been tested together in a group
bysort ID_tg: generate groupsize=_N if ID_tg>0 & !missing(ID_tg)
summarize groupsize

```

```
** create duration of math test; to achieve this, you first have to edit
** both start and end variables (which are stored in time format h:mm)

foreach var in tx80603 tx80604 { // do the following for both variables
** convert to string, add leading zero
  tostring `var', gen(`var'_str) format(%04.0f)
** generate the etc datetime (ms. since 01jan1960 00:00:00.000)
** take care of missing values!
  gen `var'_ms=clock(`var'_str,"hm") if `var'>0 & !missing(`var')
}
** now the duration is the subtraction of start from end.
** this is recoded then from miliseconds to minutes
generate duration = (tx80604_ms - tx80603_ms)/(60*1000)

summarize duration
```

## 4.2.9 pTargetCATI

[« go back to overview](#)

Description	Exemplary variables
Data from respondents CATI questionnaires	ID_t ID target
File structure	ID_i Institution ID
long format: 1 row = 1 target in 1 wave	wave Wave
ID variables needed to identify a single row	t431000 Migration sentiment
ID_t wave	t531214 Tuition loan
Other ID variables useful for linkage	t724403 Post-recording final grade
ID_i	t531250 Source of finance: family
Number of variables / number of rows in file	tg24503 Employment context doctorate
945 / 88,294	t712001 Kindergarten
Contains data from waves	t700001 Gender
1 2 3 4 5 6 7 8 9 10 11	t70000y Date of birth (year)
12 13 14	t514001 Satisfaction with life
	t514008 Satisfaction with course of study
	t741001 Size of household
	t520003 Weight in kg
Exemplary data snapshot	
ID_t wave t724403 tg24503 t700001 t70000y t514008	
7012681 12 1.7 5 [w] female 1991 6	
7014459 13 ..0 5 [m] male 1990 9	
7017128 12 1.4 4 [w] female 1991 9	
7019256 13 1.1 3 [m] male 1992 7	
7019306 10 ..0 5 [m] male 1991 7	

The data in file pTargetCATI are from computer assisted telephone interviews (CATI). As many questions are asked repeatedly over different waves, data integration follows a long data format. This means, for each wave participated, there is an additional line for each participating target in this wave. Therefore, targets are uniquely identified by ID\_t but lines are unique identified by ID\_t and wave together. As there are only lines within pTargetCATI for persons who responded, there are less lines in pTargetCATI than in CohortProfile.<sup>2</sup>

This file contains hundreds of variables, which is the gross of all items surveyed. Some of them are sociodemographic like gender (t700001), year of birth (t70000y), country of birth (t405010\_g2), or spoken languages (t414000\_g2). Others are repeatedly administered in different waves (e. g., financial means for studying (t531260), satisfaction with studies (t514008)).

<sup>2</sup> includes all students of the panel sample regardless of their questionnaire participation.

### Example 9 (Stata): Working with pTargetCATI (find R example [here](#))

```
** open the CohortProfile
use ${datapath}/SC5_CohortProfile_D_${version}.dta, clear

** merge some variable from pTargetCATI
merge 1:1 ID_t wave using ${datapath}/SC5_pTargetCATI_D_${version}.dta, ///
    keepusing(t400500_g1 t525204) nogen assert(master match)

** note that this information is now available only in waves which have
** surveyed the topic
tab wave t400500_g1

** if it makes sense, you can copy this information to cells of other waves.
** This copies information downwards (i.e., to later waves), unless a new
** value has been reported (which is usually what you want in a panel study)
bysort ID_t (wave): replace t400500_g1=t400500_g1[_n-1] ///
    if t400500_g1==54 | missing(t400500_g1)

tab wave t400500_g1
```

## 4.2.10 pTargetCAWI

[« go back to overview](#)

## Description

Data from respondents CAWI questionnaires

## File structure

long format: 1 row = 1 target in 1 wave

ID variables needed to identify a single row

ID\_t wave

Other ID variables useful for linkage

ID\_i

Number of variables / number of rows in file

1,408 / 54,468

Contains data from waves



## Exemplary variables

ID_t	ID target
wave	Wave
ID_i	Institution ID
t242020	Quality Facilities/equipment: Literature
t242107	Higher education institution activities: sport
t289902	Living in shared living
t514001	Satisfaction with life
t272061	Motivation for courses/trainings
t30300b	Amount of rent
tg51004	Degree course canceled/interrupted/completed
tg74011	Time budget: doctorate work
t241011	Time budget semester: Courses
t241012	Time budget semester: Self-study
tg72313	Discourse participation: lectures

## Exemplary data snapshot

ID_t	wave	t289902	t514001	t30300b	tg51004
7003919	8	1	3	150	2
7009386	8	1	9	300	2
7013307	8	1	8	350	2
7020948	8	1	8	455	3
7027276	8	1	5	300	3

Apart from computer assisted telephone interviews (CATIs), data collection via computer assisted web interviews (CAWIs) has been conducted. pTargetCAWI also covers similar constructs collected in the CATI. There are items related to the amount of rent (t30300b), satisfaction with life (t514001), having a roommate (t289902), and there are also variables to help you to identify if a target is currently studying (tg51000, tg51001, tg51004). In contrast to CATIs, CAWIs are self-administered. Furthermore, biographical data such as episodes of employment or episode of vocational training were not collected.

Note for variables tg5911\* (screen size): please find more information about those variables via codebook, infoquery, or NEPSplorer (see section 1.2 and section 1.8).

**Example 10 (Stata):** Working with pTargetCAWI (find R example here)

```
** open pTargetCAWI
```

```
use ${datapath}/SC5_pTargetCAWI_D_${version}.dta, clear

** only keep a single variable, and IDs
keep ID_t wave t289902

** suppose you want to know if somebody ever lived with roommates.
** Then you could make use of the expression "t289902==1", which is true (1)
** if there has been a roommate, or false (0) otherwise. The maximum of
** this expression over waves results in 1 if any wave ever evaluated to true,
** and 0 otherwise.
egen roommate = max(t289902==1), by(ID_t)

** only keep this variable; as all waves contain the same information, we
** can fall back to cross-sectional structure
keep ID_t roommate
duplicates drop
tempfile room
save `room', replace

** finally, open CohortProfile and merge this variable
use ${datapath}/SC5_CohortProfile_D_${version}.dta, clear
merge m:1 ID_t using "`room'", nogen

tab wave roommate
```

## 4.2.11 pTargetMicrom

[« go back to overview](#)

## Description

Small-scale regional indicators on respondents' place of residence

## File structure

panel format: 1 row = 1 regional level in 1 wave of 1 respondent

ID variables needed to identify a single row

ID\_t wave regio

Other ID variables useful for linkage

ID\_regio

Number of variables / number of rows in file

188 / 197,552

Contains data from waves



## Exemplary variables

ID_t	ID target
wave	Wave
regio	Indicator for enrichment level
ID_regio	System-free ID of enrichment level
mso_k_ausland	Share foreigners
mso_k_familie	Family structure
mbe_k_haustyp	Type of house
mgm_k_dom	Dominant microm geo milieu®
mgs_k_dom	Dominant geo-submilieu
mmo_k_volumen	Move volume
mpi_k_dichte	Car density
mas_k_berufsuv	Occupational disability insurance
mas_k_krankzuv	Additional health insurance
mlt_k_primit	Primary Limbic Type
kkw_w_summe	Total purchasing power in euro

## Exemplary data snapshot

ID_t	wave	regio	ID_regio	mso_k_ausland	mbe_k_haustyp	mpi_k_dichte
7009879	7	1	145167	8	6	1
7009879	7	2	239686	7	.	2
7009879	7	3	305174	8	.	2
7009879	7	4	426799	7	.	.
7009879	7	5	503553	9	.	2

The data file pTargetMicrom is only available via **On-site** access. The file is not included in the Download and Remote versions of the Scientific Use File.

The data include details about the respondent's residence at five different regional levels, distinguishable by the variable regio: house area, street section, postal code, postal code 8, municipality. All these levels are available for each respondent and each wave. Numerous regional indicators are provided, e.g. the percentage of foreigners, unemployment rate, family and age structure, milieu types, car type density, distribution of insurances, etc. To clarify, this information does **not** refer to individuals, but to regional units to which respondents belong via their



place of residence. Accordingly, the unemployment rate, for example, indicates the proportion of unemployed people in the population of a given region.

Please note that a separate documentation exists for this data file on the website (see section 1.2), which not only lists all variables, but also explains the background of the data.

### Example 11 (Stata): Working with pTargetMicrom (find R example here)

```
** open Microm datafile. Note that this data file is only available OnSite!
use ${datapath}/${cohort}_pTargetMicrom_0_${version}.dta, clear
label language en

** additionally to ID_t and wave, line identification in this file is done
** via variable regio, denoting the regional level of information
isid ID_t wave regio

** tabulating wave against regio shows availability of all levels
** in wave 5 and 7, but only the most detailed level available
** in wave 1 and 3 (usually housing level)
tab wave regio

** only keep housing level
keep if regio==1

** now you can enhance CohortProfile with regional data
merge 1:1 ID_t wave using ${datapath}/${cohort}_CohortProfile_0_${version}.dta
```

## 4.2.12 spChild

[« go back to overview](#)

## Description

information about all children of respondent

## File structure

entity format: 1 row = 1 child of 1 respondent

ID variables needed to identify a single row

ID\_t child subspell

Other ID variables useful for linkage

wave

Number of variables / number of rows in file

49 / 10,871

Contains data from waves



## Exemplary variables

ID_t	ID target
child	Child number
subspell	Number of subspell
wave	Wave
ts3320m	month of child's birth/ year of child's birth
ts3320y	year of child's birth
ts33203	Gender of the child
ts33204	Biological, adoptive or foster child
ts33209	Employment Child
ts33216	Vocational training Child

## Exemplary data snapshot

ID_t	child	subspell	wave	ts3320y	ts33203	ts33204
7002513	1	1	7	2014	[w] female	Biological child
7004924	1	1	3	2012	[w] female	Biological child
7005474	1	1	7	2014	[m] male	Biological child
7009819	1	1	12	2016	[w] female	Biological child
7012966	2	1	1	2003	[m] male	Biological child

This module contains information on all biological, foster, and adopted children of the respondent, and any other child that currently lives or has ever lived together with the respondent (e. g., children of former and current partners). In cases of twins and higher orders of multiple births, separate episodes are generated for each child. Episodes generally refer to the periods in which the respondent and the child shared a household. The enumerator variable `child` identifies children within respondents. Note that a child episode was skipped in the interview if the respondent reported that the child was deceased. Spell data on cohabitation with children is stored in file `spChildCohab` and spell data on parental leaves relating to children is stored in `spParLeave`.

**Example 12 (Stata):** Working with spChild (find R example here)

```
** open the data file
use ${datapath}/SC5_spChild_D_${version}.dta, clear
```

```
** switch to english language
label language en

** only keep full or harmonized episodes
keep if subspell==0

** generate the total count of children for each respondent
** you can do this either by taking the maximum child number:
bysort ID_t: egen children=max(child)
** or counting the number of rows:
bysort ID_t: gen children2=_N
** which both computes the same result
assert children==children2

** recode rough values (e.g. end of year) to real months
replace ts3320m=ts3320m-20 if ts3320m>20

** compute the age of one's children today
** first, create a Stata monthly date (months since 1960m1) of the birth variables
generate birth_ym =ym(ts3320y,ts3320m)
** then, create the same for the current date
gen now_ym=mofd(date(c(current_date), "DMY"))
** the age is then easily computed
gen age=(now_ym-birth_ym)/12

summarize age
```

## 4.2.13 spChildCohab

[« go back to overview](#)

## Description

file listing cohabitation spells with children

## File structure

spell format: 1 row = 1 cohabitation time of 1 respondent

## ID variables needed to identify a single row

ID\_t spell subspell

## Other ID variables useful for linkage

child wave

## Number of variables / number of rows in file

20 / 3,492

## Contains data from waves



## Exemplary variables

ID_t	ID target
child	Child number
spell	Spell number cohabitation with child
subspell	Number of subspell
wave	Wave
ts3331m	Start date Living together Child (month)
ts3331y	Start date Living together Child (year)
ts3332m	End date Living with child
ts3332y	End date Living with child
ts3332c	Currently living together with child
ts33308	Episode update Living together with child

## Exemplary data snapshot

ID_t	child	spell	subspell	wave	ts3331y	ts3332y
7004555	1	101	2	5	2011	2013
7005697	2	202	1	3	2011	2012
7012202	1	101	2	5	2011	2013
7014421	1	101	1	3	2011	2012
7018070	1	101	1	3	2011	2012

If a respondent lives together with children, durations are registered in spChildCohab. Cohabitation spells are related to children by the child number. Please note that those durations do not necessarily match birth and death events; rather see spChild for direct information on children.

**Example 13 (Stata):** Working with spChildCohab (find R example here)

```

** open the data file
use ${datapath}/SC5_spChildCohab_D_${version}.dta, clear

** switch to english language
label language en

** only keep full or harmonized episodes

```

```
keep if subspell==0

** recode rough values (e.g. end of year) to real months
replace ts3331m=ts3331m-20 if ts3331m>20
replace ts3332m=ts3332m-20 if ts3332m>20

** generate the following durations in months:
* a) the total duration of a cohabitation episode
gen cohab_duration = ym(ts3332y,ts3332m) - ym( ts3331y, ts3331m)
* b) the total duration a respondent lived together with specific child
bysort ID_t child (spell): egen total_duration_per_child = total(cohab_duration)
* c) the total duration a respondent lived together with any child
bysort ID_t (child spell): egen total_duration_per_target = total(cohab_duration)

** to work with the latter information in other files, you could do
** which gives you a cross-sectional display of cohabitation time for every
    respondent
keep ID_t total_duration_per_target
duplicates drop
```

## 4.2.14 spCourses

[« go back to overview](#)

Description

dynamic course module

File structure

spell format: 1 row = 1 episode of 1 respondent

ID variables needed to identify a single row

ID\_t wave splink

Other ID variables useful for linkage

sptype course\_w1 course\_w2 course\_w3

Number of variables / number of rows in file

31 / 5,402

Contains data from waves

1

2

3

4

5

6

7

8

9

10

11

12

13

14

Exemplary data snapshot

ID_t	wave	splink	sptype	course_w1	course_w2	course_w3
7009062	10	260004	26	1001	1002	1003
7010948	13	260007	26	1301	1302	1303
7013761	12	260003	26	1201	1202	1203
7017216	12	260009	26	1201	1202	1203
7018676	5	260006	26	501	502	503

Exemplary variables

ID\_tID target

waveWave

splinkLink for spell merging

sptypeSpell type

t271001Total duration of training courses

course\_w1Course number

t271011\_w1Course duration

course\_w2Course number

t271011\_w2Course duration

course\_w3Course number

t27800aStart date episode (month)

t27800bStart date episode (end)

t27800cEnd date Episode (month)

t27800dEnd date Episode (year)

This module comprises courses and trainings attended within the past 12 months during episodes of employment (spEmp), unemployment (spUnemp), parental leave (spParLeave), military, or civilian service (spMilitary), as well as episodes from the spGap module. The starting and end dates of the spells in this module represent the original episodes (in which a course was taken) from those modules. For each of these episodes, information on up to three courses is included in wide format. spCourses comprises all spells from the past 12 months that were recorded in the modules mentioned above. Spells may also be included if no course was taken during this episode. The only criterion for inclusion in the module is that a person provided information on at least one course. Note that in spCourses, the course enumerator is stored in wide format (course\_w1, course\_w2, and course\_w3), whereas in the other course modules (spFurtherEdu1 and spFurtherEdu2) there is only a single enumerator (course). Please note that this information has been integrated into datafile Education. If your interest in this data is not too profound, you are best advised to use Education instead.

### Example 14 (Stata): Working with spCourses (find R example [here](#))

```
** open the data file
use ${datapath}/SC5_spCourses_D_${version}.dta, clear

** check which modules provided course information
tab sptype

** only keep courses from employment spells
keep if sptype==26

** save this datafile for later usage
tempfile courses
save `courses'

** open the employment module
use ${datapath}/SC5_spEmp_D_${version}.dta, clear

** add the temporary datafile from above;
** note that this is an m:1 merge, as there are still subspells in spEmp
merge m:1 ID_t wave splink using `courses', assert(master match) nogenerate

** you now have the spEmp datafile, enhanced with information from spCourses,
** and can proceed with this in the usual way
```

## 4.2.15 spEmp

[« go back to overview](#)

## Description

spell data on employment episodes

## File structure

spell format: 1 row = 1 episode of 1 respondent

## ID variables needed to identify a single row

ID\_t spell subspell

## Other ID variables useful for linkage

wave

## Number of variables / number of rows in file

166 / 132,838

## Contains data from waves



## Exemplary variables

ID_t	ID target
splink	Link for spell merging
subspell	Number of subspell
spell	Spell number
wave	Wave
ts23222	In partial retirement, (active phase)
ts2311y	Start Employment episode (year)
ts2312y	End Employment episode (year)
ts23410	Net income, open
ts23228	Type of education required
ts23201_g1	Professional title (KldB 1988)
ts23201_g2	Professional title (KldB 2010)
ts23201_g3	Professional title (ISCO-88)

## Exemplary data snapshot

ID_t	subspell	spell	ts2311y	ts2312y	ts23410	ts23228
7008128	1	7	2015	2016	2200	8
7010229	5	2	2008	2016	400	1
7015780	1	4	2015	2016	1200	9
7017751	4	5	2014	2018	1800	3
7018811	1	5	2016	2016	2000	8

This extensive module covers all spells of regular employment, including traineeships. Information on second jobs is only collected for activities that continue up to the date of the interview. Vacation jobs, volunteering, and internships are not included. New episodes are created at the following events:

- Change of employer
- Change of occupation
- Interruption of employment (e. g., unemployment or military service)

The file comprises information like professional position (ts23203), net income (ts23410), relevance to degree course (tg26190), or permanent contract (ts23320).

**Example 15 (Stata):** Working with spEmp (find R example here)



```
** open the data file
use ${datapath}/SC5_spEmp_D_${version}.dta, clear

** only keep full or harmonized episodes
keep if subspell==0

** save this file temporarily
tempfile tmp
save `tmp'

** open the Biography data file
use ${datapath}/SC5_Biography_D_${version}.dta, clear

** merge the previously created temporary data file to this
merge 1:1 ID_t splink using `tmp' , keep(master match)

** you now have an enhanced version of Biography, enriched by
** information from the spell module. The number of total episodes
** (i.e., the amount of rows in the Biography file) did not change.
** Verify this by tabulating the spell type by the merging variable
** generated during the merge process.
tab sptype _merge
```

## 4.2.16 spFurtherEdu1

[« go back to overview](#)

## Description

information about additional courses

## File structure

entity format: 1 row = 1 course of 1 respondent

## ID variables needed to identify a single row

ID\_t course

## Other ID variables useful for linkage

wave

## Number of variables / number of rows in file

11 / 5,192

## Contains data from waves



## Exemplary variables

ID_t	ID target
wave	Wave
course	Course number
t271048	Course is ongoing
t271049	Termination Course
t272000_O	Course content other course
t271050	other courses 2
t271051	other course
t272000_g13	Content other course (course ID)
t271043	Duration of course
tx80211	Survey/Test instrument

## Exemplary data snapshot

ID_t	wave	course	t271048	t271050	t271051
7007152	13	1302	no	no	no
7013532	13	1307	no	no	no
7015467	5	501	no	no	yes
7017824	13	1304	no	no	yes
7018312	12	1204	yes	no	no

This module contains information on further courses (also private courses) attended within the past 12 months that have not been reported in spCourses or in spVocTrain. These include both professional trainings (similar to those from spCourses) and courses attended for private purposes (e. g., cookery course, yoga course, fortune telling, NLP coaching). Please note that this information has been integrated into datafile Education. If your interest in this data is not too profound, you are best advised to use Education instead.

**Example 16 (Stata):** Working with spFurtherEdu1 (find R example [here](#))

```

** open the datafile
use ${datapath}/SC5_spFurtherEdu1_D_${version}.dta, clear

** one row contains information for one course.
** The only possibility to use this file is to merge it to the data for this
** respondents wave (we use CohortProfile). So first, we have to remodel
** the file so one row contains one wave. We do this by Stata's reshape command

```

```
bysort ID_t wave (course): gen course_nr=_n
reshape wide course t*, i(ID_t wave) j(course_nr)

** create a temporary datafile for later merge
tempfile spfurther
save `spfurther'

** open CohortProfile
use `${datapath}/SC5_CohortProfile_D_${version}.dta, clear

** merge data
merge 1:1 ID_t wave using `spfurther', assert(master match) nogen

** Please note that you now have multiple variables added to CohortProfile,
** one set of variables for each course reported in spFurtherEdu1
```

## 4.2.17 spFurtherEdu2

[« go back to overview](#)

## Description

information about courses

## File structure

entity format: 1 row = 1 course of 1 respondent

ID variables needed to identify a single row

ID\_t course

Other ID variables useful for linkage

wave

Number of variables / number of rows in file

28 / 11,053

Contains data from waves



## Exemplary data snapshot

ID_t	wave	course	t279046	t279041	t272043
7013271	13	1301	fully	a lot of effort	Confirmation of attendance
7013844	13	1301	fully	no effort at all	Certificate
7014150	5	501	partially	a lot of effort	Confirmation of attendance
7014551	5	501	fully	some effort	Confirmation of attendance
7015030	12	1202	fully	some effort	Confirmation of attendance

## Exemplary variables

ID_t	ID target
wave	Wave
course	Course number
t279040	Professional/personal reasons
t279046	Course costs Employer
t272040	Provider
t279041	Motivation for course attendance
t272043	Certificate
t272003	Course assessment: learned new things
t274022	Course assessment: teacher patient
t279047	Course costs Employment agency

The survey instrument randomly selected two courses from the spCourses and spFurtherEdu1 modules, collecting additional information on these courses (e. g., costs incurred by employer t279046, motivation t279041, and certificates t272043). These data are included in spFurtherEdu2. Please note that this information has been integrated into datafile Education. If your interest in this data is not too profound, you are best advised to use Education instead.

**Example 17 (Stata):** Working with spFurtherEdu2 (find R example here)

```

** Two possibilities to use spFurtherEdu2

** A) Merge data to spCourses

** open spCourses datafile
use ${datapath}/SC5_spCourses_D_${version}.dta, clear

```

```
** one row contains information for up to three courses.  
** To make merging possible, you first have to reshape the datafile  
** so one row contains only one course  
reshape long course_w, i(ID_t wave splink) j(course_nr)  
rename course_w course  
  
** merge spFurtherEdu2 using ID_t and course  
merge m:1 ID_t course using ${datapath}/SC5_spFurtherEdu2_D_${version}.dta, keep(  
    master match)  
  
** ----  
** B) merge to spFurtherEdu1  
  
** open spFurtherEdu1 datafile  
use "${datapath}/SC5_spFurtherEdu1_D_${version}.dta", clear  
  
** merge spFurtherEdu2 using ID_t and course  
merge 1:1 ID_t course using ${datapath}/SC5_spFurtherEdu2_D_${version}.dta, keep(  
    master match)
```

## 4.2.18 spGap

[« go back to overview](#)

## Description

reported gap episodes

## File structure

spell format: 1 row = 1 gap of 1 respondent

ID variables needed to identify a single row

ID\_t spell subspell

Other ID variables useful for linkage

wave splink

Number of variables / number of rows in file

27 / 15,157

Contains data from waves



## Exemplary variables

ID_t	ID target
splink	Link for spell merging
spell	Spell number
subspell	Number of subspell
wave	Wave
spms	Check module: Type of event (edited)
ts29901	Auxiliary variable current gap episode
ts29300	Episode mode
ts2911m	Start date Gap (month)
ts2911y	Start date Gap (year)
ts2912m	End date Gap (month)
ts2912y	End date Gap (year)
ts2912c	Gap ongoing
ts29201	Training course during gap
ts29101	Type of gap episode

## Exemplary data snapshot

ID_t	spell	subspell	wave	ts29901	ts2911y	ts2912y
7003330	1	1	1	1	2011	2011
7004355	2	0	1	1	2010	2010
7004528	2	1	1	1	2011	2011
7005015	2	1	1	1	2011	2011
7005064	4	1	1	1	2011	2011

Gaps in individual life courses are identified by a check module. Such gap episodes are included in the spGap module. The spells in this file refer to different types of gaps that can be distinguished by the variable ts29101 (Type of gap episode). The most common gap episode is (extended) holidays.

**Example 18 (Stata):** Working with spGap (find R example here)

```

** open the data file
use ${datapath}/SC5_spGap_D_${version}.dta, clear

** only keep full or harmonized episodes
keep if subspell==0

** save this file temporarily
tempfile tmp

```

```
save `tmp'

** open the Biography data file
use ${datapath}/SC5_Biography_D_${version}.dta, clear

** merge the previously created temporary data file to this
merge 1:1 ID_t splink using `tmp' , keep(master match)

** you now have an enhanced version of Biography, enriched by
** information from the spell module. The number of total episodes
** (i.e., the amount of rows in the Biography file) did not change.
** Verify this by tabulating the spell type by the merging variable
** generated during the merge process.
tab sptype _merge
```

## 4.2.19 spInternship

[« go back to overview](#)

## Description

reported internship episodes

## File structure

spell format: 1 row = 1 internship episode of 1 respondent

## ID variables needed to identify a single row

ID\_t spell subspell

## Other ID variables useful for linkage

wave splink

## Number of variables / number of rows in file

45 / 38,963

## Contains data from waves



## Exemplary variables

ID_t	ID target
splink	Link for spell merging
spell	Spell number
subspell	Number of subspell
wave	Wave
tg3607m	Start month internship episode
tg3607y	Start year internship episode
tg3608m	End month Internship episode
tg3608y	End year Internship episode
tg36109	Continuing of internship episode
tg36110	Type of internship
tg36111	Average working hours Internship
tg36119	Placement as an intern
t265321	Learning content: Autonomy 1
t264300	Support: Supervision

## Exemplary data snapshot

ID_t	spell	subspell	wave	tg3607y	tg3608y	tg36111
7003864	1	1	3	2012	2012	40
7005610	3	1	3	2011	2012	6
7012075	4	1	9	2015	2015	35
7013949	3	2	7	2013	2013	12
7018159	1	2	7	2013	2013	45

As internships during studies are regarded as central to professional success, both compulsory and voluntary internships have been surveyed and made available in this datafile. Information about duration, remuneration, learning content, and other key aspects have been surveyed.

**Example 19 (Stata):** Working with spInternship (find R example here)

```

** open the data file
use ${datapath}/SC5_spInternship_D_${version}.dta, clear

** only keep full or harmonized episodes
keep if subspell==0

** save this file temporarily
tempfile tmp
save `tmp'

```



```
** open the Biography data file
use ${datapath}/SC5_Biography_D_${version}.dta, clear

** merge the previously created temporary data file to this
merge 1:1 ID_t splink using `tmp' , keep(master match)

** you now have an enhanced version of Biography, enriched by
** information from the spell module. The number of total episodes
** (i.e., the amount of rows in the Biography file) did not change.
** Verify this by tabulating the spell type by the merging variable
** generated during the merge process.
tab sptype _merge
```

## 4.2.20 spMilitary

[« go back to overview](#)

## Description

military / civilian service and voluntary gap years

## File structure

spell format: 1 row = 1 episode of 1 respondent

ID variables needed to identify a single row

ID\_t spell subspell

Other ID variables useful for linkage

wave splink

Number of variables / number of rows in file

25 / 5,123

Contains data from waves



## Exemplary variables

ID_t	ID target
splink	Link for spell merging
subspell	Number of subspell
spell	Spell number
wave	Wave
ts21201	Type of military service episode
ts2111m	Start Military service episode - month
ts2111y	Start Military service episode - year
ts2112m	End Military service episode - month
ts2112y	End Military service episode - year
ts21202	Attendance of training courses/courses during military service
ts2111y_g1	Check module: start date (year, edited)
ts2112y_g1	Check module: end date (year, edited)

## Exemplary data snapshot

ID_t	splink	subspell	spell	wave	ts2111y	ts2112y
7007254	250002	1	2	10	2016	2016
7008772	250001	1	1	1	2009	2011
7011558	250002	1	2	5	2012	2013
7012831	250001	2	1	3	2009	2012
7014759	250001	1	1	3	2011	2012

This module includes episodes of military or civilian service as well as gap years taken to do voluntary work in the social or environmental sector. Regular or professional soldiers are considered employed and are therefore included in the employment module.

**Example 20 (Stata):** Working with spMilitary (find R example here)

```
** open the data file
use ${datapath}/SC5_spMilitary_D_${version}.dta, clear
```

```
** only keep full or harmonized episodes
keep if subspell==0

** save this file temporarily
tempfile tmp
save `tmp'

** open the Biography data file
use `${datapath}/SC5_Biography_D_${version}.dta, clear

** merge the previously created temporary data file to this
merge 1:1 ID_t splink using `tmp' , keep(master match)

** you now have an enhanced version of Biography, enriched by
** information from the spell module. The number of total episodes
** (i.e., the amount of rows in the Biography file) did not change.
** Verify this by tabulating the spell type by the merging variable
** generated during the merge process.
tab sptype _merge
```

## 4.2.21 spParLeave

[« go back to overview](#)

## Description

episodes of parental leave

## File structure

spell format: 1 row = 1 parental leave episode of 1 respondent

## ID variables needed to identify a single row

ID\_t spell subspell

## Other ID variables useful for linkage

wave child splink

## Number of variables / number of rows in file

27 / 2,442

## Contains data from waves



## Exemplary variables

ID_t	ID target
child	Child number
spell	Spell number
subspell	Number of subspell
wave	Wave
ts2711m	Start Parental leave
ts2711y	Start Parental leave
ts2712m	End Parental leave (month)
ts2712y	End Parental leave (year)
ts2712c	Ongoing of parental leave

## Exemplary data snapshot

ID_t	child	spell	subspell	wave	ts2711y	ts2712y
7002912	1	101	4	13	2014	2017
7005423	2	202	2	13	2016	2017
7006121	1	101	2	7	2012	2014
7008138	2	202	3	12	2015	2016
7013977	1	101	2	5	2011	2012

For each child in spChild (except for deceased children), information is collected on whether the respondent took a parental leave. Each parental leave episode contributes one record to spParLeave. Parental leaves do not include maternity protection. These periods are added to the corresponding employment episode. As a result, an employment spell is not interrupted if the mother only takes the maternity leave without an additional parental leave.

**Example 21 (Stata):** Working with spParLeave (find R example here)

```

** open the data file
use ${datapath}/SC5_spParLeave_D_${version}.dta, clear

** only keep full or harmonized episodes
keep if subspell==0

```

```
** save this file temporarily
tempfile tmp
save `tmp'

** open the Biography data file
use ${datapath}/SC5_Biography_D_${version}.dta, clear

** merge the previously created temporary data file to this
merge 1:1 ID_t splink using `tmp' , keep(master match)

** you now have an enhanced version of Biography, enriched by
** information from the spell module. The number of total episodes
** (i.e., the amount of rows in the Biography file) did not change.
** Verify this by tabulating the spell type by the merging variable
** generated during the merge process.
tab sptype _merge
```

## 4.2.22 spPartner

[« go back to overview](#)

## Description

history of partners in the household

## File structure

entity format: 1 row = 1 partner of 1 respondent

## ID variables needed to identify a single row

ID\_t partner subspell

## Other ID variables useful for linkage

wave

## Number of variables / number of rows in file

109 / 67,167

## Contains data from waves



## Exemplary variables

ID_t	ID target
partner	Partner number
subspell	Number of subspell
ts31204	Partner: born Germany/abroad
ts31211	Partner German
ts31203	Gender of partner
ts3141m	Marriage date (month)
ts3141y	Marriage date (year)
ts3120y	year of birth partner
tg2811m	Start date Partnership month
tg2811y	Start date Partnership year
tg2804m	End date Partnership episode (month)
tg2804y	End date Partnership episode (year)
ts31206	Age at the time of moving to Germany Partner
ts31207	Place of birth of partner's father
ts31209	Place of birth of partner's mother

## Exemplary data snapshot

ID_t	partner	subspell	ts31203	ts3120y	tg2811m	tg2811y	tg2804m	tg2804y
7013068	1	1	[m] male	1968	7	1994	6	2011
7016246	2	0	[w] female	1988	12	2013	6	2014
7018201	1	0	[m] male	1990	2	2012	6	2012
7011886	1	0	[m] male	1983	5	2012	3	2013
7005584	1	0	[m] male	1990	5	2011	2	2012

This module covers the partnership history of the respondent. Respondents' subjective reports define whether they live in a relationship and whether they cohabit or not. A comprehensive set of additional questions refers to the present partner. For earlier partners, only information on the year of birth and education is available. Information on the current partner is collected regardless of the cohabitation status, whereas previous partners are only included if they cohabitated with the respondent. The enumerator variable partner identifies partners *within* respondents. This variable is coded 1 for the first partner and counts upwards until the last (current) partner.

**Example 22 (Stata):** Working with spPartner (find R example here)

```
** open the data file
use ${datapath}/SC5_spPartner_D_${version}.dta, clear

** switch to english language
label language en

** only keep full or harmonized episodes
keep if subspell==0

** to find out if a respondent is or was ever been married,
** check out if the indicating variable ever stated a marriage
bysort ID_t: egen married = max(ts31410==1)

** look at the data
list ID_t partner ts31410 married in 1/20, sepby(ID_t)

** reduce the datafile, so you have one single row for each respondent
keep ID_t married
duplicates drop

** you now can save this datafile and merge it to, e.g., CohortProfile
tempfile married
save `married'
use ${datapath}/SC5_CohortProfile_D_${version}.dta, clear
merge m:1 ID_t using `married', nogen keep(master match)
```

#### 4.2.23 spSchool

[« go back to overview](#)

The diagram illustrates the relationship between the data structure, variables, and an exemplary data snapshot. It is divided into three main sections: Description, Exemplary variables, and Exemplary data snapshot.

**Description:**

- general schooling history**
- File structure**
- spell format: 1 row = 1 school episode of 1 respondent**
- ID variables needed to identify a single row**
- ID\_t spell subspell**
- Other ID variables useful for linkage**
- wave splink**
- Number of variables / number of rows in file**
- 78 / 47,058**
- Contains data from waves**

**Exemplary variables:**

- ID\_t**: ID target
- splink**: Link for spell merging
- subspell**: Number of subspell
- spell**: Spell number
- wave**: Wave
- ts11204**: Type of school
- ts1111m**: Start date month School episode
- ts1111y**: Start date year School episode
- ts1112m**: End month School episode
- ts1112y**: End year School episode
- ts11209**: School-leaving qualification
- ts11214**: Intended school-leaving qualification
- ts11218**: Final grade school-leaving certificate
- t724801**: 1st 'Abitur' subject
- t724802**: 2nd Abitur subject

**Exemplary data snapshot:**

ID_t	splink	subspell	spell	wave	ts1111y	ts1112y
7005793	220001	0	1	1	1998	2002
7010981	220002	0	2	1	2002	2010
7014127	220005	2	5	7	2012	2013
7019077	220003	0	3	1	2001	2008
7019499	220004	2	4	5	2011	2012

The diagram shows how the variables listed in the 'Exemplary variables' section map to the columns in the 'Exemplary data snapshot' table. For example, 'ID\_t' maps to the 'ID\_t' column, 'splink' maps to the 'splink' column, 'subspell' maps to the 'subspell' column, 'spell' maps to the 'spell' column, 'wave' maps to the 'wave' column, 'ts1111y' maps to the 'ts1111y' column, and 'ts1112y' maps to the 'ts1112y' column. The 't724801' and 't724802' variables are not represented in the snapshot table.

This module covers each respondent's general education history from school entry until the date of (anticipated) completion, including

- episodes of elementary schooling,
- completed episodes of secondary schooling that led to a school leaving certificate, and
- incomplete episodes of schooling that would have led to a school leaving certificate if they had been completed.

A new episode is generated only if the school type changes. That is, a change from one Gymnasium to another is not recorded. As a result, a single schooling episode may take place at more than one location. In such cases, only information on the last location is included. A new episode is generated at each school type change even if both schools offer the same certificate.



### Example 23 (Stata): Working with spSchool (find R example [here](#))

```
** open the data file
use ${datapath}/SC5_spSchool_D_${version}.dta, clear

** only keep full or harmonized episodes
keep if subspell==0

** save this file temporarily
tempfile tmp
save `tmp'

** open the Biography data file
use ${datapath}/SC5_Biography_D_${version}.dta, clear

** merge the previously created temporary data file to this
merge 1:1 ID_t splink using `tmp' , keep(master match)

** you now have an enhanced version of Biography, enriched by
** information from the spell module. The number of total episodes
** (i.e., the amount of rows in the Biography file) did not change.
** Verify this by tabulating the spell type by the merging variable
** generated during the merge process.
tab sptype _merge
```

## 4.2.24 spSchoolExtExam

[« go back to overview](#)

## Description

school exam certificates acquired outside of the regular German educational system

## File structure

entity format: 1 row = 1 exam of 1 respondent

ID variables needed to identify a single row

ID\_t exam

Other ID variables useful for linkage

wave

Number of variables / number of rows in file

28 / 812

Contains data from waves



## Exemplary variables

ID_t	ID target
wave	Wave
exam	Exam number
ts11300	Awarded qualification in Germany?
ts1130m	Date: month qualification was awarded
ts1130y	Date: year qualification was awarded
ts11302	Awarded school-leaving qualification
ts11300_g1	Awarded qualification in Germany? (edited)
ts11301_g1R	Country of school-leaving qualification
ts11301_g2	Country of awarded school-leaving qualification (categorized)

## Exemplary data snapshot

ID_t	wave	exam	ts11300	ts1130y	ts11302	ts11300_g1
7003025	1	1	1	2007	.	1
7004477	1	1	1	2005	.	1
7008285	1	2	1	2006	.	1
7010437	1	1	1	2003	.	1
7014263	7	2	1	2013	4	1

The file spSchoolExtExam comprises information about school exam certifications that have not been acquired through “regular” schooling in the German educational system. These can consist of:

- certificates that have been acquired abroad and were accredited by German authorities
- certificates that have been acquired in a German school as external examinee (i. e., without attending class lessons)
- certificates that are automatically awarded by advancing through grades in upper secondary education

**Example 24 (Stata):** Working with spSchoolExtExam (find R example here)

```
** aim of this example is to evaluate the age of the respondent
** at the exam

** first, we have to get the birth date of the respondent
use ${datapath}/SC5_pTargetCATI_D_${version}.dta, clear
keep if wave==1 // only first wave as this data is time-invariant
keep ID_t t70000m t70000y
label language en
tempfile temp
save `temp'

** now, open the data file
use ${datapath}/SC5_spSchoolExtExam_D_${version}.dta, clear
label language en

** merge the previously extracted birth dates
merge m:1 ID_t using `temp', keep(master match) nogenerate

** recode the two date variables (year, month) into one:
gen exam_date=ym(ts1130y,ts1130m)
gen birth_date=ym(t70000y,t70000m)
format *_date %tm

** calculate the age (in years)
gen age=(exam_date-birth_date)/12

** recode missings to .a, b,... (not necessarily needed)
nepsmis ts11302

** show some deviation
tabulate ts11302, summarize(age)
```

## 4.2.25 spSibling

[« go back to overview](#)

## Description

siblings of respondent

## File structure

entity format: 1 row = 1 sibling of 1 respondent

ID variables needed to identify a single row

ID\_t sibling

Other ID variables useful for linkage

wave

Number of variables / number of rows in file

11 / 26,933

Contains data from waves



## Exemplary variables

ID_t	ID target
wave	Wave
sibling	Sibling number
tx80211	Survey/Test instrument
tg3270m	month of sibling's birth
tg3270y	year of sibling's birth
tg32706	Is sibling still alive?
tg32708	Employment status siblings
tg32709	Unemployment Siblings
tg32711	Highest school-leaving qualification Siblings
tg32724	Sibling lives with parents

## Exemplary data snapshot

ID_t	wave	sibling	tg3270m	tg3270y	tg32708	tg32711
7003757	1	3	5	1986	full-time employed	5
7008347	1	1	4	1975	full-time employed	3
7008623	1	1	27	1980	full-time employed	5
7013003	1	3	6	1985	full-time employed	3
7015608	1	2	8	1991	full-time employed	3

The file spSibling contains all reported siblings of the respondent. Each sibling is stored in one row, containing information about the date of birth (tg3270m/y), employment status (tg32708), and highest degree (tg32711).

**Example 25 (Stata):** Working with spSibling (find R example here)

```

** aim of this example is to evaluate the number of older and younger
** siblings of a respondent

** first, we have to get the birth date of the respondent
use ${datapath}/SC5_pTargetCATI_D_${version}.dta, clear
keep if wave==1 // only first wave as this data is time-invariant
keep ID_t t70000m t70000y
label language en
tempfile temp
save `temp'

```

```
** now, open the spSibling data file
use ${datapath}/SC5_spSibling_D_${version}.dta, clear
label language en

** merge the previously extracted birth dates
merge m:1 ID_t using `temp', keep(master match) nogenerate

** recode the two date variables (year, month) into one:
gen sibling_bdate=ym(tg3270y,tg3270m)
gen target_bdate=ym(t70000y,t70000m)
format *_bdate %tm

** check the difference between the two
gen older=.
replace older=0 if sibling_bdate>target_bdate
replace older=1 if sibling_bdate<target_bdate
replace older=. if missing(sibling_bdate) | missing(target_bdate)

** care about twins. As we do not know the day (or even the hour),
** we can not know which is older. We set this for a missing thus.
replace older=. if (sibling_bdate==target_bdate)

** generate the total amount of older siblings
bysort ID_t: egen total_older=total(older)
** generate the total amount of younger siblings
bysort ID_t: egen total_younger=total(1-older)

** aggregate to a single line for each respondent.
** the file then is cross-sectional with ID_t the sole identifier
keep ID_t total*
duplicates drop
```

## 4.2.26 spUnemp

[« go back to overview](#)

## Description

spell data on unemployment episodes

## File structure

spell format: 1 row = 1 episode of 1 respondent

ID variables needed to identify a single row

ID\_t spell subspell

Other ID variables useful for linkage

wave splink

Number of variables / number of rows in file

33 / 17,294

Contains data from waves



## Exemplary variables

ID_t	ID target
subspell	Number of subspell
spell	Spell number
wave	Wave
ts2511m	Start Unemployment episode (month)
ts2511y	Start Unemployment episode (year)
ts2512m	End Unemployment episode
ts2512y	End Unemployment episode
ts25202	Receipt of unemployment benefits or support at the beginning
ts25203	registered unemployment at present/at end
ts25205	Number Job applications
ts25206	Invitation to job interviews
ts25207	Number Interviews

## Exemplary data snapshot

ID_t	subspell	spell	wave	ts2511m	ts2511y	ts2512m	ts2512y
7007302	1	1	9	10	2014	6	2015
7012349	2	2	13	8	2017	8	2017
7014067	2	1	12	3	2016	6	2016
7014796	2	1	12	4	2016	6	2016
7027070	1	3	7	10	2013	6	2014

This module includes all episodes of unemployment irrespective of whether a person was registered as unemployed or not. Questions on registration of unemployment and receipt of benefits refer to both the beginning and the end of an unemployment spell.

**Example 26 (Stata):** Working with spUnemp (find R example here)

```

** open the data file
use ${datapath}/SC5_spUnemp_D_${version}.dta, clear

** only keep full or harmonized episodes
keep if subspell==0

** save this file temporarily
tempfile tmp
save `tmp'

```

```
** open the Biography data file
use ${datapath}/SC5_Biography_D_${version}.dta, clear

** merge the previously created temporary data file to this
merge 1:1 ID_t splink using `tmp' , keep(master match)

** you now have an enhanced version of Biography, enriched by
** information from the spell module. The number of total episodes
** (i.e., the amount of rows in the Biography file) did not change.
** Verify this by tabulating the spell type by the merging variable
** generated during the merge process.
tab sptype _merge
```

## 4.2.27 spVocExtExam

[« go back to overview](#)

## Description

vocational education certificates acquired outside of the regular German educational system

## File structure

entity format: 1 row = 1 exam of 1 respondent

## ID variables needed to identify a single row

ID\_t exam

## Other ID variables useful for linkage

wave

## Number of variables / number of rows in file

30 / 2,406

## Contains data from waves



## Exemplary variables

ID_t	ID target
wave	Wave
exam	Exam number
ts15301_g1	Professional/specialization title (KldB 1988)
ts15301_g4	Professional/specialization title (ISCO-08)
ts15301_g6	Professional/specialization title (SIOPS-88)
ts1530m	End month External examination
ts1530y	End year External examination
ts15304	External examination qualification
ts15302	External examination in Germany/abroad
th28370	External examination preparation done abroad for at least one month

## Exemplary data snapshot

ID_t	wave	exam	ts1530m	ts1530y	ts15304
7002915	13	1	7	2017	Second Staatsexamen (in teaching)
7004203	13	1	4	2018	Second Staatsexamen (in teaching)
7009102	13	1	1	2018	Second Staatsexamen (in teaching)
7009105	13	1	3	2018	other type of leaving certificate from a Fachschule
7009138	13	1	11	2017	Second/Third Staatsexamen (not in teaching)

The file spVocExtExam comprises information about vocational training certifications that have not been received by “regularly” passing through the German vocational training system. These can consist of:

- certificates that have been acquired abroad and were accredited by German authorities
- certificates that have been acquired in a German vocational training exam as external examinee (i. e., without attending lessons or courses registered with German authorities)

This especially includes second and third state examinations for alumni of medicine and law studies.



### Example 27 (Stata): Working with spVocExtExam (find R example [here](#))

```
** aim of this example is to evaluate the age of the respondent
** at the exam

** first, we have to get the birth date of the respondent
use ${datapath}/SC5_pTargetCATI_D_${version}.dta, clear
keep if wave==1 // only first wave as this data is time-invariant
keep ID_t t70000m t70000y
label language en
tempfile temp
save `temp'

** now, open the data file
use ${datapath}/SC5_spVocExtExam_D_${version}.dta, clear
label language en

** merge the previously extracted birth dates
merge m:1 ID_t using `temp', keep(master match) nogenerate

** recode the two date variables (year, month) into one:
gen exam_date=ym(ts1530y,ts1530m)
gen birth_date=ym(t70000y,t70000m)
format *_date %tm

** calculate the age (in years)
gen age=(exam_date-birth_date)/12

** recode missings to .a, b,... (not necessarily needed)
nepsmis ts15304

** show some deviation
tabulate ts15304, summarize(age)
```

## 4.2.28 spVocPrep

[« go back to overview](#)

## Description

vocational preparation schemes

## File structure

spell format: 1 row = 1 episode of 1 respondent

ID variables needed to identify a single row

ID\_t spell subspell

Other ID variables useful for linkage

wave

Number of variables / number of rows in file

62 / 473

Contains data from waves



## Exemplary variables

ID_t	ID target
splink	Link for spell merging
spell	Spell number
subspell	Number of subspell
spgen	Generated spell
wave	Wave
ts13103	Program type
ts1311m	Start Vocational preparation (month)
ts1311y	Start Vocational preparation (year)
ts1312m	End month Vocational preparation
ts1312y	End year Vocational preparation
ts1312c	Ongoing of the vocational preparatory year
ts13201	Termination Vocational preparation

## Exemplary data snapshot

ID_t	spell	subspell	wave	ts1311m	ts1311y	ts1312m	ts1312y
7004950	1	2	10	6	2015	8	2015
7004953	2	2	9	7	2014	8	2014
7009208	1	1	9	9	2014	7	2015
7010126	1	1	3	9	2011	5	2012
7010246	1	1	3	3	2012	4	2012

This module comprises episodes of vocational preparation after general education, including

- pre-training courses,
- basic vocational training years, and
- work preparation courses of the employment agency.

Data were collected on the duration from taking up until completing a vocational preparation scheme, including possible intermissions.

**Example 28 (Stata):** Working with spVocPrep (find R example here)

```
** open the data file
use ${datapath}/SC5_spVocPrep_D_${version}.dta, clear
```

```
** only keep full or harmonized episodes
keep if subspell==0

** save this file temporarily
tempfile tmp
save `tmp'

** open the Biography data file
use ${datapath}/SC5_Biography_D_${version}.dta, clear

** merge the previously created temporary data file to this
merge 1:1 ID_t splink using `tmp' , keep(master match)

** you now have an enhanced version of Biography, enriched by
** information from the spell module. The number of total episodes
** (i.e., the amount of rows in the Biography file) did not change.
** Verify this by tabulating the spell type by the merging variable
** generated during the merge process.
tab sptype _merge
```

## 4.2.29 spVocTrain

[« go back to overview](#)

### Description

vocational education history

### File structure

spell format: 1 row = 1 episode of 1 respondent

### ID variables needed to identify a single row

ID\_t spell subspell

### Other ID variables useful for linkage

wave splink

### Number of variables / number of rows in file

191 / 127,356

### Contains data from waves



### Exemplary variables

ID_t	ID target
spell	Spell number
subspell	Number of subspell
ts15201	Type of vocational training
ts1511m	Start month apprenticeship episode
ts1511y	Start year apprenticeship episode
ts1512m	End month Vocational training episode
ts1512y	End year Vocational training episode
ts15215	Company size of training company
ts15219	Vocational qualification
ts15221	Intended vocational qualification

### Exemplary data snapshot

ID_t	spell	subspell	ts1511m	ts1511y	ts1512m	ts1512y
7003932	2	1	6	2012	5	2013
7010545	1	2	10	2010	7	2012
7004110	1	4	10	2010	10	2013
7012896	2	2	10	2012	7	2014
7016829	3	2	9	2014	6	2018

This module covers all further trainings, vocational and/or academic, that a respondent ever attended:

- vocational training and retraining
- training at technical schools such as schools of public health, full-time vocational schools (excluding basic vocational training years), other vocational schools, and master craftsmen's colleges
- training in specialized fields of medicine
- accredited training courses to receive licenses
- conferral of a doctorate or postdoctoral thesis

- tertiary education at universities, specialized colleges for higher education, colleges of advanced vocational studies, and colleges of advanced administrative and commercial studies. Note: Only the main subjects are surveyed. New episodes are generated if

- a main subject changes over the course of studies, or
- the attainable degree changes over the course of studies (e. g., from MA to teaching certification).

Episodes are continued in case of location changes unless the main subjects change as well.

Training courses for licenses are comparable to courses in the spCourses, spFurtherEdu1, and spFurtherEdu2 modules and can therefore be identified by the spell indicator course. This enumerator allows linking information about the few courses included in this module to the courses in those modules. Interruptions of vocational training spells, so-called vocational interruption episodes, are stored in wide format (be aware of this when working with harmonized spell data!).

### Example 29 (Stata): Working with spVocTrain (find R example here)

```
** open the data file
use ${datapath}/SC5_spVocTrain_D_${version}.dta, clear

** only keep full or harmonized episodes
keep if subspell==0

** save this file temporarily
tempfile tmp
save `tmp'

** open the Biography data file
use ${datapath}/SC5_Biography_D_${version}.dta, clear

** merge the previously created temporary data file to this
merge 1:1 ID_t splink using `tmp' , keep(master match)

** you now have an enhanced version of Biography, enriched by
** information from the spell module. The number of total episodes
** (i.e., the amount of rows in the Biography file) did not change.
** Verify this by tabulating the spell type by the merging variable
** generated during the merge process.
tab sptype _merge
```

## 4.2.30 Weights

[« go back to overview](#)

Description	Exemplary variables
Sample weights for various occasions	
File structure	
wide format: 1 row = 1 target	
ID variables needed to identify a single row	
ID_t	ID_t ID target
Other ID variables useful for linkage	ID_i Institution ID of sampling
ID_i	sLevel_R Stratification second level
Number of variables / number of rows in file	ID_cl Cluster ID
Contains data from waves	w_h Weight stratification first level
1 2 3 4 5 6 7 8 9 10 11	w_t1 Cross-sectional weight for wave 1
12 13 14	w_t2 Cross-sectional weight for wave 2
	w_t3 Cross-sectional weight for wave 3
	w_allWaves Longitudinal weight wave 1 to 12
	w_allCATI Longitudinal weight wave 1,3,5,7,9,10,12
	w_allCAWI Longitudinal weight wave 2,4,6,8,11
Exemplary data snapshot	
ID_t	ID_i
ID_cl	w_h
w_h	w_t1
w_t1	w_t2
w_t2	w_allWaves
w_allWaves	
7012677	1002094
7010868	1002042
7015369	1002011
7017316	1002090
7007500	1002042
80	6.28600
212	1.66700
132	6.28600
161	6.28600
34	6.28600
1.01588	0.77769
0.24525	0.19608
1.71001	1.78825
0.88280	0.96650
1.24167	1.21356
	0.51145
	0.15421
	1.84627
	0.97222
	1.50226

Weighting variables (starting with w\_) are included in the Weights dataset. Also, you find cluster (ID\_cl) and stratification (stratum) identifiers here. Given the quite complex structure of the sample, no final recommendations are at hand concerning the use of design and adjusted weights. More information about weight estimation can be found in Zinn et al., 2017. There are no general rules available on how the use of design or adjusted weights render any possible analysis more stable. Weights may possibly help to highlight important features of the analysis, or at least serve as a robustness check for the performed analysis.

**Example 30 (Stata):** Working with Weights (find R example here)

```

** open Weights datafile
use ${datapath}/SC5_Weights_D_${version}.dta, clear

```

```
** note that this file is cross-sectional, although the weights
** seem to contain panel logic
d w_t*

** only keep weight corresponding to all waves
keep ID_t w_t123456789

** create a "panel" logic, i.e., clone each row
expand 9

** then create a wave variable
bysort ID_t: gen wave=_n

** save as temporary file
tempfile weights
save `weights', replace

** open CohortProfile
use `${datapath}/SC5_CohortProfile_D_${version}.dta, clear

** and merge weight
merge 1:1 ID_t wave using `weights', nogen

** note that this weight is only non-zero if respondents participated in
** all waves
tab wave tx80220 if w_t123456789!=0
```

## 4.2.31 xEcoCAPI

[« go back to overview](#)

## Description

additional competencies for students of economics and business administration

## File structure

wide format: 1 row = 1 student

ID variables needed to identify a single row

ID\_t

Other ID variables useful for linkage

wave ID\_int

Number of variables / number of rows in file

136 / 600

Contains data from waves



## Exemplary variables

ID_t	ID target
wave	Wave
tx80921	Participation status: Economics-subsample
testm	Test: Survey day (month)
testy	Test: Survey day (year)
bas7mar1_c	Economic competence: Marketing 1
bas7_sc1	Economic competence: WLE
bas7_sc2	Economic competence: SE(WLE)
ID_int	Interviewer: ID
tg90308	Number Semesters Economics
tg24160_g2	Subject group Subject 1 (destatis 2010/11)
tx80200	Interview: Number of all contact attempts
tx80302	Interviewer: Age group
tx80430	Interview: Location

## Exemplary data snapshot

ID_t	wave	tx80921	bas7_sc1	bas7_sc2	tg90308	tg24160_g2	tx80200
7002199	7	6	0.69803	0.51537	7	3	3
7003481	7	6	0.80859	0.41941	2	3	2
7007758	7	6	0.33918	0.38575	5	3	4
7008650	7	6	1.94842	0.58791	7	3	3
7018018	7	6	0.05863	0.37350	6	3	2

Apart from the basic CATI-data collection in wave 7, additional data was collected for students of economics and business administration. A paper-based competency test containing questions specifically for the target's field of study was embedded within a short computer assisted personal interview (CAPI).

This data was part of pTargetCATI and xTargetCompetencies in releases prior to data version 10-0-0. To emphasize the focus on this small subgroup of targets, all this information is now gathered in xEcoCAPI. As this file contains data from wave 7 only, ID\_t is a unique identifier in this wide-format dataset. To make things simpler, participation in CAPI, CATI, and competency testing is indicated by tx80921. Additional methods data – like number of contact tries



(tx80200) and reasons for item-nonresponse in testing (e.g., tx80411) – are available as well. CAPI data are basically focussing on the student's area of studies (e.g., tg24160\_g2).

### Example 31 (Stata): Working with xEcoCAPI

```
** open the CohortProfile
use ${datapath}/SC5_CohortProfile_D_${version}.dta, clear

** merge some variables from xEcoCAPI
merge 1:1 ID_t wave using ${datapath}/SC5_xEcoCAPI_D_${version}.dta, ///
    keepusing(bas7_sc1 bas7_sc2) nogen assert(master match)

** note that this information is now available only in waves which have
** surveyed the topic
tab wave bas7_sc1
```

## 4.2.32 xInstitution

[« go back to overview](#)

## Description

context information about the institution

## File structure

wide format: 1 row = 1 area of studies in 1 institution

## ID variables needed to identify a single row

ID\_i tg04001\_g7

## Other ID variables useful for linkage

none

## Number of variables / number of rows in file

127 / 3,480

## Contains data from waves



## Exemplary data snapshot

ID_i	tg04001_g7	tg92104_O	tg92301_O	tg92601_R
1002057	8	1	1	3
1002008	8	1	1	3
1002069	4	1	1	3
1002017	6	1	1	3
1002106	2	1	1	3

## Exemplary variables

ID_i	Institution ID
tg04001_g7	Subject area WT 2010 (for merging with context data)
tg91102_R	HEI region: BLK-region type
tg92104_O	HEI: Winner Cluster of excellence 2006 or 2007
tg92301_O	HEI: Funding body
tg92601_R	HEI: Students 2010 total (aggr. Tercent. universities)
tg93204_O	SG: Students 2010: male
tg93205_O	SG: Students 2010: female
tg93601_O	SG: Students per professor
tg93602_R	SG: Students per lecturer (aggr. by terc. of all HEI, sep. U/UAS)

Data file xInstitution has been constructed during data edition of the first wave. At this time, information about the participating institutions (e. g., universities) has been collected. The file contains data on 10 area of studies for 322 institutions, e. g., about the university region, if the university has been winner or nominee of different prizes, the funding body, and number of students, lecturers, and professors. Note that due to data protection issues, this file is not available in the Download version of SUF. You find it in **RemoteNEPS** and **Onsite**. Please note that higher education context data are only available for winterterm 2010/11. The provision of panel data on higher education contexts is currently not planned.

**Example 32 (Stata):** Working with xInstitution (find R example here)

```

** open datafile
use ${datapath}/SC5_pTargetCATI_O_${version}.dta, clear

foreach var in ID_i tg04001_g7 { // do the following for both variables

```

```
** copy the information from the first wave downwards for each target,
** unless a new value has been reported
bysort ID_t: replace `var' = `var'[_n-1] ///
            if `var' == -54|missing(`var')
}
** drop all observations where no satisfaction with studies was reported
drop if t514008 == -98|t514008 == -97|t514008 == -93|t514008 == -54|missing(t514008)

** some respondents reported satisfaction with studies in 7th and in 9th waves
** to keep the latest information, create a seq and a max variables
bysort ID_t: gen seq = _n
bysort ID_t: gen max = _N
** only keep the latest reported information
keep if seq == max
** only keep the variables relevant for the merge and the analysis
keep ID_t ID_i tg04001_g7 t514008

** merge two variables from xInstitution
merge m:1 ID_i tg04001_g7 using ${datapath}/SC5_xInstitution_0_${version}.dta, ///
      keepusing(tg92601_R tg92104_0) nogen assert(master match)

** assuming that the less students at university the more intensive the support by
the
** university staff per student and the more satisfied are students with their
studies
** tabulate Satisfaction with studies by Students 2010 total
** note that the following analysis is feasible in both, RemoteNEPS and Onsite
tab t514008 tg92601_R, col

** assuming that students at excellence universities are more satisfied with
** their studies, tabulate the distribution of satisfaction by tg92104_0
** note that the following analysis is only feasible in the Onsite version of SUF,
** since the variable tg92104_0 is anonymized in RemoteNEPS
tab t514008 tg92104_0, col
```

## 4.2.33 xPlausibleValues

[« go back to overview](#)

## Description

Plausible Values of competence data

## File structure

wide format: 1 row = 1 respondent

ID variables needed to identify a single row

ID\_t

Other ID variables useful for linkage

wave\_w\*

Number of variables / number of rows in file

114 / 11,740

Contains data from waves



## Exemplary variables

ID_t	ID target
wave_w1	Row contains data from wave 1 (2010/2011 (CATI+competencies))
wave_w5	Row contains data from wave 5
wave_w12	Row contains data from wave 12
mas1_pv1	Math: cross-sectional plausible value 1
mas1_pv2	Math: cross-sectional plausible value 2
mas1_pv10	Math: cross-sectional plausible value 10
mas1_pv1u	Math: longitudinal plausible value 1
mas1_pv2u	Math: longitudinal plausible value 2
mas1_pv10u	Math: longitudinal plausible value 10

## Exemplary data snapshot

ID_t	wave_w1	mas1_pv1	mas1_pv2	mas1_pv10	mas1_pv1u
7003772	1	0.59433	0.73821	0.44791	0.73766
7003026	1	2.34065	1.36901	0.73805	2.88913
7011527	1	0.11436	0.85125	0.67075	0.80017
7003213	1	0.41126	1.53214	0.57179	1.96270
7015603	1	0.14944	0.73821	0.84500	1.05429

Plausible Values (PV) are a way of describing the competencies of individuals at the group level. They allow (unbiased) estimates of effects at the population level that are adjusted for measurement errors. In contrast to point estimators such as Weighted Likelihood Estimates (WLE), the use of Plausible Values is suitable for more precise inferential statistical tests in correlation and mean value analyses.

Plausible Values are based on the individual answers in the competence tests and additional background characteristics (e.g. gender, age, socioeconomic status). For each person, the probability distribution of his or her competence is first determined and then several values are randomly drawn from it (hence *Plausible Values*). Hypothesis tests for the specific question of interest are calculated for each of these values and combined into an overall result.

Please find more information on Plausible Values in the corresponding NEPS Survey Paper (Scharl, Carstensen, and Gnambs, 2020) and on our website:

→ [www.neps-data.de](http://www.neps-data.de) > Data Center > Overview and Assistance > Plausible Values

### Example 33 (Stata): Working with xPlausibleValues

```
** open datafile.
use ${datapath}/${cohort}_xPlausibleValues_D_${version}.dta, clear
label language en

** as the 'x' in the filename indicates, this is a cross sectional file
** (no wave structure). You can verify this by asking if one row is
** solely identified by the respondents ID
isid ID_t

** note that competence testing has been conducted in multiple waves.
** An indicator marks if a row contains information for a specific wave.
tab1 wave_w*

** see more on how to work with this data in the Survey Paper mentioned above!
```

## 4.2.34 xTargetCompetencies

[« go back to overview](#)

## Description

Test data of respondents

## File structure

wide format: 1 row = 1 target

ID variables needed to identify a single row

ID\_t

Other ID variables useful for linkage

wave\_w\*

Number of variables / number of rows in file

371 / 11,811

Contains data from waves



## Exemplary variables

ID_t	ID target
wave_w1	Row contains data from wave 1 (2010/2011 (CATI+competencies))
wave_w5	Row contains data from wave 5
mas1r092_c	Mathematical competence: Item 2
mas1_sc1	Mathematical competence: WLE
mas1_sc2	Mathematical competence: SE(WLE)
res1_sc1	Reading competence: WLE
res1_sc2	Reading competence: SE(WLE)
rsci0051_c	Reading speed: Item 51
rss1_sc3	Reading speed: Sum
ics3_sc1	ICT-Literacy WLE
ics3_sc2	ICT-Literacy SE of WLE

## Exemplary data snapshot

ID_t	wave_w1	wave_w5	mas1_sc1	mas1_sc2
7015422	1	1	1.88074	0.75218
7006171	1	1	0.43645	0.54009
7017846	1	1	0.60330	0.53224
7017369	1	1	0.44488	0.54007
7017733	1	1	1.30443	0.76587

File xTargetCompetencies contains data from competence assessments conducted. Scored item variables as well as scale variables are available in a cross-sectional format. Note that not all respondents took part in the assessment. Since assessments were conducted in CATI mode, those persons who were interviewed in CATI-mode have been excluded from testing. Additionally, those who had severe visual impairments or were even blind were excluded from the assessment.

**Example 34 (Stata):** Working with xTargetCompetencies (find R example here)

```

** open datafile
use ${datapath}/SC5_xTargetCompetencies_D_${version}.dta, clear

** change language to english (defaults to german)
label language en

```

```
** as the 'x' in the filename indicates, this is a cross sectional file
** (no wave structure). You can verify this by asking if one row is
** solely identified by the respondents ID
isid ID_t

** note that competence testing has been conducted in multiple waves
** an indicator marks if a row contains information for a specific wave
tab1 wave_w*

** to work with competence data, you might want to merge it to CohortProfile.
** if you want to keep the panel logic (and not only add all competencies
** to every wave), you need a mergeable wave variable in xTargetCompetencies.
** in this example, we focus on math competencies, which have been tested in wave 1.
generate wave=1

** now, remove cases which did not took part in the testing
drop if wave_w1==0

** and reduce the dataset to the relevant variables
keep ID_t wave mas1_sc1 mas1_sc2

** save a temporary datafile
tempfile tmp
save `tmp'

** and merge this to CohortProfile
use ${datapath}/SC5_CohortProfile_D_${version}.dta, clear
merge 1:1 ID_t wave using `tmp', nogen
```

## 5 Special Issues

### 5.1 Service Variables (Area of studies, ISCED-97 subject)

**subject of study** The variables tg2416\* were edited due to discrepancies between subspells. Subjects are filled for the first explicit mention only, missing information was labeled accordingly.

Currently the code -29 “*Value from last-mentioned sub-episode*” describes two cases: missing information can be found in the previous sub-spell or in the previous spell (the latter means a person started a new study-episode but claims that the subject is still the same as in the previously recorded episode).

The missing code -28 “*Value from recruitment pTargetCATI*” denotes that the missing information can be found in the recruitment data in file pTargetCATI.

The service variables tg2417\* contain the respective subject of study, thus the variables tg24170\_g1-5, tg24173\_g1-5, tg24176\_g1-5 provide complete subject information for all study episodes. Working with the service variables is recommended.

**type of university** The variable tg01003\_g1 (*type of university*, four levels) is originally a part of the first wave recruitment information contained in dataset pTargetCATI. The variable ts15201 (*type of vocational training program*, twenty-four levels) is part of the core education questionnaire and is recorded for each educational spell; it is part of spVocTrain. The service variable tg01003\_ha (*type of university*) provides an aggregated version of ts15201 in spVocTrain partly using information from tg01003\_g1 for first wave spells, as seen in table 7.

**Table 7:** Harmonization of type of university

tg01003_ha		tg01003_g1		ts15201	
1	University of applied sciences (incl. cooperative state university)	1	University of applied sciences (incl. cooperative state university)	6	Degree course at an administration and business academy (VWA)
				7	Degree course at a Berufsakademie/cooperative state university
				8	Degree course at a college of public administration
				9	Degree course at a university of applied sciences (not a college of public administration)
2	University	2	University	10	Degree course at a university, including college of education, art college, music college



**vocational education history** In waves 3, 5, and 7, an attempt has been made to retrospectively gather additional information about vocational education episodes that were concurrent with the first study episode of the winter term 2010/11. This has led to duplicate and/or right-censored episodes in the dataset spVocTrain. In order to deal with those episodes, the variable tx20100 was introduced to give a recommendation which episodes should be used for analyses. The rule applies that episodes from wave 1 are always recommended when the start date lies at or before the beginning of the first study episode of the winter term 2010/11. Episodes from wave 1 are never recommended when the start date lies after the beginning of the first study episode of the winter term 2010/11.

## 5.2 Coding subject of study

### 5.2.1 Recruitment

**data collection** Information on subject of study of initial studies was collected in PAPI and CATI mode (for information on sampling in SC5, see Aßmann et al., 2011, and Zinn et al., 2017). PAPI questionnaires were typewritten and delivered to NEPS by the data collecting institute (infas). Information on subject of study collected in first CATI was delivered to NEPS as original string variable.

**coding** Coding of subject of study was provided by the NEPS department *From Higher Education to the Labor Market* at DZHW Hannover (formerly HIS), based on data delivered by the data collecting institute (infas) from both modes (CATI and PAPI). The coding process faced a few challenges due to a change of the destatis classification between recruitment and first wave data collection: sampling was based on the destatis-classification of 2009/10 while the coding of recruitment information was based on the destatis-classification of 2010/11.

Coding was done manually by occasionally using additional information when a decision could not be taken only based on the string variable.

**classification used** The classification used for coding the recruitment information on subject of study is based on the destatis classification of 2010/11. Coding decisions can differ from destatis recommendations for coding degree programs into subjects of study due to individual decisions based on extensive research.

### 5.2.2 Panel Waves

**data collection** For higher education episodes reported after recruitment, the subject of study has been recorded using lists – in CATI as well as in online surveys. In cases where inter-

viewers were unable to fit a respondents answer into the respective list, the subject of study been recorded as an open string. Both in CATI and online panel waves, the lists are based on the destatis classification 2010/11 and the recruitment information.

To facilitate the allocation of respondent answers, the CATI-list has been continuously extended with supplementary information (based on open responses and changes in the academic landscape in Germany); the online list has remained the same.

Up until wave 13 subjective decisions in the maintenance of the CATI-lists and technical restrictions have led to deviations from the original classification. In some cases, subjects of study were assigned to different codes within the list. The idea behind this was for the other subjects within the same code to serve as covariates, so interviewers (and respondents) could choose the *right* list entry. Starting with wave 13, the CATI-lists will only be extended in the sense that new subject names will be added to the existing subject groups corresponding to a code if those subject names are not already listed under another code. The allocation will follow the coding rules described below to ensure consistency and transparency. This way, the list documented below will not be changed but will be enhanced over time. Starting with wave 14, online waves will use the CATI-list of the previous CATI-wave to harmonize the recording of subject of study in CATI and online mode.

**coding** Coding of open responses on subject of study has been provided by the NEPS department *From Higher Education to the Labor Market* for all panel waves so far. Since SUF 6.0.0 all strings that have been coded once have been collected in a reference list with their corresponding code by the LIfBi Research Data Center to avoid inconsistencies. In the following waves, open strings have been matched with that list first and strings in the list automatically get assigned the same code. Open strings that have been reported for the first time were coded manually until SUF 9.0.0. Starting with SUF 10.0.0, coding has followed a set of standardized rules and the software CODI has been used.

**classification used** Data collection and coding of subject of study largely follows the destatis 2010/11 classification of subjects of study.

**derivation of SUF-variables** In the Scientific Use File, several alternative variables containing subject of study are offered. Variables with the suffix `_g1R` and `_g2` contain the first four digits of the seven digit destatis 2010/11 classification (“Studienbereich” and “Fächergruppe”), `_g3R`, `_g4R` and `_g5` contain derivations of the destatis classification into different levels of the ISCED 97 classification. All derivations are based on the seven digit version of the destatis classification, using a transcoding table supplied by the Federal Statistical Office.

## 6 References

- Aßmann, C., Steinhauer, H. W., Kiesel, H., Koch, S., Schönberger, B., Müller-Kuller, A., ... Blossfeld, H.-P. (2011). Sampling designs of the National Educational Panel Study: Challenges and Solutions, 51–65. doi:10.1007/s11618-011-0181-8
- Blossfeld, H.-P., Roßbach, H. G., & von Maurice, J. (Eds.). (2011). Education as a Lifelong Process: The German National Educational Panel Study (NEPS). [*Special Issue*] *Zeitschrift für Erziehungswissenschaft*: 14.
- Dahm, G. (2014). *Starting Cohort 5 - Dokumentation der Variable tg24150\_g2 "NTS" (Nicht-traditionelle Studierende)*. DZHW - Deutsches Zentrum für Hochschul- und Wissenschaftsforschung GmbH.
- FDZ-LifBi. (2020). *Data Manual NEPS Starting Cohort 5– First-Year Students, From Higher Education to the Labor Market, Scientific Use File Version 14.0.0*. Bamberg, Germany: Leibniz Institute for Educational Trajectories, National Educational Panel Study.
- Lauterbach, O. (2015). *Erfassung wirtschaftswissenschaftlicher Fachkompetenzen von Studierenden in Startkohorte 5 des Nationalen Bildungspanels - Technischer Bericht* (NEPS Working Paper No. 51). Bamberg: Leibniz-Institut für Bildungsverläufe, Nationales Bildungspanel.
- NEPS (Ed.). (2020). *Starting Cohort 5: First-Year Students (SC5), Wave 14, Questionnaires (SUF Version 14.0.0)*. Bamberg, Germany: Leibniz Institute for Educational Trajectories, National Educational Panel Study.
- Pohl, S., & Carstensen, C. H. (2012). *NEPS Technical Report – Scaling the Data of the Competence Tests* (NEPS Working Paper No. 14). German National Educational Panel Study (NEPS). Bamberg.
- Scharl, A., Carstensen, C. H., & Gnambs, T. (2020). *Estimating Plausible Values with NEPS Data: An Example Using Reading Competence in Starting Cohort 6* (NEPS Survey Paper No. -). Bamberg, Germany: Leibniz Institute for Educational Trajectories, National Educational Panel Study.
- Schönberger, K., & Koberg, T. (2017). *Regional Data: Microm*. Bamberg, Germany: Leibniz Institute for Educational Trajectories, National Educational Panel Study.
- Steinwede, J., & Aust, F. (2012). *Methodenbericht, NEPS Startkohorte 5 – CATI-Haupterhebung Herbst 2010, B52*. Bonn, Germany: infas.
- Wenzig, K. (2012). *NEPS-Daten mit DOIs referenzieren*. RatSWD Working Paper Series. Rat für Sozial- und Wirtschaftsdaten, Berlin.
- Zielonka, M., & Pelz, S. (2015). *Implementation of the ISCED-97, CASMIN and Years of Education. Classification Schemes in SUF Starting Cohort 6*. Bamberg, Germany: Leibniz Institute for Educational Trajectories, National Educational Panel Study.

## References

- Zinn, S., Steinhauer, H. W., & Aßmann, C. (2017). *Samples, Weights, and Nonresponse: the Student Sample of the National Educational Panel Study (Wave 1 to 8)* (NEPS Survey Paper No. 18). Bamberg, Germany: Leibniz Institute for Educational Trajectories, National Educational Panel Study.

# A Appendix

## A.1 R examples

In this Appendix, you will find R usage examples that correspond to the Stata usage examples in the main body of the data manual. Before working with R, it is recommended to set a working directory and to install the package *readstata13*:

### Example 35 (R): Setting working directory

```
setwd("C:/User/.../Desktop/R_examples")
#set working directory

install.packages("readstata13")
#install the package readstata13 that reads Stata files
library(readstata13)
#import the package readstata13 into library
```

If you'd like to work with the English version of the data, it is recommended to switch the language in Stata first, save the Stata file and then import it in R. The language can be switched by running the command `label language en` in Stata.

To import a data set, use:

### Example 36 (R): Importing the data

```
'** here based on the example of the data set spEmp:'
spEmp = read.dta13("SC5_spEmp_D_version.dta", convert.factors = T)
#convert.factors = T converts value labels from Stata into factor label in R
#i.e. "1", "2" data class: integer becomes "yes", "no" data class: factor
```

The following step is not absolutely necessary. However it is recommended, if you attach great importance to keep the variable labels handy during your analysis. After importing the data set, you can display an overview over all variable labels by running the command `varlabel(spEmp)`. However, this command doesn't work anymore after modifying the data by e.g. deleting or merging variables, since the single variable labels aren't attached to the single variable names. To prevent that, following steps are necessary:

### Example 37 (R): Assigning variable labels

```
'** here based on the example of the data set spEmp:'

#install and integrates the package "Hmisc"
install.packages("Hmisc")
library(Hmisc)
```

```
#First, create a dataframe with all variable names and labels for spEmp
spEmp_meta = data.frame(attr(spEmp,"names"),attr(spEmp,"var.labels"))

#renames the columns in "names" and "labels"
colnames(spEmp_meta) = c("names", "labels")

spEmp_meta_names = as.vector(spEmp_meta$names)
#extracts the column "names" as vector "spEmp_meta_names"

spEmp_meta_labels = as.vector(spEmp_meta$labels)
#extracts the column "labels" as vector "spEmp_meta_labels"

names(spEmp_meta_labels) = spEmp_meta_names
#assigns the names to the labels, so that the vector "spEmp_meta_labels" is now a
  named vector
#this procedure produces the same result as the following command:
#spEmp_meta_labels = c(ID_t = "Target-ID", splink = "Link für Spell-Merging",
  subspell = "Teilepisodenummer", ... for all variables)

for(i in seq_along(spEmp)){
  label(spEmp[,i]) = spEmp_meta_labels[i]
}
#assigns variable labels that are stored in spEmp_meta_labels to the single columns

label(spEmp)
label(spEmp$subspell)
#Now the variable labels are assigned to the single columns
```

### Example 38 (R): Working with Basics

```
'** import the data files'
CohortProfile =
  read.dta13("SC5_CohortProfile_D_version.dta",
    convert.factors = T)

Basics =
  read.dta13("SC5_Basics_D_version.dta",
    convert.factors = T)

'** merge the data from Basics, enhancing every entry in CohortProfile'
CohortProfile = merge(CohortProfile, Basics, by = "ID_t", all = TRUE)
#The option all = TRUE makes sure that both, matched AND unmatched cases are kept
  during the merging process

'** tabulate gender by wave'
addmargins(table(Data$wave, Data$t700001))
```

### Example 39 (R): Working with Biography

```
'** import the data file'
Biography =
```

```
read.dta13("SC6_Biography_D_9-0-0.dta",
           convert.factors = T)

'** check out which spell modules you can merge to this file'
addmargins(table(Biography$sptype))

'** check that you will need splink to merge information
** from other modules to this file'
anyDuplicated(Biography[,c("ID_t","splink")])
#returns "0" if there are no duplicates.
#If there are duplicates this command returns the index of the first duplicate
```

### Example 40 (R): Working with CohortProfile

```
'** import the data file'
CohortProfile =
  read.dta13("SC5_CohortProfile_D_version.dta",
             convert.factors = T)

'** how many different respondents are there?'
length(unique(CohortProfile$ID_t))
#number of distinct ID_t

'** respondents in each wave'
cbind(addmargins(table(CohortProfile$wave)),
      addmargins(prop.table(table(CohortProfile$wave))))

'** check participation status by wave'
cbind(addmargins(table(CohortProfile$wave, CohortProfile$tx80220)))
```

### Example 41 (R): Working with Education

```
'** we want to merge the school type from spSchool to this datafile.
** For this to work, we first have to prepare spSchool and keep only
** harmonized episodes (subspell == 0)'
spSchool =
  read.dta13("SC5_spSchool_D_version.dta",
             convert.factors = T)

spSchool = subset(spSchool, spSchool$subspell == 0)

'** open the Education data file'
Education =
  read.dta13("SC5_Education_D_version.dta",
             convert.factors = T)

'** check which spell modules you can merge to this file'
table(Education$tx28100)

'** check that you will need splink to merge information
** from other modules to this file'
```

```
anyDuplicated(Education[,c("ID_t", "splink")])
#returns "0" if there are no duplicates.
#If there are duplicates this command returns the index of the first duplicate

'** merge spSchool to Education'
#After merging, Stata merge has one variable more than R, because in Stata
#a merge indicator is produced during the merging process and in R isn't.
#Since we need a merge indicator here, the merge command has to be extended:
Education = transform(merge(
  x = cbind(Education, source = "master"),
  #x contains the Education data set plus one extra column "source",
  #where source = "master"
  y = cbind(spSchool[,c("ID_t", "splink", "ts11204")], source = "using"),
  # y contains only the columns ID_t, splink and ts11204 from spSchool
  # plus one extra column "source" where source = "using"
  all.x = TRUE, by = c("ID_t", "splink")),
  # merges x and y by ID_t and splink
  source = ifelse(!is.na(source.x) & !is.na(source.y), "both",
    ifelse(!is.na(source.x), "master", "using")),
  # in the merged dataset, source = "both" if the observations is in x AND in y
  #otherwise, source = "master" if the obs. is only in x
  #and source = "using" if the obs. is only in y
  source.x = NULL,
  source.y = NULL
  # the columns "source" in x and y are deleted
)

'** see that this only added information to the rows corresponding to spSchool'
cbind(addmargins(table(Education$tx28100, Education$source)))
```

### Example 42 (R): Working with MethodsCATI

```
'** import the data file'
MethodsCATI =
  read.dta13("SC5-MethodsCATI_D_version.dta",
    convert.factors = T)

'** check out participation status by wave'
cbind(addmargins(table(MethodsCATI$wave, MethodsCATI$tx80220)))

'** how many different interviewers did CATI surveys?'
length(unique(MethodsCATI$ID_int))

'** create one single variable containing the interview date'
Sys.setlocale("LC_TIME", "C")
#turns off the location-specific language, so that the english months are recognized.

MethodsCATI$intdate =
  as.Date(paste(MethodsCATI$intm, MethodsCATI$intd, MethodsCATI$inty, sep = '-'),
    "%B-%d-%Y")
#binds the three columns "intm", "intd" and "inty" into one new column "intdate"
```



```
head(MethodsCATI[c("intd", "intm", "inty", "intdate"]), 10)
#displays first 10 rows of intd, intm, inty and intdate
```

### Example 43 (R): Working with MethodsCompetencies

```
'** open the data file'
MethodsCompetencies =
  read.dta13("SC5_MethodsCompetencies_D_version.dta",
    convert.factors = T)

'** how many respondents have been tested together in a group'
MethodsCompetencies = within(MethodsCompetencies,{
  groupsize = ave(ID_tg, ID_tg, FUN = length)})
#creates a new variable "groupsize" and counts the observations in each ID_tg group

#Problem: NEPS-Missings are also counted as regular values and summarized in groups
for (i in 1:length(MethodsCompetencies$ID_tg)) {
  if(!is.na(MethodsCompetencies$ID_tg[i]) & MethodsCompetencies$ID_tg[i] < 0){
    MethodsCompetencies$groupsize[i] = NA
    #sets all observations to NA for which ID_tg < 0 (here -55 and -54)
  }
}

summary(MethodsCompetencies$groupsize)
#displays Min, Max and Mean for "groupsize"
sd(MethodsCompetencies$groupsize, na.rm = TRUE)
#displays Std.Dev. for "groupsize"
length(MethodsCompetencies$groupsize[!is.na(MethodsCompetencies$groupsize)])
#displays the number of observations in "groupsize" without NA

'** create duration of math test'
for (t in names(MethodsCompetencies[,c(38, 39)])) {
  # run over columns 38 and 39 (variables tx80603 and tx80804)
  for (i in 1:length(MethodsCompetencies[[t]])) {
    #runs over every single observation
    if(nchar(MethodsCompetencies[[t]][i]) == 3 & MethodsCompetencies[[t]][i] > 0) {
      #if the observation length is 3 and positive (e.g., "923", but not "-54")
      MethodsCompetencies[[t]][i] = paste("0", MethodsCompetencies[[t]][i], sep = "")
      #adds a leading 0 character, such that 923 becomes 0923
    }
  }
}

install.packages("chron")
library(chron)
#package for creating chronological objects

for (i in names(MethodsCompetencies[,c(38, 39)])){
  MethodsCompetencies[[paste(i, 't', sep = "_")]] =
    #creates new variables tx80603_t and tx80604_t
```

```
times((strftime(strptime(MethodsCompetencies[[i]], format = "%H%M"), "%H:%M:%S")))
#assigns the values from tx80603 and tx80604 in time format to them
}

MethodsCompetencies$duration =
  MethodsCompetencies$tx80604_t - MethodsCompetencies$tx80603_t
#creates a new variable "duration", subtracting start time from end time

summary(MethodsCompetencies$duration)
#displays Min, Max and Mean for "duration" in time format
mean(MethodsCompetencies$duration) * 60 * 24
#displays the mean in minutes format
#one unit equals one day, thus it has to be multiplied by 60 minutes and 24 hours

sd(MethodsCompetencies$duration, na.rm = TRUE) * 60 * 24
#displays Std.Dev. for "duration" in minutes format
times(sd(MethodsCompetencies$duration, na.rm = TRUE))
#displays Std.Dev. in time format

length(MethodsCompetencies$duration[!is.na(MethodsCompetencies$duration)])
#displays the number of observations in "duration" without NA
```

### Example 44 (R): Working with pTargetCATI

```
'** open the CohortProfile dataset'
CohortProfile =
  read.dta13("SC5_CohortProfile_D_version.dta",
    convert.factors = T)

'** merge some variable from pTargetCATI'

pTargetCATI =
  read.dta13("SC5_pTargetCATI_D_version.dta",
    convert.factors = T)
#imports the pTargetCATI dataset

CohortProfile =
  merge(x = CohortProfile,
    y = pTargetCATI[,c("ID_t", "wave", "t400500_g1", "t525204")],
    by = c("ID_t", "wave"), all.x = TRUE)
#merges only variables "t400500_g1" and "t525204" from pTargetCATI to CohortProfile

'** note: this information is available only in waves which have surveyed the topic'
addmargins(table(CohortProfile$wave, CohortProfile$t400500_g1))

'** if it makes sense, you can copy this information to cells of other waves.
** This copies information downwards (i.e., to late waves), unless a new
** value has been reported (which is usually what you want in a panel study'
for (i in 2:length(CohortProfile$ID_t)) {
  if(CohortProfile$ID_t[i] == CohortProfile$ID_t[i-1]) {
    if(is.na(CohortProfile$t400500_g1[i]) |
      CohortProfile$t400500_g1[i] == "Missing by design") {
      CohortProfile$t400500_g1[i] = CohortProfile$t400500_g1[i-1]
```

```

    }
  }
}

addmargins(table(CohortProfile$wave, CohortProfile$t400500_g1))

```

### Example 45 (R): Working with pTargetCAWI

```

'** open the pTargetCAWI dataset'
pTargetCAWI = read.dta13("SC5_pTargetCAWI_D_version.dta", convert.factors = T)

'** only keep single variables and IDs'
pTargetCAWI = subset(pTargetCAWI, select = c(ID_t, wave, t289902))

'** suppose you want to know if somebody ever lived with roommates.
** t289902 == "Specified" if there has been a roommate,
** and t289902 == "Not specified" otherwise. The maximum of
** this expression over waves results in 1 if any wave ever evaluated to true,
** and 0 otherwise.'
for (i in 1:length(pTargetCAWI$ID_t)){
  if(pTargetCAWI$t289902[i] == "Specified")pTargetCAWI$roommate[i] = 1
  else pTargetCAWI$roommate[i] = 0
}

pTargetCAWI = within(pTargetCAWI, {roommate = ave(roommate, ID_t, FUN = max)})
#for every ID_t with at least one roommate == 1, all other roommate observations
#are also replaced by 1 within this ID_t.

'** only keep this variable; as all waves contain the same information, we
** can fall back to cross-sectional structure'
pTargetCAWI = subset(pTargetCAWI, select = c(ID_t, roommate))
pTargetCAWI = pTargetCAWI[!duplicated(pTargetCAWI),]

'** finally, open CohortProfile and merge this variable'
CohortProfile = read.dta13("SC5_CohortProfile_D_version.dta", convert.factors = T)
CohortProfile = merge(CohortProfile, pTargetCAWI, by = c("ID_t"), all = TRUE)
addmargins(table(CohortProfile$wave, CohortProfile$roommate))

```

### Example 46 (R): Working with pTargetMicrom

```

'** open pTargetMicrom datafile. Note that this data file is only available OnSite!'
pTargetMicrom = read.dta13("SC6_pTargetMicrom_O_version.dta", convert.factors = T)

'** additionally to ID_t and wave, line identification in this file is done
** via variable regio, denoting the regional level of information'
anyDuplicated(pTargetMicrom[,c("ID_t", "wave", "regio")])
#returns 0 if there are no duplicates
#If there are duplicates this command returns the index of the first duplicate

'** tabulating wave against regio shows availability of all levels

```

```

** in wave 5 and 7, but only the most detailed level available
** in wave 1 and 3 (usually housing level)'
addmargins(table(pTargetMicrom$wave, pTargetMicrom$regio))

'** only keep housing level'
pTargetMicrom = subset(pTargetMicrom, pTargetMicrom$regio == 1)

'** now you can enhance CohortProfile with regional data'
CohortProfile = read.dta13("SC6_CohortProfile_O_version.dta", convert.factors = T)
pTargetMicrom = merge(CohortProfile, pTargetMicrom, by = c("ID_t", "wave"), all =
  TRUE)

```

### Example 47 (R): Working with spChild

```

'** open the data file'
spChild = read.dta13("SC5_spChild_D_version.dta", convert.factors = T)

'** only keep full or harmonized episodes'
spChild = subset(spChild, spChild$subspell == 0)

'** generate the total count of children for each respondent
** you can do this either by taking the maximum child number:'
spChild = within(spChild, {children = ave(child, ID_t, FUN = max)})

'** or counting the number of rows:'
spChild = within(spChild, {children2 = ave(ID_t, ID_t, FUN = length)})

'** which both computes the same result'
identical(spChild$children, spChild$children2)

'** recode rough values (e.g., end of year) to real months'
levels(spChild$ts3320m)[levels(spChild$ts3320m) == "Beginning of the year/winter"] =
  "January"
levels(spChild$ts3320m)[levels(spChild$ts3320m) == "Spring/Easter"] = "April"
levels(spChild$ts3320m)[levels(spChild$ts3320m) == "Mid-Year/Summer"] = "July"
levels(spChild$ts3320m)[levels(spChild$ts3320m) == "Fall"] = "October"
levels(spChild$ts3320m)[levels(spChild$ts3320m) == "End of year"] = "December"

'** compute the age of 'ones children today
** first, create a date of the birth variables'
spChild$ts3320m = match(spChild$ts3320m, month.name)

install.packages("zoo")
library(zoo)
#the zoo package is needed to transform time data

#transforms month names into month numbers
spChild$birth_ym = as.yearmon(paste(spChild$ts3320y, spChild$ts3320m), "%Y %m")

'** then, create the same for the current date'
spChild$today_ym = as.yearmon(rep(cut(Sys.Date(), "month"), length(spChild$ID_t)))

'** the age is then easily computed'

```

```
spChild$age = (spChild$today_ym - spChild$birth_ym)

summary(spChild$age)
# displays Min, Max and Mean of "age"
sd(spChild$age, na.rm = TRUE)
# displays Std.Dev. of "age"
length(spChild$age[!is.na(spChild$age)])
# displays the number of observations in "age" without NA
```

### Example 48 (R): Working with spChildCohab

```
'** open the data file'
spChildCohab = read.dta13("SC5_spChildCohab_D_version.dta", convert.factors = T)

'** only keep full or harmonized episodes'
spChildCohab = subset(spChildCohab, spChildCohab$subspell == 0)

'** recode rough values (e.g., end of year) to real months'
for (i in names(spChildCohab[c(16, 18)])){
  #run over the variables ts3331m and ts3332m in columns 16 and 18
  levels(spChildCohab[[i]])[levels(spChildCohab[[i]]) == "Beginning of the year/
winter"] = "January"
  levels(spChildCohab[[i]])[levels(spChildCohab[[i]]) == "Spring/Easter"] = "April"
  levels(spChildCohab[[i]])[levels(spChildCohab[[i]]) == "Mid-Year/Summer"] = "July"
  levels(spChildCohab[[i]])[levels(spChildCohab[[i]]) == "Fall"] = "October"
  levels(spChildCohab[[i]])[levels(spChildCohab[[i]]) == "End of year"] = "December"
}

'** generate the following durations in months:
* a) the total duration of a cohabitation episode'
for (i in names(spChildCohab[c(16, 18)])) {
  spChildCohab[[i]] = match(spChildCohab[[i]], month.name)
  #transforms month names into month numbers
}

install.packages("zoo")
library(zoo)
#the zoo package is needed to transform time data

spChildCohab$cohab_start =
  as.yearmon(paste(spChildCohab$ts3331y, spChildCohab$ts3331m), "%Y %m")
spChildCohab$cohab_end =
  as.yearmon(paste(spChildCohab$ts3332y, spChildCohab$ts3332m), "%Y %m")

spChildCohab$cohab_duration =
  (spChildCohab$cohab_end - spChildCohab$cohab_start)*12

'* b) the total duration a respondent lived together with specific child'
spChildCohab = within(spChildCohab,
  {total_duration_per_child =
    ave(cohab_duration, ID_t, child, FUN =
      function(x) round(sum(x, na.rm = TRUE)))})
```

```
'* c) the total duration a respondent lived together with any child'
spChildCohab = within(spChildCohab,
  {total_duration_per_target =
    ave(cohab_duration, ID_t, FUN =
      function(x) round(sum(x, na.rm = TRUE)))})

'** to work with the latter information in other files, you could do
** which gives you a cross-sectional display of cohabitation time per respondent'
spChildCohab = subset(spChildCohab, select = c("ID_t", "total_duration_per_target"))
spChildCohab = spChildCohab[!duplicated(spChildCohab),]
```

### Example 49 (R): Working with spCourses

```
'** open the data file'
spCourses = read.dta13("SC5_spCourses_D_version.dta", convert.factors = T)

'** check which modules provided course information'
cbind(addmargins(table(spCourses$sptype)))

'** only keep courses from employment spells'
spCourses = subset(spCourses, spCourses$sptype == "Emp")

'** open the employment module'
spEmp = read.dta13("SC5_spEmp_D_version.dta", convert.factors = T)

'** merge spCourses to spEmp
** note that this is an m:1 merge, as there are still subspells in spEmp'
#Since the variable tx80211 is in both data sets spCourses AND spEmp
intersect(names(spCourses), names(spEmp))
#and since the variable is not one of the merging variables, both versions
#are contained in the new merged data set as tx80211.x and tx80211.y.

#To avoid that there are two possibilities:

#1. You can include the variable in the merging process by:
spEmp =
  merge(spEmp, spCourses, by = c("ID_t", "wave", "splink", "tx80211"), all.x = TRUE)
# In that case the version from the master data set, here spEmp, is kept

#OR

#2. If you'd like to compare the both versions first, you can merge the
#data sets as usual by:
spEmp =
  merge(spEmp, spCourses, by = c("ID_t", "wave", "splink"), all.x = TRUE)

#compare the two versions of the variable tx80211 by:
addmargins(table(spEmp$tx80211.x, spEmp$tx80211.y))

#and then drop one of the variables by:
spEmp$tx80211.y = NULL
```

```
'** you now have the spEmp datafile, enhanced with information from spCourses,
** and can proceed with this in the usual way'
```

### Example 50 (R): Working with spEmp

```
'** open the data file'
spEmp = read.dta13("SC5_spEmp_D_version.dta", convert.factors = T)

'** only keep full or harmonized episodes'
spEmp = subset(spEmp, spEmp$subspell == 0)

'** open the Biography data file'
Biography = read.dta13("SC5_Biography_D_version.dta", convert.factors = T)

'** merge the spEmp to Biography'
#After merging, Stata merge has one variable more than R, because in Stata
#a merge indicator is produced during the merging process and in R isn't.
#Since we need a merge indicator here, the merge command has to be extended:
Biography = transform(merge(
  x = cbind(Biography, source = "master"),
  #x contains the Biography data set plus one extra column "source",
  #where source = "master"
  y = cbind(spEmp, source = "using"),
  #y contains the spEmp data set plus one extra column "source",
  #where source = "using"
  all.x = TRUE, by = c("ID_t", "splink")),
  #merges x and y by ID_t and splink
  source = ifelse(!is.na(source.x) & !is.na(source.y), "both",
    ifelse(!is.na(source.x), "master", "using")),
  #in the merged dataset, source = "both" if the observations is in x AND in y
  #otherwise, source = "master" if the obs. is only in x
  #and source = "using" if the obs. is only in y
  source.x = NULL,
  source.y = NULL
  #the columns "source" in x and y are deleted
)

#Since the variables wave and spms are in both data sets, Biography AND spEmp
#check before merging by: intersect(names(Biography), names(spEmp))
#and since the variables are not part of the merging process,
#both versions are contained in the new merged data set as
#wave.x/wave.y and spms.x/spms.y
#For each variable, one of the versions can be dropped by:
Biography$wave.y = NULL
Biography$spms.y = NULL

'** you now have an enhanced version of Biography, enriched by
** information from the spell module. The number of total episodes
** (i.e., the amount of rows in the Biography file) did not change.
** Verify this by tabulating the spell type by the merging variable
** generated during the merge process.'
addmargins(table(Biography$sptype, Biography$source))
```

### Example 51 (R): Working with spFurtherEdu1

```
'** open the datafile'
spFurtherEdu1 = read.dta13("SC5_spFurtherEdu1_D_version.dta", convert.factors = T)

'** one row contains information for one course.
** The only possibility to use this file is to merge it to the data for this
** respondents wave (we use CohortProfile). So first, we have to remodel
** the file so one row contains one wave.'
spFurtherEdu1$course_nr = ave(spFurtherEdu1$ID_t, spFurtherEdu1$ID_t,
                             spFurtherEdu1$wave, FUN = seq_along)

spFurtherEdu1 = reshape(data = spFurtherEdu1,
                        #data in long format
                        idvar = c("ID_t", "wave"),
                        #idvar is/are the variable/s that need/s to be left unaltered
                        v.names = names(spFurtherEdu1[,3:11]),
                        #v.names contains names of variables in the long format that
                        #correspond to multiple variable in the wide format
                        timevar = "course_nr",
                        #timevar is/are the variable/s that need/s to be converted to
                        #wide format
                        direction = "wide")
                        #direction is to which format the data needs to be transformed

'** open CohortProfile'
CohortProfile = read.dta13("SC5_CohortProfile_D_version.dta", convert.factors = T)

'** merge the data'
CohortProfile =
    merge(CohortProfile, spFurtherEdu1, by = c("ID_t", "wave"), all.x = TRUE)
'** Please note that you now have multiple variables added to CohortProfile,'
'** one set of variables for each course reported in spFurtherEdu1'
```

### Example 52 (R): Working with spFurtherEdu2

```
'** Two possibilities to use spFurtherEdu2'

'-----'
'** A) Merge data to spCourses'

'** open spCourses datafile'
spCourses = read.dta13("SC5_spCourses_D_version.dta", convert.factors = T)

'** one row contains information for up to three courses.
** To make merging possible, you first have to reshape the datafile
** so one row contains only one course'
spCourses = reshape(data = spCourses,
                    # data in wide format
                    idvar = c("ID_t", "wave", "splink"),
                    #idvar is/are the variable/s that need/s to be left unaltered
```



```

varying = c("course_w1","course_w2","course_w3"),
#varying are the variables that need to be converted from
#wide to long
v.names = c("course"),
#v.names defines the name of the variable in that the in
#varying defined variables are summarized
times = c(1,2,3),
#new variable "time" is created with levels 1, 2 and 3
#for the three courses
new.row.names = 1:100000,
#sets row names as numeric
direction = "long"
##direction is to which format the data needs to be transformed
)

names(spCourses)[names(spCourses) == "time"] <- "course_nr"
#renames the variable "time" to "course_nr"

'** merge spFurtherEdu2 using ID_t and course'
#open spFurtherEdu2 datafile
spFurtherEdu2 = read.dta13("SC5_spFurtherEdu2_D_version.dta", convert.factors = T)

intersect(names(spCourses), names(spFurtherEdu2))
#common variables in the both data sets are "ID_t", "wave", "tx80211" and "course"
#Since the variables "wave" and "tx80211" are not part of the merging process,
#both versions are contained in the new merged data set
#as wave.x/wave.y and tx80211.x/tx80211.y.

'**To avoid that, there are two merging options:'
#1. You can include the variables in the merging process by:
spCourses =
  merge(spCourses, spFurtherEdu2,
        by = c("ID_t", "course", "wave", "tx80211"), all.x = TRUE)
#In that case the versions from the master data set are kept (wave.x and tx80211.x)

#OR

#2. If you'd like to compare the both versions first,
#you can merge the data sets as usual by:
spCourses = merge(spCourses, spFurtherEdu2, by = c("ID_t", "course"), all.x = TRUE)

#compare the two versions of the variables by:
addmargins(table(spCourses$wave.x, spCourses$wave.y))
addmargins(table(spCourses$tx80211.x, spCourses$tx80211.y))

#and then drop one of the versions by:
spCourses$wave.y = NULL
spCourses$tx80211.y = NULL
'-----'

'-----'

'** B) merge to spFurtherEdu1'

```

```
'** open spFurtherEdu1 and FurtherEdu2 datafiles'
spFurtherEdu1 = read.dta13("SC5_spFurtherEdu1_D_version.dta", convert.factors = T)
spFurtherEdu2 = read.dta13("SC5_spFurtherEdu2_D_version.dta", convert.factors = T)

'** merge spFurtherEdu2 using ID_t and courses'

intersect(names(spFurtherEdu1), names(spFurtherEdu2))
#common variables in the both data sets are "ID_t", "wave", "course" and "tx80211"
#Since the variables "wave" and "tx80211" are not part of the merging process,
#both versions are contained in the new merged data set
#as wave.x/wave.y and tx80211.x/tx80211.y.

'**To avoid that, there are two merging options:'
#1. You can include the variables in the merging process by:
spFurtherEdu1 =
    merge(spFurtherEdu1, spFurtherEdu2,
          by = c("ID_t", "course", "wave", "tx80211"), all.x = TRUE)
#In that case the versions from the master data set are kept (wave.x and tx80211.x)

#OR

#2. If you'd like to compare the both versions first,
#you can merge the data sets as usual by:
spFurtherEdu1 =
    merge(spFurtherEdu1, spFurtherEdu2,
          by = c("ID_t", "course"), all.x = TRUE)

#compare the two versions of the variables by:
addmargins(table(spFurtherEdu1$wave.x, spFurtherEdu1$wave.y))
addmargins(table(spFurtherEdu1$tx80211.x, spFurtherEdu1$tx80211.y))

#and then drop one of the versions by:
spFurtherEdu1$wave.y = NULL
spFurtherEdu1$tx80211.y = NULL
'-----'
```

### Example 53 (R): Working with spGap

```
'** open the data file'
spGap = read.dta13("SC5_spGap_D_version.dta", convert.factors = T)

'** only keep full or harmonized episodes'
spGap = subset(spGap, spGap$subspell == 0)

'** open the Biography data file'
Biography = read.dta13("SC5_Biography_D_version.dta", convert.factors = T)

'** merge the spGap to Biography'

#After merging, Stata merge has one variable more than R, because in Stata
#a merge indicator is produced during the merging process and in R isn't.
#Since we need a merge indicator here, the merge command has to be extended:
```

```
Biography = transform(merge(
  x = cbind(Biography, source = "master"),
  #x contains the Biography data set plus one extra column "source",
  #where source = "master"
  y = cbind(spGap, source = "using"),
  #y contains the spGap data set plus one extra column "source",
  #where source = "using"
  all.x = TRUE, by = c("ID_t", "splink")),
  #merges x and y by ID_t and splink
  source = ifelse(!is.na(source.x) & !is.na(source.y), "both",
    #in the merged dataset, source = "both" if the observations is in x AND in y
    ifelse(!is.na(source.x), "master", "using")),
    #otherwise, source = "master" if the obs. is only in x
    #and source = "using" if the obs. is only in y
  source.x = NULL,
  source.y = NULL
  #the columns "source" in x and y are deleted
)

#Since the variables wave and spms are in both data sets, Biography AND spGap
#check before merging by: intersect(names(Biography), names(spGap))
#and since the variables are not part of the merging process,
#both versions are contained in the new merged data set as
#wave.x/wave.y and spms.x/spms.y
#For each variable, one of the versions can be dropped by:
Biography$wave.y = NULL
Biography$spms.y = NULL

'** you now have an enhanced version of Biography, enriched by
** information from the spell module. The number of total episodes
** (i.e., the amount of rows in the Biography file) did not change.
** Verify this by tabulating the spell type by the merging variable
** generated during the merge process.'
addmargins(table(Biography$sptype, Biography$source))
```

### Example 54 (R): Working with spInternship

```
'** open the data file'
spInternship = read.dta13("SC5_spInternship_D_version.dta", convert.factors = T)

'** only keep full or harmonized episodes'
spInternship = subset(spInternship, spInternship$subspell == 0)

'** open the Biography data file'
Biography = read.dta13("SC5_Biography_D_version.dta", convert.factors = T)

'** merge spInternship to Biography'
#After merging, Stata merge has one variable more than R, because in Stata
#a merge indicator is produced during the merging process and in R isn't.
#Since we need a merge indicator here, the merge command has to be extended:
Biography = transform(merge(
  x = cbind(Biography, source = "master"),
  #x contains the Biography data set plus one extra column "source",
```

```
#where source = "master"
y = cbind(spInternship,source = "using"),
#y contains the spInternship data set plus one extra column "source",
#where source = "using"
all.x = TRUE, by = c("ID_t", "splink")),
#merges x and y by ID_t and splink
source = ifelse(!is.na(source.x) & !is.na(source.y), "both",
#in the merged dataset, source = "both" if the observations is in x AND in y
ifelse(!is.na(source.x), "master", "using")),
#otherwise, source = "master" if the obs. is only in x
#and source = "using" if the obs. is only in y
source.x = NULL,
source.y = NULL
#the columns "source" in x and y are deleted
)

#Since the variables wave and spms are in both data sets, Biography AND spInternship
#check before merging by: intersect(names(Biography), names(spInternship))
#and since the variables are not part of the merging process,
#both versions are contained in the new merged data set as
#wave.x/wave.y and spms.x/spms.y
#For each variable, one of the versions can be dropped by:
Biography$wave.y = NULL
Biography$spms.y = NULL

'** you now have an enhanced version of Biography, enriched by
** information from the spell module. The number of total episodes
** (i.e., the amount of rows in the Biography file) did not change.
** Verify this by tabulating the spell type by the merging variable
** generated during the merge process.'
addmargins(table(Biography$sptype, Biography$source))
```

### Example 55 (R): Working with spMilitary

```
'** open the data file'
spMilitary = read.dta13("SC5_spMilitary_D_version.dta", convert.factors = T)

'** only keep full or harmonized episodes'
spMilitary = subset(spMilitary, spMilitary$subspell == 0)

'** open the Biography data file'
Biography = read.dta13("SC5_Biography_D_version.dta", convert.factors = T)

'** merge spMilitary to Biography'

#After merging, Stata merge has one variable more than R, because in Stata
#a merge indicator is produced during the merging process and in R isn't.
#Since we need a merge indicator here, the merge command has to be extended:
Biography = transform(merge(
  x = cbind(Biography,source = "master"),
  #x contains the Biography data set plus one extra column "source",
  #where source = "master"
  y = cbind(spMilitary,source = "using"),
```

```

#y contains the spMilitary data set plus one extra column "source",
#where source = "using"
all.x = TRUE, by = c("ID_t", "splink")),
#merges x and y by ID_t and splink
source = ifelse(!is.na(source.x) & !is.na(source.y), "both",
#in the merged dataset, source = "both" if the observations is in x AND in y
            ifelse(!is.na(source.x), "master", "using")),
            #otherwise, source = "master" if the obs. is only in x
            #and source = "using" if the obs. is only in y
source.x = NULL,
source.y = NULL
#the columns "source" in x and y are deleted
)

#Since the variables wave and spms are in both data sets, Biography AND spMilitary
#check before merging by: intersect(names(Biography), names(spMilitary))
#and since the variables are not part of the merging process,
#both versions are contained in the new merged data set as
#wave.x/wave.y and spms.x/spms.y
#For each variable, one of the versions can be dropped by:
Biography$wave.y = NULL
Biography$spms.y = NULL

'** you now have an enhanced version of Biography, enriched by
** information from the spell module. The number of total episodes
** (i.e., the amount of rows in the Biography file) did not change.
** Verify this by tabulating the spell type by the merging variable
** generated during the merge process.'
addmargins(table(Biography$sptype, Biography$source))

```

### Example 56 (R): Working with spParLeave

```

'** open the data file'
spParLeave = read.dta13("SC5_spParLeave_D_version.dta", convert.factors = T)

'** only keep full or harmonized episodes'
spParLeave = subset(spParLeave, spParLeave$subspell == 0)

'** open the Biography data file'
Biography = read.dta13("SC5_Biography_D_version.dta", convert.factors = T)

'** merge spParLeave to Biography'

#After merging, Stata merge has one variable more than R, because in Stata
#a merge indicator is produced during the merging process and in R isn't.
#Since we need a merge indicator here, the merge command has to be extended:
Biography = transform(merge(
  x = cbind(Biography, source = "master"),
  #x contains the Biography data set plus one extra column "source",
  #where source = "master"
  y = cbind(spParLeave, source = "using"),
  #y contains the spParLeave data set plus one extra column "source",
  #where source = "using"

```

```

all.x = TRUE, by = c("ID_t", "splink")),
#merges x and y by ID_t and splink
source = ifelse(!is.na(source.x) & !is.na(source.y), "both",
#in the merged dataset, source = "both" if the observations is in x AND in y
             ifelse(!is.na(source.x), "master", "using")),
             #otherwise, source = "master" if the obs. is only in x
             #and source = "using" if the obs. is only in y
source.x = NULL,
source.y = NULL
#the columns "source" in x and y are deleted
)

#Since the variables wave and spms are in both data sets, Biography AND spParLeave
#check before merging by: intersect(names(Biography), names(spParLeave))
#and since the variables are not part of the merging process,
#both versions are contained in the new merged data set as
#wave.x/wave.y and spms.x/spms.y
#For each variable, one of the versions can be dropped by:
Biography$wave.y = NULL
Biography$spms.y = NULL

'** you now have an enhanced version of Biography, enriched by
** information from the spell module. The number of total episodes
** (i.e., the amount of rows in the Biography file) did not change.
** Verify this by tabulating the spell type by the merging variable
** generated during the merge process.'
addmargins(table(Biography$sptype, Biography$source))

```

### Example 57 (R): Working with spPartner

```

'** open the data file'
spPartner = read.dta13("SC5_spPartner_D_version.dta", convert.factors = T)

'** only keep full or harmonized episodes'
spPartner = subset(spPartner, spPartner$subspell == 0)

'** to find out if a respondent has ever been lived together with a partner,
** you could'
cbind(addmargins(table(spPartner$t733030)),
      addmargins(prop.table(table(spPartner$t733030))))

```

### Example 58 (R): Working with spSchool

```

'** open the data file'
spSchool = read.dta13("SC5_spSchool_D_version.dta", convert.factors = T)

'** only keep full or harmonized episodes'
spSchool = subset(spSchool, spSchool$subspell == 0)

'** open the Biography data file'
Biography = read.dta13("SC5_Biography_D_version.dta", convert.factors = T)

```

```
'** merge spSchool to Biography'

#After merging, Stata merge has one variable more than R, because in Stata
#a merge indicator is produced during the merging process and in R isn't.
#Since we need a merge indicator here, the merge command has to be extended:
Biography = transform(merge(
  x = cbind(Biography,source = "master"),
  #x contains the Biography data set plus one extra column "source",
  #where source = "master"
  y = cbind(spSchool,source = "using"),
  #y contains the spSchool data set plus one extra column "source",
  #where source = "using"
  all.x = TRUE, by = c("ID_t", "splink")),
  #merges x and y by ID_t and splink
  source = ifelse(!is.na(source.x) & !is.na(source.y), "both",
    #in the merged dataset, source = "both" if the observations is in x AND in y
    ifelse(!is.na(source.x), "master", "using")),
    #otherwise, source = "master" if the obs. is only in x
    #and source = "using" if the obs. is only in y
  source.x = NULL,
  source.y = NULL
  #the columns "source" in x and y are deleted
)

#Since the variables wave and spms are in both data sets, Biography AND spSchool
#check before merging by: intersect(names(Biography), names(spSchool))
#and since the variables are not part of the merging process,
#both versions are contained in the new merged data set as
#wave.x/wave.y and spms.x/spms.y
#For each variable, one of the versions can be dropped by:
Biography$wave.y = NULL
Biography$spms.y = NULL

'** you now have an enhanced version of Biography, enriched by
** information from the spell module. The number of total episodes
** (i.e., the amount of rows in the Biography file) did not change.
** Verify this by tabulating the spell type by the merging variable
** generated during the merge process.'
addmargins(table(Biography$sptype, Biography$source))
```

### Example 59 (R): Working with spSchoolExtExam

```
'** aim of this example is to evaluate the age of the respondent
** at the exam'

'** first, we have to get the birth date of the respondent'
#open pTargetCATI
pTargetCATI = read.dta13("SC5_pTargetCATI_D_version.dta", convert.factors = T)

#display value labels
levels(pTargetCATI$wave)

#keep only the first wave as this data is time-invariant
```

```

pTargetCATI =
  subset(pTargetCATI, pTargetCATI$wave == "2010/2011 (CATI+competencies)")

#keep only ID_t, t70000m and t70000y from pTarget
pTargetCATI =
  subset(pTargetCATI, select = c("ID_t", "t70000m", "t70000y"))

'** now, open the data file spSchoolExtExam'
spSchoolExtExam =
  read.dta13("SC5_spSchoolExtExam_D_version.dta", convert.factors = T)

'** merge the previously extracted birth dates in pTargetCATI to spSchoolExtExam'
spSchoolExtExam = merge(spSchoolExtExam, pTargetCATI, by = c("ID_t"), all.x = TRUE)

'** recode the two date variables (year, month) into one:'

Sys.setlocale("LC_TIME", "C")
#turns off the location-specific language, such that the english month names
#are recognized as months.

spSchoolExtExam$ts1130m = match(spSchoolExtExam$ts1130m, month.name)
spSchoolExtExam$t70000m = match(spSchoolExtExam$t70000m, month.name)
#transforms month names into month numbers

install.packages("zoo")
library(zoo)
#the zoo package is needed to transform time data

spSchoolExtExam$exam_date =
  as.yearmon(paste(spSchoolExtExam$ts1130y, spSchoolExtExam$ts1130m), "%Y %m")
spSchoolExtExam$birth_date =
  as.yearmon(paste(spSchoolExtExam$t70000y, spSchoolExtExam$t70000m), "%Y %m")
#recode the two date variables (year, month) into one

'** calculate the age (in years)'
spSchoolExtExam$age = (spSchoolExtExam$exam_date - spSchoolExtExam$birth_date)

'** show some deviation'
aggregate(spSchoolExtExam$age, by = list(spSchoolExtExam$ts11302),
  FUN = function(x)
    c(mean = mean(x, na.rm = TRUE),
      sd = sd(x, na.rm = TRUE), Freq = length(x)))
#displays mean and sd of age by school-leaving qualification

summary(spSchoolExtExam$age)
#display mean of age in general

sd(spSchoolExtExam$age, na.rm = TRUE)
#display sd of age in general

```



**Example 60 (R): Working with spSibling**

```

'** aim of this example is to evaluate the number of older and younger
** siblings of a respondent'

'** first, we have to get the birth date of the respondent'
#open pTargetCATI
pTargetCATI = read.dta13("SC5_pTargetCATI_D_version.dta", convert.factors = T)

#display value labels
levels(pTargetCATI$wave)

#keep only the first wave as this data is time-invariant
pTargetCATI =
  subset(pTargetCATI, pTargetCATI$wave == "2010/2011 (CATI+competencies)")

#keep only ID_t, t70000m and t70000y from pTarget
pTargetCATI = subset(pTargetCATI, select = c("ID_t", "t70000m", "t70000y"))

'** now, open the data file spSibling'
spSibling = read.dta13("SC5_spSibling_D_version.dta", convert.factors = T)

'** merge the previously extracted birth dates in pTargetCATI to spSibling'
spSibling = merge(spSibling, pTargetCATI, by = c("ID_t"), all.x = TRUE)

'** recode the two date variables (year, month) into one:'
Sys.setlocale("LC_TIME", "C")
#turns off the location-specific language, such that the english month names are
  recognized as months.

spSibling$tg3270m = match(spSibling$tg3270m, month.name)
spSibling$t70000m = match(spSibling$t70000m, month.name)
#transforms month names into month numbers

install.packages("zoo")
library(zoo)
#the zoo package is needed to transform time data

spSibling$sibling_bdate =
  as.yearmon(paste(spSibling$tg3270y, spSibling$tg3270m), "%Y %m")
spSibling$target_bdate =
  as.yearmon(paste(spSibling$t70000y, spSibling$t70000m), "%Y %m")
#recode the two date variables (year, month) into one

'** check the difference between the two'

spSibling$older = rep(NA, times = length(spSibling$ID_t))
#create an empty variable "older"

#check the difference between the two bdates:
for (i in 1:length(spSibling$older)) {

```

```

if(!is.na(spSibling$sibling_bdate[i]) & !is.na(spSibling$target_bdate[i]) &
    spSibling$sibling_bdate[i] > spSibling$target_bdate[i]) {
  spSibling$older[i] = 0
} else {
  if (!is.na(spSibling$sibling_bdate[i]) & !is.na(spSibling$target_bdate[i]) &
      spSibling$sibling_bdate[i] < spSibling$target_bdate[i]) {
    spSibling$older[i] = 1
  } else {
    spSibling$older[i] = NA
  }
}
}
}

'** generate the total amount of older siblings'
spSibling =
  within(spSibling, {total_older = ave(older, ID_t,
    FUN = function(x) sum(x, na.rm = TRUE))})

'** generate the total amount of younger siblings'
spSibling =
  within(spSibling, {total_younger = ave(older, ID_t,
    FUN = function(x) sum(1-x, na.rm = TRUE))})

'** aggregate to a single line for each respondent.
** the file then is cross-sectional with ID_t the sole identifier'

spSibling = subset(spSibling, select = c("ID_t", "total_older", "total_younger"))
#keep only the variables ID_t, total_older and total_younger

spSibling = unique(spSibling)
#drops duplicate rows from spSibling

```

### Example 61 (R): Working with spUnemp

```

'** open the data file'
spUnemp = read.dta13("SC5_spUnemp_D_version.dta", convert.factors = T)

'** only keep full or harmonized episodes'
spUnemp = subset(spUnemp, spUnemp$subspell == 0)

'** open the Biography data file'
Biography = read.dta13("SC5_Biography_D_version.dta", convert.factors = T)

'** merge spUnemp to Biography'

#After merging, Stata merge has one variable more than R, because in Stata
#a merge indicator is produced during the merging process and in R isn't.
#Since we need a merge indicator here, the merge command has to be extended:
Biography = transform(merge(
  x = cbind(Biography, source = "master"),
  #x contains the Biography data set plus one extra column "source",
  #where source = "master"

```

```

y = cbind(spUnemp, source = "using"),
#y contains the spUnemp data set plus one extra column "source",
#where source = "using"
all.x = TRUE, by = c("ID_t", "splink")),
#merges x and y by ID_t and splink
source = ifelse(!is.na(source.x) & !is.na(source.y), "both",
#in the merged dataset, source = "both" if the observations is in x AND in y
ifelse(!is.na(source.x), "master", "using")),
#otherwise, source = "master" if the obs. is only in x
#and source = "using" if the obs. is only in y
source.x = NULL,
source.y = NULL
#the columns "source" in x and y are deleted
)

#Since the variables wave and spms are in both data sets, Biography AND spUnemp
#check before merging by: intersect(names(Biography), names(spUnemp))
#and since the variables are not part of the merging process,
#both versions are contained in the new merged data set as
#wave.x/wave.y and spms.x/spms.y
#For each variable, one of the versions can be dropped by:
Biography$wave.y = NULL
Biography$spms.y = NULL

'** you now have an enhanced version of Biography, enriched by
** information from the spell module. The number of total episodes
** (i.e., the amount of rows in the Biography file) did not change.
** Verify this by tabulating the spell type by the merging variable
** generated during the merge process.'
addmargins(table(Biography$sptype, Biography$source))

```

### Example 62 (R): Working with spVocExtExam

```

'** aim of this example is to evaluate the age of the respondent
** at the exam'

'** first, we have to get the birth date of the respondent'
#open pTargetCATI
pTargetCATI = read.dta13("SC5_pTargetCATI_D_version.dta", convert.factors = T)

#display value labels
levels(pTargetCATI$wave)

#keep only the first wave as this data is time-invariant
pTargetCATI =
  subset(pTargetCATI, pTargetCATI$wave == "2010/2011 (CATI+competencies)")

#keep only ID_t, t70000m and t70000y from pTarget
pTargetCATI = subset(pTargetCATI, select = c("ID_t", "t70000m", "t70000y"))

'** open the data file spVocExtExam'
spVocExtExam = read.dta13("SC5_spVocExtExam_D_version.dta", convert.factors = T)

```

```
'** merge the previously extracted birth dates in pTargetCATI to spVocExtExam'
spVocExtExam = merge(spVocExtExam, pTargetCATI, by = c("ID_t"), all.x = TRUE)

'** recode the two date variables (year, month) into one:'
Sys.setlocale("LC_TIME", "C")
#turns off the location-specific language, such that the english month names are
  recognized as months.

spVocExtExam$ts1530m = match(spVocExtExam$ts1530m, month.name)
spVocExtExam$t70000m = match(spVocExtExam$t70000m, month.name)
#transforms month names into month numbers

install.packages("zoo")
library(zoo)
#the zoo package is needed to transform time data

spVocExtExam$exam_date =
  as.yearmon(paste(spVocExtExam$ts1530y, spVocExtExam$ts1530m), "%Y %m")
spVocExtExam$birth_date =
  as.yearmon(paste(spVocExtExam$t70000y, spVocExtExam$t70000m), "%Y %m")
#recode the two date variables (year, month) into one

'** calculate the age (in years)'
spVocExtExam$age = (spVocExtExam$exam_date - spVocExtExam$birth_date)

'** show some deviation'
aggregate(spVocExtExam$age, by = list(spVocExtExam$ts15304),
  FUN = function(x)
    c(mean = mean(x, na.rm = TRUE),
      sd = sd(x, na.rm = TRUE), Freq = length(x)))
#displays mean and sd of age by school-leaving qualification

summary(spVocExtExam$age)
#displays mean of age in general

sd(spVocExtExam$age, na.rm = TRUE)
#displays sd of age in general
```

### Example 63 (R): Working with spVocPrep

```
'** open the data file'
spVocPrep = read.dta13("SC5_spVocPrep_D_version.dta", convert.factors = T)

'** only keep full or harmonized episodes'
spVocPrep = subset(spVocPrep, spVocPrep$subspell == 0)

'** open the Biography data file'
Biography = read.dta13("SC5_Biography_D_version.dta", convert.factors = T)

'** merge spVocPrep to Biography'

#After merging, Stata merge has one variable more than R, because in Stata
#a merge indicator is produced during the merging process and in R isn't.
```

```
#Since we need a merge indicator here, the merge command has to be extended:
Biography = transform(merge(
  x = cbind(Biography,source = "master"),
  #x contains the Biography data set plus one extra column "source",
  #where source = "master"
  y = cbind(spVocPrep,source = "using"),
  #y contains the spVocPrep data set plus one extra column "source",
  #where source = "using"
  all.x = TRUE, by = c("ID_t", "splink")),
  #merges x and y by ID_t and splink
  source = ifelse(!is.na(source.x) & !is.na(source.y), "both",
  #in the merged dataset, source = "both" if the observations is in x AND in y
    ifelse(!is.na(source.x), "master", "using")),
    #otherwise, source = "master" if the obs. is only in x
    #and source = "using" if the obs. is only in y
  source.x = NULL,
  source.y = NULL
  #the columns "source" in x and y are deleted
)

#Since the variables wave and spms are in both data sets, Biography AND spVocPrep
#check before merging by: intersect(names(Biography), names(spVocPrep))
#and since the variables are not part of the merging process,
#both versions are contained in the new merged data set as
#wave.x/wave.y and spms.x/spms.y
#For each variable, one of the versions can be dropped by:
Biography$wave.y = NULL
Biography$spms.y = NULL

'** you now have an enhanced version of Biography, enriched by
** information from the spell module. The number of total episodes
** (i.e., the amount of rows in the Biography file) did not change.
** Verify this by tabulating the spell type by the merging variable
** generated during the merge process.'
addmargins(table(Biography$sptype, Biography$source))
```

### Example 64 (R): Working with spVocTrain

```
'** open the data file'
spVocTrain = read.dta13("SC5_spVocTrain_D_version.dta", convert.factors = T)

'** only keep full or harmonized episodes'
spVocTrain = subset(spVocTrain, spVocTrain$subspell == 0)

'** open the Biography data file'
Biography = read.dta13("SC5_Biography_D_version.dta", convert.factors = T)

'** merge spVocTrain to Biography'

#After merging, Stata merge has one variable more than R, because in Stata
#a merge indicator is produced during the merging process and in R isn't.
#Since we need a merge indicator here, the merge command has to be extended:
Biography = transform(merge(
```

```
x = cbind(Biography, source = "master"),
#x contains the Biography data set plus one extra column "source",
#where source = "master"
y = cbind(spVocTrain, source = "using"),
#y contains the spVocTrain data set plus one extra column "source",
#where source = "using"
all.x = TRUE, by = c("ID_t", "splink")),
#merges x and y by ID_t and splink
source = ifelse(!is.na(source.x) & !is.na(source.y), "both",
  ifelse(!is.na(source.x), "master", "using")),
#in the merged dataset, source = "both" if the observations is in x AND in y
#otherwise, source = "master" if the obs. is only in x
#and source = "using" if the obs. is only in y
source.x = NULL,
source.y = NULL
#the columns "source" in x and y are deleted
)

#Since the variables wave and spms are in both data sets, Biography AND spVocTrain
#check before merging by: intersect(names(Biography), names(spVocTrain))
#and since the variables are not part of the merging process,
#both versions are contained in the new merged data set as
#wave.x/wave.y and spms.x/spms.y
#For each variable, one of the versions can be dropped by:
Biography$wave.y = NULL
Biography$spms.y = NULL

'** you now have an enhanced version of Biography, enriched by
** information from the spell module. The number of total episodes
** (i.e., the amount of rows in the Biography file) did not change.
** Verify this by tabulating the spell type by the merging variable
** generated during the merge process.'
addmargins(table(Biography$sptype, Biography$source))
```

### Example 65 (R): Working with Weights

```
'** open the data file'
Weights = read.dta13("SC5_Weights_D_version.dta", convert.factors = T)

'** note that this file is cross-sectional,
**although the weights seem to contain panel logic'
attr(Weights, "var.labels")

'** only keep weights corresponding to all waves'
Weights = subset(Weights, select = c(ID_t, w_t123456789))

'** create a "panel" logic, i.e., clone each row'
Weights = Weights[rep(seq_len(nrow(Weights)), each = 9),]

'** then create a wave variable'
Weights$wave = ave(Weights$ID_t, Weights$ID_t, FUN = seq_along)

'** open CohortProfile'
```

```
CohortProfile = read.dta13("SC5_CohortProfile_D_version.dta", convert.factors = T)

#value labels of "wave" in "CohortProfile" and "Weights"
#have to be synchronized before merging
levels((CohortProfile$wave))
levels(Weights$wave)

Weights$wave = as.factor(Weights$wave)
#sets "wave" in "Weights" as factor

for (i in 1:9) {
  levels(Weights$wave)[i] = levels(CohortProfile$wave)[i]
  #assigns the same value labels to "wave" in "Weights" as in "CohortProfile"
}

'** and merges Weights to CohortProfile'
CohortProfile = merge(CohortProfile, Weights, by = c("ID_t", "wave"), all = TRUE)

'** note that this weight is only nonzero if respondents participated in all waves'
with(subset(CohortProfile, w_t123456789 != 0), addmargins(table(wave, tx80220)))
```

### Example 66 (R): Working with xInstitution

```
'** open datafile pTargetCATI'
pTargetCATI = read.dta13("SC5_pTargetCATI_D_version.dta", convert.factors = T)

'** copy the information from the first wave downwards for each target,
** unless a new value has been reported'
for (t in names(pTargetCATI[c("ID_i", "tg04001_g7")])) {
  #run over variables ID_i and tg04001_g7
  for (i in 2:length(pTargetCATI$ID_t)) {
    #run over all observations
    if(pTargetCATI$ID_t[i] == pTargetCATI$ID_t[i-1]){
      #for the same ID_t, check...
      if(is.na(pTargetCATI[[t]][i]) | pTargetCATI[[t]][i] == "Missing by design"){
        #...whether missing value or -54(Missing by design)
        pTargetCATI[[t]][i] = pTargetCATI[[t]][i-1]
        #copy information downwards, unless a new value has been reported
      }
    }
  }
}

'** drop all observations where no satisfaction with studies was reported'
levels(pTargetCATI$t514008)

#remove observations with NA in t514008
pTargetCATI = pTargetCATI[!(is.na(pTargetCATI$t514008)),]

#remove observations with other missings in t514008
pTargetCATI = subset(pTargetCATI, !(t514008 == "Don't know"
| t514008 == "Refused"
| t514008 == "Does not apply"))
```

```

| t514008 == "Missing by design"))

'** some respondents reported satisfaction with studies in 7th and in 9th waves
** to keep the latest information, create a seq and a max variables'
pTargetCATI = within(pTargetCATI,{seq = ave(ID_t, ID_t, FUN = seq_along)})
pTargetCATI = within(pTargetCATI,{max = ave(ID_t, ID_t, FUN = length)})

'** only keep the latest reported information'
pTargetCATI =
  subset(pTargetCATI, pTargetCATI$seq == pTargetCATI$max)

'** only keep the variables relevant for the merge and the analysis'
pTargetCATI =
  subset(pTargetCATI, select = c("ID_t", "ID_i", "tg04001_g7", "t514008"))

'** merge two variables from xInstitution'

#open datafile xInstitution
xInstitution = read.dta13("SC5_xInstitution_0_version.dta", convert.factors = T)

#merge xInstitution to pTargetCATI
pTargetCATI =
  merge(x = pTargetCATI,
        y = xInstitution[,c("ID_i", "g04001_g7", "tg92601_R", "tg92104_0")],
        by = c("ID_i", "g04001_g7"), all.x = TRUE)

'** assuming that the less students at university the more intensive the support by
** the university staff per student and the more satisfied are students with their
** studies tabulate Satisfaction with studies by Students 2010 total
** note that the following analysis is feasible in both, RemoteNEPS and Onsite'
cbind(addmargins(table(pTargetCATI$t514008, pTargetCATI$tg92601_R)))
cbind(addmargins(prop.table(table(pTargetCATI$t514008, pTargetCATI$tg92601_R))))

'** assuming that students at excellence universities are more satisfied with
** their studies, tabulate the distribution of satisfaction by tg92104_0
** note that the following analysis is only feasible in the Onsite version of SUF,
** since the variable tg92104_0 is anonymized in RemoteNEPS'
cbind(addmargins(table(pTargetCATI$t514008, pTargetCATI$tg92104_0)))
cbind(addmargins(prop.table(table(pTargetCATI$t514008, pTargetCATI$tg92104_0))))

```

### Example 67 (R): Working with xTargetCompetencies

```

'** open the data file xTargetCompetencies'
xTargetCompetencies =
  read.dta13("SC5_xTargetCompetencies_D_version.dta", convert.factors = T)

'** as the x in the filename indicates, this is a cross sectional file
** (no wave structure). You can verify this by asking if one row is
** solely identified by the respondents ID'
anyDuplicated(xTargetCompetencies[,c("ID_t")])
#returns "0" if there are no duplicates.

```



```
#If there are duplicates this command returns the index of the first duplicate

'** note that competence testing has been conducted in multiple waves
** an indicator marks if a row contains information for a specific wave'
table(xTargetCompetencies$wave_w1)
table(xTargetCompetencies$wave_w5)
table(xTargetCompetencies$wave_w7)

'** to work with competence data, you might want to merge it to CohortProfile.
** if you want to keep the panel logic (and not only add all competencies
** to every wave), you need a mergeable wave variable in xTargetCompetencies.
** here, we focus on math competencies, that have been tested in wave 1.'
xTargetCompetencies$wave =
  rep(levels(CohortProfile$wave)[1],length(xTargetCompetencies$ID_t))
xTargetCompetencies$wave = as.factor(xTargetCompetencies$wave)

'** now, keep cases which did took part in the testing'
xTargetCompetencies = subset(xTargetCompetencies, wave_w1 == "ja")

'** and reduce the dataset to the relevant variables'
xTargetCompetencies =
  subset(xTargetCompetencies, select = c(ID_t, wave, mas1_sc1, mas1_sc2))

'** and merge this to CohortProfile'

#open the data file Cohort Profile
CohortProfile = read.dta13("SC5_CohortProfile_D_version.dta", convert.factors = T)

#look for common variables in both data sets
intersect(names(CohortProfile), names(xTargetCompetencies))

#merge CohortProfile with xTargetCompetencies
CohortProfile =
  merge(CohortProfile, xTargetCompetencies, by = c("ID_t", "wave"), all = TRUE)
```

### A.2 Release notes

The following is the release note taken from the documentation page at the time this document has been computed:

```
=====
**
** NEPS STARTING COHORT 5 – RELEASE NOTES a.k.a CHANGE LOG
** changes and updates for release NEPS SC5 13.0.0
** (doi:10.5157/NEPS:SC5:13.0.0)
**
=====

* Known Issues *

MethodsCAWI:
    - waves 6 and 8 erroneously added to MethodsCAWI, data in those lines is
      completely missing, please drop these waves

=====
* Changes introduced to NEPS:SC5 by version 13.0.0 *
=====

General remarks:
    - some versionized variables were dropped, some were introduced ...
    - some intro variables are back again ...
    - some variable labels were corrected
    - supplemental meta information on several variables was added

CohortProfile:
    - some checks on plausibility and smoothing on ID_i has been done

EditionsBackup:
    - new dataset since release 12-0-0: provides original data prior to coding and
      smoothing during the process of data preparation

pTargetCATI:
    - variable tg26390_g1 "Spell number with reference to transition questions (
      from spEmp)" was generated for merging information from spEmp to
      information on transitions into the labormarket in pTargetCATI

spVocTrain:
    - variable t724401 (grades of academic degrees) dropped – information is
      integrated into variable ts15265 (the variable concerning grades of
      vocational qualification)
    - variable ts15219_g1 dropped – information provided in variable ts15219_g1 is
      redundant to information provided by variable ts15219

=====
**
** NEPS STARTING COHORT 5 – RELEASE NOTES a.k.a CHANGE LOG
** changes and updates for release NEPS SC5 12.0.0
** (doi:10.5157/NEPS:SC5:12.0.0)
**
=====

* Known Issues *
```

=====  
\* Changes introduced to NEPS:SC5 by version 12.0.0 \*  
=====

General remarks:

- for several variables information of the respective \_v-variables was integrated into the variables without suffixes. The respective \_v-variables were dropped.
- intro-variables were dropped, except for intro-variables in spChild and spPartner.

pTargetCATI:

- variable tg24201\_g1, tg24202\_g2 and tg02001\_ha were generated to provide detailed information on teaching degrees gathered in wave 1. For further information, see Data Manual (5.4 Teacher Education Students and Teachers)
- variable tg12001\_g2 was generated to provide missing information on the desired subject of study for target persons who claim to study in their desired subject. Therefore it combines information from variable tg04001 and tg12003.

pTargetCAWI:

- for several variables open answers were (belatedly) coded.

spChild:

- variable ts33204\_g1 was generated to provide information on the status of the child. Therefore category "other child in household" was added.

spEmp:

- variable tg2608a "student or other occupation" has been renamed to ts23256 to match the corresponding variable's name in other starting cohorts.

spSchoolExtExam:

- additional information on external examinations from wave 1 and 3 was gathered from file spSchool.

spVocExtExam:

- additional information on external examinations from wave 1 and 3 was gathered from file spVocTrain.

spVocPrep:

- variable ts13101 was deleted by mistake. Please use information on the program type for wave 1 and 3 from earlier SUF releases.

spVocTrain:

- variables tg24162\_g1, tg24165\_g1 and tg24168\_g1 were generated to provide information on major or minor subjects for each subspell of an episode. For further information, see Data Manual (5.1 service variables).
- information on external examinations from waves 1 and 3 was removed and integrated in file spVocExtExam.
- variable ts15221\_g1 was edited to provide (the revised) information on the intended vocational qualification for all target persons and for all subspells of an episode. For further information, see Data Manual (5.1 service variables).
- variable tg01003\_ha was edited and now excludes administration and business academies.
- servicevariables with information on subject of studies (tg2417\*) were revised.

=====

\* Changes introduced to NEPS:SC5 by version 11.0.0 \*  
=====

General remarks:

- several variables surveyed have been renamed to \*\_v1 and \*\_v2 in prior releases;  
this has been improved by renaming some variables with suffix \_v1 to variable names without suffixes  
and some variables with suffix \_v2 to suffix \_v1;  
a detailed list and comparison of \_v1 variables can be found in the Data Manual (Appendix A.3).

CohortProfile:

- testy testm testd erroneously coded to –56 for testing data in wave 7 have now been coded with correct dates

pTargetCAWI:

- there have been changes during the field phase regarding interviewer instructions in variable "tg51001";  
the new indicator variable "Version\_tg51001" contains information about the version of the survey instrument

MethodsCAWI:

- a new data file including para data from the CAWI interviews has been added

=====

\* Changes introduced to NEPS:SC5 by version 10.0.0 \*  
=====

General remarks:

- several variables surveyed prior to wave 10 have been renamed to \*\_v1 and \*\_v2,  
as wording of question texts has changed in recent survey instruments

CohortProfile:

- testy testm testd erroneously had been coded to –56 even though tx80522==1; this has been fixed
- new indicator variable tx80121 has been introduced: subsample "students of economics"
- tx80921 has been revised

xEcoCAPI:

- new dataset featuring items from CAPI–shortquestionnaire , economics–competency –test and the corresponding methods data that has been administered to students of economics in wave 7; all of these data has been removed from pTargetCATI , xTargetcompetencies , and MethodsCompetencies , respectively , for this subsample

=====

\* Changes introduced to NEPS:SC5 by version 9.0.0 \*  
=====

pTargetCATI:

- ts15911 (highest degree obtained) was falsely programmed in wave 9. Therefore ts15911\_g1 was generated for all participants.

spVocTrain:

- original variables tg2416\* (subjects) were edited due to discrepancies between subspells. Subsequently, subjects are filled for the first explicit mention only. Missing information was labeled accordingly. Working with service variables is recommended.
- service variables tg2417\* (subjects) have been revised so that each subspell of a corresponding spell is now filled with the first information available, still variables tg24170\_g1-g5, tg24173\_g1-g5 and tg24176\_g1-g5 provide complete information for all study episodes.
- ts15221 (qualification sought) was falsely derived in some cases. Therefore, ts15221\_g1 was generated for the affected episodes

```
=====
* Changes introduced to NEPS:SC5 by version 8.0.0 *
=====
```

General remarks on harmonization of variables concerning subjects, type of university and type of vocational training program:

- harmonization of type of university-variable: tg01003\_g1(pTargetCATI) >> tg01003\_ha (spVocTrain, considering values of ts15201)
- harmonized service variables on subjects: tg24160\_g\*, tg24163\_g\*, tg24166\_g\* (spVocTrain) >> tg24170\_g\*, tg24173\_g\*, tg24176\_g\* in spVocTrain (considering values of tg04001\_g1-5, tg04004\_g1-5, tg04007\_g1-5 in pTargetCATI)
- harmonization provides valid values for type of university and subjects where information on study episode from winter term 2010/11 was missing
- missing codes -28, -29 were introduced in the original variables tg24160\_g\*, tg24163\_g\*, tg24166\_g\*, tg01003\_g1, ts15201

CohortProfile:

- tx80951 indicates the participation status for students of economics in wave 7. Besides CATI survey and competency testing, these students had also the possibility of taking parting in a short CAPI questionnaire as well.

pTargetCATI:

- the concept of reflecting migrational background in NEPS SUFs has been improved in order to also represent migrants in 3.75th generation; thus, the older variables on migrational background [t400500\_g1,t400500\_g2, t400500\_g3] in the pTargetCATI dataset have been renamed using the "v1" suffix [t400500\_g1v1,t400500\_g2v1,t400500\_g3v1], and the new ones have been introduced
- variables of students of economics who took part in a short CAPI questionnaire were added to pTargetCATI

spVocTrain:

- service variables tg2417\* (subjects) and tg01003\_ha (type of university)\* were introduced to simplify working with the dataset. Small discrepancies from the original variables (tg2416\*) cannot be ruled out and have to be considered by the user.
- each subspell of a corresponding spell was filled with the most recent information available, so that the variables tg24170\_g1-5, tg24173\_g1-5, tg24176\_g1-5 provide complete information for all study episodes.

```
=====
* Changes introduced to NEPS:SC5 by version 6.0.0 *
=====
```

### General:

- starting with this release, all NEPS Scientific Use Files will ship with an additional, unicode-enabled Stata data set version;  
this version is only readable in Stata version 14 or younger, and is placed in the subdirectory "Stata14"
- translation for all meta data (variable and value labels, question texts, etc) have been revised and completed
- meta data for all variables have been revised and updated where appropriate
- additional waves 5 (CAWI) and 6 (CATI/CAPI) have been incorporated into the data
- the subspell harmonization routine in all spell datasets ("sp\*") has been updated, leading to more accurate harmonized subspell information (subspell==0) for panel continuation spells
- staff from NEPS stage 7 at the DZHW excessively reviewed and overworked all syntax for generated tg\*-variables, which may lead to slightly different contents
- staff from NEPS stage 7 at the DZHW reviewed the cohorts' sample frame in consultation with NEPS methods department, leading to 3 observations removed from the SUF
- all datasets from version 4.0.0 did not reflect the correct doi in their dataset labels; the correct doi would have been "10.5157/NEPS:SC5:4.0.0", not "none";  
this issue has been fixed and all datasets of version 6.0.0 correctly are labeled with doi:10.5157/NEPS:SC5:6.0.0

### xTargetCompetencies:

- all variables of domains "maths" and "reading" erroneously contained the missing value -54 ("missing by design") in versions 4.0.0 and 3.1.0;  
as there were no additional competency assessments in wave 4, it was safe to use the xTargetCompetencies dataset file from version 3.0.0  
instead without missing any information; this has been fixed

### pTargetCATI:

- variables "Specialized fair/congress: professional/personal reasons" [t272802\_w1] and "Specialized fair/congress: Learned something new" [t272802\_w1]  
as well as the corresponding variables for "Lectures" [t272802\_w2, t272802\_w2] and "Self-instruction programs" [t272802\_w3, t272802\_w3] in version 4.0.0 and earlier  
erroneously are not filled for all interviewees reporting the specific further education activity; this has been fixed
- variable names of variables "Father's mother: Country of birth" [t405240\*] and "Mother's father: Country of birth" [t405230\*] in dataset pTargetCATI erroneously had been flipped in version 4.0.0, also leading to slight inconsistencies in generated variables for migrational background; this has been fixed

### spChild:

- all wide variables documenting cohabitation (\*\_w\*) in version 4.0.0 and earlier with the focal child have been extracted and are now saved in the separate dataset "spChildCohab"

### spChildCohab:

- new dataset containing child cohabitation spells that formerly had been saved in wide format inside of spChild

spEmp:

- version 4.0.0 and earlier did not contain coded occupational information for studentical employment episodes reported in wave 1; this has been fixed

Biography:

- additional spells of type "data edition gap" have been inserted to fill gaps between
  - (a) the eighth birth day and the first reported episode and
  - (b) the most recently reported episode and the most recent interview date

```
=====
* Changes introduced to NEPS:SC5 by version 4.0.0 *
=====
```

General:

- full translations have been added
- wave 4 (online survey in semester 5) has been added
- several minor bug fixes to data edition scripts have been introduced

pTargetCATI:

- when generating variable "Global self-esteem" [t66003a\_g1] in the pTargetCATI dataset, variable "Global self-esteem: competence" [t66003d] erroneously had been ignored; this has been fixed; t66003a\_g1 can be re-generated in 3.1.0 using the following Stata syntax:

```
* -----BEGIN Stata-----
local target_variable t66003a_g1
nepsmis t66003a t66003b t66003c t66003d t66003e t66003f t66003g
t66003h t66003i t66003j
tempvar t66003b_r t66003e_r t66003f_r t66003h_r t66003i_r rowmissings
recode t66003b (1=5) (2=4) (3=3) (4=2) (5=1), generate('t66003b_r')
recode t66003e (1=5) (2=4) (3=3) (4=2) (5=1), generate('t66003e_r')
recode t66003f (1=5) (2=4) (3=3) (4=2) (5=1), generate('t66003f_r')
recode t66003h (1=5) (2=4) (3=3) (4=2) (5=1), generate('t66003h_r')
recode t66003i (1=5) (2=4) (3=3) (4=2) (5=1), generate('t66003i_r')
egen 'rowmissings'=rowmiss(t66003a 't66003b_r' t66003c t66003d ///
't66003e_r' 't66003f_r' t66003g 't66003h_r' 't66003i_r' t66003j)
egen 'target_variable'=rowtotal(t66003a 't66003b_r' t66003c t66003d ///
't66003e_r' 't66003f_r' t66003g 't66003h_r' 't66003i_r' t66003j) if '
rowmissings'==0 & wave==3
replace 'target_variable'=-54 if wave!=3
label variable 'target_variable' "Global self-esteem"
replace 'target_variable'=-55 if missing('target_variable')
* -----END Stata-----
```

xTargetCAWI:

- as wave 3 data makes this a panel dataset, the filename has changed from "xTargetCAWI" to "pTargetCAWI"

```
=====
* Changes introduced to NEPS:SC5 by version 3.1.0 *
=====
```

### General:

- meta data in all datasets have been revised and updated where appropriate
- English translation for all datasets except xTargetCAWI have been introduced to the data
- end dates in episodes neglected in the panel interview erroneously contained the interview date of the panel wave instead of the first interview's date; this has been fixed
- 185 duplicate respondents have been identified by the survey institute; the redundant observations have been dropped from the data, resulting in slightly smaller number of cases

### pTargetCATI:

- variables indicating migrational background (t400500\_g1 through \_g3) have been added

### spVocTrain:

- spell integration and recommendation (via variable tx20100) was erroneous; this has been fixed
- spell linkage between waves 1 and 3 was erroneous; this has been fixed

### spEmp:

- spell linkage between waves 1 and 3 was erroneous; this has been fixed

### Weights:

- dataset containing weighting variables has been added

### Basics:

- dataset containing oversimplified, "flat" cross-sectional data on the cohort has been added; use for orientation, not for analyses!

### xInstitution:

- dataset containing detailed information on the targets' institutions has been added for onsite access in Bamberg



## A.3 Comparison of \_v1 variables

The following tables shows all changes of variables where construction of a \_v1-variable seemed necessary. Note that by v1, we generally mean *first version* or *version one*. Thus, this usually is the old variant of a variable, which has been updated in a later wave. Small arrows indicate if an entry belongs to the old version («) or if it is an update (»). Grayed out entries did not change between the versions, and are printed for your orientation only.

## pTargetCATI

	«	t516201_v1		pTargetCATI		t516201	»
Label	«	Party election					
	»	Parliamentary elections: Party election					
Text		If parliamentary elections were to be held tomorrow, which party would you give your second vote to?					
-98		Don't know					
-97		Refused					
-93	»	Does not apply					
-55	«	Not determinable					
-54		Missing by design					
-21	»	Would not vote					
-20		Not entitled to vote, because no German citizenship					
1		CDU or CSU					
2		SPD - Social Democratic Party of Germany					
3	«	FDP (political party)					
	»	FDP - Free Democratic Party					
4		Bündnis 90/Die Grünen [green political party]					
5		Die Linke - Left Party					
6		NPD - National Democratic Party of Germany					
7	«	Die Republikaner - The Republicans					
8		other party					
9		Would not vote					
10		Piratenpartei - Pirate Party					
11	»	AfD					

	« <b>t525008_v1</b>   <b>pTargetCATI</b>   <b>t525008</b> »
Label	« smoking status » Smoking status
Text	« Did you smoke in the past or do you currently smoke? » Do you currently smoke - even if only occasionally?
-98	« Don't know
-97	« Refused
-54	Missing by design
1	« never smoked » yes, daily
2	« did smoke before » yes, occasionally
3	« currently smoke occasionally » no, not anymore
4	« currently smoke every day » have never smoked

	« <b>t525209_v1</b>   <b>pTargetCATI</b>   <b>t525209</b> »
Label	« Consumption of alcohol » Alcohol consumption Frequency Last 12 months
Text	« How often do you consume alcoholic drinks? » How often do you drink alcoholic beverages? Think about the average of the last 12 months.
-98	« Don't know
-97	« Refused
-54	Missing by design
1	« (almost) never » never
2	once a month or less
3	twice or three times a month
4	once a week
5	several times a week
6	« (almost) every day » daily

	« <b>tg2450a_v1</b>   <b>pTargetCATI</b>   <b>tg2450a</b> »	
Label	«	Doctorate context - Research project higher education institution
	»	Doctorate context - Third-party funded position higher education institution
Text	«	[MF] We have noted that you have begun a doctorate. In the following, we would like to ask you a few questions about this doctorate. Under what circumstances are you currently studying for a doctorate?
	»	We have noted that you have started to do your doctorate. In the following we would like to ask you some questions. A doctoral thesis can be written in very different institutional contexts, e.g. at a higher education institution or a research institution as a research assistant, in a structured doctoral program or as a freelance doctoral student without institutional involvement. In which institutional context are you currently doing your doctorate?
-98		Don't know
-97		Refused
-93	»	Does not apply
-92	«	Question erroneously not asked
-54		Missing by design
-52	«	Implausible value removed
-20	»	none of this
0		not specified
1		specified

	« <b>tg2450b_v1</b>   <b>pTargetCATI</b>   <b>tg2450b</b> »
Label	« Doctorate context - Chair higher education institution
	» Doctorate context - Budget funded position higher education institution
Text	« [MF] We have noted that you have begun a doctorate. In the following, we would like to ask you a few questions about this doctorate. Under what circumstances are you currently studying for a doctorate?
	» We have noted that you have started to do your doctorate. In the following we would like to ask you some questions. A doctoral thesis can be written in very different institutional contexts, e.g. at a higher education institution or a research institution as a research assistant, in a structured doctoral program or as a freelance doctoral student without institutional involvement. In which institutional context are you currently doing your doctorate?
-98	Don't know
-97	Refused
-93	» Does not apply
-92	« Question erroneously not asked
-54	Missing by design
-52	Implausible value removed
-20	» none of this
0	not specified
1	specified

	« <b>tg2450c_v1</b>   <b>pTargetCATI</b>   <b>tg2450c</b> »
Label	Doctorate context - Non-university research institution
Text	« [MF] We have noted that you have begun a doctorate. In the following, we would like to ask you a few questions about this doctorate. Under what circumstances are you currently studying for a doctorate?
	» We have noted that you have started to do your doctorate. In the following we would like to ask you some questions. A doctoral thesis can be written in very different institutional contexts, e.g. at a higher education institution or a research institution as a research assistant, in a structured doctoral program or as a freelance doctoral student without institutional involvement. In which institutional context are you currently doing your doctorate?
-98	Don't know
-97	Refused
-93	» Does not apply
-92	« Question erroneously not asked
-54	Missing by design
-52	Implausible value removed
-20	» none of this
0	not specified
1	specified

	« <b>tg2450d_v1</b>   <b>pTargetCATI</b>   <b>tg2450d</b> »
Label	Doctorate context - Doctoral program
Text	« [MF] We have noted that you have begun a doctorate. In the following, we would like to ask you a few questions about this doctorate. Under what circumstances are you currently studying for a doctorate? » We have noted that you have started to do your doctorate. In the following we would like to ask you some questions. A doctoral thesis can be written in very different institutional contexts, e.g. at a higher education institution or a research institution as a research assistant, in a structured doctoral program or as a freelance doctoral student without institutional involvement. In which institutional context are you currently doing your doctorate?
-98	Don't know
-97	Refused
-93	» Does not apply
-92	« Question erroneously not asked
-54	Missing by design
-52	Implausible value removed
-20	» none of this
0	not specified
1	specified

	« <b>tg2450e_v1</b>   <b>pTargetCATI</b>   <b>tg2450e</b> »
Label	« Doctorate context - Doctorate course of study » Doctorate context - Scholarship program
Text	« [MF] We have noted that you have begun a doctorate. In the following, we would like to ask you a few questions about this doctorate. Under what circumstances are you currently studying for a doctorate? » We have noted that you have started to do your doctorate. In the following we would like to ask you some questions. A doctoral thesis can be written in very different institutional contexts, e.g. at a higher education institution or a research institution as a research assistant, in a structured doctoral program or as a freelance doctoral student without institutional involvement. In which institutional context are you currently doing your doctorate?
-98	Don't know
-97	Refused
-93	» Does not apply
-92	« Question erroneously not asked
-54	Missing by design
-52	Implausible value removed
-20	» none of this
0	not specified
1	specified

	« <b>tg2450f_v1</b>   <b>pTargetCATI</b>   <b>tg2450f</b> »
Label	« Doctorate context - Private sector/industry
	» Doctorate context - Private sector (industrial research and development)
Text	« [MF] We have noted that you have begun a doctorate. In the following, we would like to ask you a few questions about this doctorate. Under what circumstances are you currently studying for a doctorate?
	» We have noted that you have started to do your doctorate. In the following we would like to ask you some questions. A doctoral thesis can be written in very different institutional contexts, e.g. at a higher education institution or a research institution as a research assistant, in a structured doctoral program or as a freelance doctoral student without institutional involvement. In which institutional context are you currently doing your doctorate?
-98	Don't know
-97	Refused
-93	» Does not apply
-92	« Question erroneously not asked
-54	Missing by design
-52	Implausible value removed
-20	» none of this
0	not specified
1	specified

	« <b>tg2450g_v1</b>   <b>pTargetCATI</b>   <b>tg2450g</b> »
Label	Doctorate context - Alongside studies
Text	« [MF] We have noted that you have begun a doctorate. In the following, we would like to ask you a few questions about this doctorate. Under what circumstances are you currently studying for a doctorate?
	» We have noted that you have started to do your doctorate. In the following we would like to ask you some questions. A doctoral thesis can be written in very different institutional contexts, e.g. at a higher education institution or a research institution as a research assistant, in a structured doctoral program or as a freelance doctoral student without institutional involvement. In which institutional context are you currently doing your doctorate?
-98	Don't know
-97	Refused
-93	» Does not apply
-92	« Question erroneously not asked
-54	Missing by design
-52	Implausible value removed
-20	» none of this
0	not specified
1	specified

	« tg2450h_v1   pTargetCATI   tg2450h »	
Label	«	Doctorate context - Without institutional integration
	»	Doctorate context - Without institutional integration, free doctorate student
Text	«	[MF] We have noted that you have begun a doctorate. In the following, we would like to ask you a few questions about this doctorate. Under what circumstances are you currently studying for a doctorate?
	»	We have noted that you have started to do your doctorate. In the following we would like to ask you some questions. A doctoral thesis can be written in very different institutional contexts, e.g. at a higher education institution or a research institution as a research assistant, in a structured doctoral program or as a freelance doctoral student without institutional involvement. In which institutional context are you currently doing your doctorate?
-98		Don't know
-97		Refused
-93	»	Does not apply
-92	«	Question erroneously not asked
-54		Missing by design
-52		Implausible value removed
-20	»	none of this
0		not specified
1		specified

« tg2450i_v1   pTargetCATI   tg2450i »		
Label		Doctorate context - other
Text	«	[MF] We have noted that you have begun a doctorate. In the following, we would like to ask you a few questions about this doctorate. Under what circumstances are you currently studying for a doctorate?
	»	We have noted that you have started to do your doctorate. In the following we would like to ask you some questions. A doctoral thesis can be written in very different institutional contexts, e.g. at a higher education institution or a research institution as a research assistant, in a structured doctoral program or as a freelance doctoral student without institutional involvement. In which institutional context are you currently doing your doctorate?
-98		Don't know
-97		Refused
-93	»	Does not apply
-92	«	Question erroneously not asked
-54		Missing by design
-52		Implausible value removed
-20	»	none of this
0		not specified
1		specified

« tg60013_v1   pTargetCATI   tg60013 »		
Label		Auxiliary variable: phase of teacher education and employment (CATI)
Text		[AUTO] Auxiliary variable: Teaching groups, current status
-54		Missing by design
0		no teaching reference or status unknown
1		first phase teacher training not yet completed
2		completed teaching degree course and Referendariat is intended or completed teaching degree course and employment as a teacher is intended
3		ongoing Referendariat
4		completed Referendariat and employment as a teacher is intended
5		Employment as teacher

« tg60031_v1   pTargetCATI   tg60031 »		
Label		Preload Completed teaching degree course
Text		[AUTO] Preload Completed teaching degree course
-54		Missing by design
0		no teaching degree course completed
1		teaching degree course completed



	« <b>ts15911_v1</b>   <b>pTargetCATI</b>   <b>ts15911</b> »
Label	« Graduate
	» Auxiliary variable: Highest degree
Text	« [AUX]
	» [AUX] Highest degree
-54	Missing by design
0	« no higher education qualification
	» No degree
1	« BA, MA, Diploma, Staatsexamen
	» BA
2	« Doctorate
	» MA, Diploma, Staatsexamen
3	» Doctorate

### pTargetCAWI

	« <b>t289900_v1</b>   <b>pTargetCAWI</b>   <b>t289900</b> »
Label	Type of accommodation
Text	Now we would like to ask you a few questions about your living situation and your spending. During term time, do you stay primarily...
-99	« Filtered
-97	Refused
-92	Question erroneously not asked
-91	Survey aborted
-54	Missing by design
1	« with parents or relatives?
	» with your parents?
2	in a dormitory?
3	« in some other rental accommodation?
	» in another type of rented apartment?/in a rented apartment?
4	« in an apartment/house that you own?
	» in a condo/own house?
5	with private individuals for subtenancy?
6	» with relatives?

« tg51101_v1   pTargetCAWI   tg51101 »	
Label	Curr. activity: Employed
Text	« [MF] Which of the following activities are your currently doing? I am currently ...
	» [MF] Which of the following positions do you currently work in? I am currently ...
-99	Filtered
-97	Refused
-92	Question erroneously not asked
-91	Survey aborted
-54	Missing by design
0	not specified
1	specified

---

« tg51102_v1   pTargetCAWI   tg51102 »	
Label	Curr. activity: Volontariat
Text	« [MF] Which of the following activities are your currently doing? I am currently ...
	» [MF] Which of the following positions do you currently work in? I am currently ...
-99	Filtered
-97	Refused
-92	Question erroneously not asked
-91	Survey aborted
-54	Missing by design
0	not specified
1	specified

---

« tg51103_v1   pTargetCAWI   tg51103 »	
Label	Curr. activity: Internship
Text	« [MF] Which of the following activities are your currently doing? I am currently ...
	» [MF] Which of the following positions do you currently work in? I am currently ...
-99	Filtered
-97	Refused
-92	Question erroneously not asked
-91	Survey aborted
-54	Missing by design
0	not specified
1	specified

	« <b>tg51104_v1</b>   <b>pTargetCAWI</b>   <b>tg51104</b> »
Label	« Curr. activity: Vocational training
	» Voc. train./further educ.: vocational training
Text	« [MF] Which of the following activities are you currently doing? I am currently ...
	» Are you currently ...?
-99	Filtered
-97	Refused
-92	Question erroneously not asked
-91	Survey aborted
-54	Missing by design
-21	» none of it
0	not specified
1	specified

	« <b>tg51108_v1</b>   <b>pTargetCAWI</b>   <b>tg51108</b> »
Label	« Curr. activity: Retraining or further education
	» Voc. train./further educ.: retraining, further education
Text	« [MF] Which of the following activities are you currently doing? I am currently ...
	» Are you currently ...?
-99	Filtered
-97	Refused
-92	Question erroneously not asked
-91	Survey aborted
-54	Missing by design
-21	» none of it
0	not specified
1	specified

	« <b>tg51109_v1</b>   <b>pTargetCAWI</b>   <b>tg51109</b> »	
Label	«	Curr. activity: (Voluntary) services, (military/alternative/community/social)
	»	Other activities: Voluntary services, (military, social)
Text	«	[MF] Which of the following activities are your currently doing? I am currently ...
	»	Are you also or exclusively doing any of the following activities? I am currently ...
-99		Filtered
-97		Refused
-92		Question erroneously not asked
-91		Survey aborted
-54		Missing by design
0		not specified
1		specified

	« <b>tg51110_v1</b>   <b>pTargetCAWI</b>   <b>tg51110</b> »	
Label	«	Curr. activity: On parental leave
	»	Other activities: Parental leave
Text	«	[MF] Which of the following activities are your currently doing? I am currently ...
	»	Are you also or exclusively doing any of the following activities? I am currently ...
-99		Filtered
-97		Refused
-92		Question erroneously not asked
-91		Survey aborted
-54		Missing by design
0		not specified
1		specified

	« <b>tg51111_v1</b>   <b>pTargetCAWI</b>   <b>tg51111</b> »
Label	« Curr. activity: Housewife/househusband
	» Other activities: Housewife/househusband
Text	« [MF] Which of the following activities are your currently doing? I am currently ...
	» Are you also or exclusively doing any of the following activities? I am currently ...
-99	Filtered
-97	Refused
-92	Question erroneously not asked
-91	Survey aborted
-54	Missing by design
0	not specified
1	specified

	« <b>tg51112_v1</b>   <b>pTargetCAWI</b>   <b>tg51112</b> »
Label	« Curr. activity: Unemployed
	» Other activities: Unemployed
Text	« [MF] Which of the following activities are your currently doing? I am currently ...
	» Are you also or exclusively doing any of the following activities? I am currently ...
-99	Filtered
-97	Refused
-92	Question erroneously not asked
-91	Survey aborted
-54	Missing by design
0	not specified
1	specified

« tg51113_v1   pTargetCAWI   tg51113 »	
Label	« Curr. activity: On sick leave
	» Other activities: On sick leave
Text	« [MF] Which of the following activities are your currently doing? I am currently ...
	» Are you also or exclusively doing any of the following activities? I am currently ...
-99	Filtered
-97	Refused
-92	Question erroneously not asked
-91	Survey aborted
-54	Missing by design
0	not specified
1	specified

---

« tg51114_v1   pTargetCAWI   tg51114 »	
Label	« Curr. activity: Other
	» Other activities: Other, namely:
Text	« [MF] Which of the following activities are your currently doing? I am currently ...
	» Are you also or exclusively doing any of the following activities? I am currently ...
-99	Filtered
-97	Refused
-92	Question erroneously not asked
-91	Survey aborted
-54	Missing by design
0	not specified
1	specified

---

« tg51115_v1   pTargetCAWI   tg51115 »	
Label	Curr. activity: Referendariat
Text	« [MF] Which of the following activities are your currently doing? I am currently ...
	» [MF] Which of the following positions do you currently work in? I am currently ...
-99	Filtered
-97	Refused
-92	Question erroneously not asked
-91	Survey aborted
-54	Missing by design
0	not specified
1	specified

« tg51116_v1   pTargetCAWI   tg51116 »	
Label	Curr. activity: Vicariate
Text	« [MF] Which of the following activities are your currently doing? I am currently ...
	» [MF] Which of the following positions do you currently work in? I am currently ...
-99	Filtered
-97	Refused
-92	Question erroneously not asked
-91	Survey aborted
-54	Missing by design
0	not specified
1	specified

---

« tg51117_v1   pTargetCAWI   tg51117 »	
Label	Curr. activity: Trainee program
Text	« [MF] Which of the following activities are your currently doing? I am currently ...
	» [MF] Which of the following positions do you currently work in? I am currently ...
-99	Filtered
-97	Refused
-92	Question erroneously not asked
-91	Survey aborted
-54	Missing by design
0	not specified
1	specified

---

« tg51118_v1   pTargetCAWI   tg51118 »	
Label	Curr. activity: Probationary year / practical year
Text	« [MF] Which of the following activities are your currently doing? I am currently ...
	» [MF] Which of the following positions do you currently work in? I am currently ...
-99	Filtered
-97	Refused
-92	Question erroneously not asked
-91	Survey aborted
-54	Missing by design
0	not specified
1	specified

	« <b>tg51300_v1</b>   <b>pTargetCAWI</b>   <b>tg51300</b> »
Label	« Change of field of study since starting university » Field of study changed since last survey
Text	« Have you changed your field of study since starting your studies in winter semester 2010/2011? » Have you changed your field of study since <h_zebePRE(Label)>?
-99	Filtered
-98	Don't know
-97	Refused
-92	Question erroneously not asked
-91	Survey aborted
-54	Missing by design
-52	Implausible value removed
1	yes
2	no

	« <b>tg51400_v1</b>   <b>pTargetCAWI</b>   <b>tg51400</b> »
Label	« Change in leaving qualification since starting university » Change Leaving qualification since last survey
Text	« Have you switched your chosen leaving qualification since the starting your studies in winter semester 2010/2011 (for example, from a bachelor's degree to a state examination)? » Have you changed the leaving qualification since <h_zebePRE(Label)> (for example, from a Bachelor degree to a Staatsexamen)?
-99	Filtered
-98	Don't know
-97	Refused
-92	Question erroneously not asked
-91	Survey aborted
-54	Missing by design
-52	Implausible value removed
1	yes
2	no



	« <b>tg51500_v1</b>   <b>pTargetCAWI</b>   <b>tg51500</b> »
Label	« Change in university after starting studies
	» Change of higher education institution since last survey
Text	« Have you changed universities since starting your studies in winter semester 2010/2011?
	» Have you changed higher education institution since <h_zebePRE(Label)>?
-99	Filtered
-98	Don't know
-97	Refused
-92	Question erroneously not asked
-91	Survey aborted
-54	Missing by design
-52	Implausible value removed
1	yes
2	no

spEmp

		« <b>ts23228_v1</b>   <b>spEmp</b>   <b>ts23228</b> »
Label		Type of education required
Text		What kind of training is usually required to do this job?
-98		Don't know
-97		Refused
-92	«	Question erroneously not asked
-54		Missing by design
1	«	no qualification
	»	No qualification
2	«	a training on the job
	»	Training on the job
3	«	a completed vocational training
	»	Completed vocational training
4	«	a completed training at a Fachschule
	»	Leaving certificate from a Fachschule
5	«	a master craftsman's/craftswoman's certificate or technician certificate
	»	a Master craftsman/craftswoman or technician certificate
6	«	a completed higher education degree (university of applied sciences or university)
7	«	a doctorate or habilitation
	»	A doctorate or habilitation
8	»	A Bachelor's degree (university of applied sciences or university)
9	»	A Master's degree or Staatsexamen, a diploma or a Magister (degrees from a university of applied sciences or university)

		« <b>ts23901_v1</b>   <b>spEmp</b>   <b>ts23901</b> »
Label		Auxiliary variable: Current employment
Text		[AUX] Auxiliary variable Current employment
-95	«	Implausible value
-54		Missing by design
1	«	currently employed
	»	Current employment
2	«	employed within the last year, but not currently
	»	completed employment
3	«	not employed within the last year / end not assignable

	« <b>ts23911_v1</b>   <b>spEmp</b>   <b>ts23911</b> »
Label	« Auxiliary variable: Type of employment
	» Auxiliary variable: Employee type
Text	« [AUX] Beschäftigtentyp
	» [AUX] Employee type
-54	Missing by design
-29	» Value from the last sub-episode
-20	« Not assignable
1	« Worker/ employee
	» worker/employee/civil servant/soldier/not classifiable
2	« Civil servants/soldiers
	» temporary/seasonal worker
3	» 2nd job market/training opportunities
4	» self-employed/assistant/ freelancer
5	« 2nd job market
6	« Freelancers
7	« Self-employed persons
8	« Positions in an assisting capacity
9	« Vocational training jobs
13	» semi-skilled or unskilled work/student assistant
14	» private student tuition/homework supervision

### spInternship

	« <b>tg36111_v1</b>   <b>spInternship</b>   <b>tg36111</b> »
Label	Average working hours Internship
Text	How many hours per week are your average working hours in this internship?
-98	Don't know
-97	» Refused
-54	Missing by design
-21	No fixed working hours
-20	« more than 50 hours per week
	» more than 90 hours per week

### spPartner

	« <b>ts31223_v1</b>   <b>spPartner</b>   <b>ts31223</b> »
Label	Employment Partner
Text	« Is your partner currently full-time employed, part-time employed or unemployed? » Is your partner currently employed full-time or part-time, has a side-job or is unemployed?
-98	Don't know
-97	Refused
-54	Missing by design
1	« primarily working » full-time employed
2	» part-time employed
3	« part-time employed » in a side job
4	unemployed

	« <b>ts31510_v1</b>   <b>spPartner</b>   <b>ts31510</b> »
Label	« Termination of partnership (separation/death, moving out without separation) » End of the partnership due to separation from or death of the partner
Text	« Have you divorced, separated or is your partner deceased? » Did you get divorced, did you split up, or did your partner die?
-98	» Don't know
-97	Refused
-54	Missing by design
1	« Divorced/civil partnership annulled » divorced / civil partnership annulled
2	« Separated » separated
3	« Partner deceased » partner deceased
4	» marital status unchanged
5	» moved back together, currently living together
6	» living apart, but still in partnership
9	« Do not live together any more, with partnership still persisting

spVocExtExam

« ts15304_v1   spVocExtExam   ts15304 »		
Label		External examination qualification
Text		What leaving qualification did you obtain?
-99	«	Filtered
-98	«	Don't know
-55	«	Not determinable
-20		no qualification
1		completed vocational training (administrative, company-based, industrial, agricultural), journeyman's/journeywoman's certificate, dual vocational training, GDR: skilled worker's certificate)
2		Leaving qualification from a school for healthcare professionals
3		Leaving certificate of Berufsfachschule or commercial school
4		other type of leaving certificate from a Fachschule
5		Master craftsman's/craftswoman's diploma
6		Technician's qualification
7	«	Diploma
8	«	Bachelor
9	«	Master
10		Diploma from a university of applied sciences (Dipl(FH))
11		Diploma from a university
12		Bachelor (in teaching)
13		Bachelor (not in teaching)
14		Master (in teaching)
15		Master (not in teaching)
16		Magister
17		First Staatsexamen (in teaching)
18		First Staatsexamen (not in teaching)
19	«	Second/Third Staatsexamen
	»	Second/Third Staatsexamen (not in teaching)
20		Doctorate
21		Habilitation
22		Medical specialist
23		Civil service examination for the subclerical class
24		Civil service examination for the clerical class
25		Civil service examination for the executive class
26		Civil service examination for the administrative class
27		IHK (Chamber of Industry and Commerce) examination
28		other qualification
29		other degree from a higher education institution (e.g., ecclesiastical examination, artistic examination)
30	»	Second Staatsexamen (in teaching)

### spVocTrain

		« <b>tg24205_v1</b>   <b>spVocTrain</b>   <b>tg24205</b> »	
Label		Point of time decision for master	
Text		When did you make the decision for your Master's degree program?	
-54		Missing by design	
1		before starting the previous higher education program	
2		during the previous higher education program	
3		after completion of the previous course of study	
		« <b>th32367_v1</b>   <b>spPartner</b>   <b>th32367</b> »	
Label		Episode update 7	
Text	«	Has your (male) partner obtained a (another) vocational qualification since our last interview?	
	»	Has your partner achieved a (additional) vocational qualification since our last interview?	
-98		Don't know	
-97		Refused	
-54		Missing by design	
1		Not acquired (any further) qualification	
2		Acquired (another) qualification	
		« <b>th32368_v1</b>   <b>spPartner</b>   <b>th32368</b> »	
Label		Episode update 8	
Text	«	In our last interview in <20101P3(intmPRE/intjPRE)> we noted that your (male) partner was working as a <28102P11> at that time.	
	»	In our last interview in <20101P3(intmPRE/intjPRE)> we noted that your partner was working as a <28102P11> at that time.	
-54		Missing by design	
1		TP does NOT disagree	
2		TP disagrees	

	« <b>ts31204_v1</b>   <b>spPartner</b>   <b>ts31204</b> »
Label	« Partner: born in Germany/abroad » Partner: born Germany/Abroad
Text	And where was he born?
-98	Don't know
-97	Refused
-54	Missing by design
1	In Germany / in the area that is present-day Germany
2	In Germany's former eastern territories
3	Abroad / in another country

	« <b>ts31206_v1</b>   <b>spPartner</b>   <b>ts31206</b> »
Label	« Partner's age on moving to Germany » Age at the time of moving to Germany partner
Text	« At what age did your (male) partner move to Germany (for the first time)? » At what age did your partner move to Germany (for the first time)?
-98	Don't know
-97	Refused
-93	Does not apply
-54	Missing by design
-20	Partner never moved to Germany



	« <b>ts3120y_v1</b>   <b>spPartner</b>   <b>ts3120y</b> »
Label	« Partner's year of birth » Year of birth partner
Text	« In what year was your (male) partner <28109> born? » In what year was your partner <28109> born?
-99	Filtered
-98	Don't know
-97	Refused
-96	Not in list
-95	Implausible value
-94	Not reached
-93	Does not apply
-92	Question erroneously not asked
-91	Survey aborted
-90	Unspecific missing
-56	Not participated
-55	Not determinable
-54	Missing by design
-53	Anonymized
-52	Implausible value removed
-51	No estimate in check module

	« <b>ts31211_v1</b>   <b>spPartner</b>   <b>ts31211</b> »
Label	Partner German
Text	Does your partner <28109> have German citizenship?
-98	Don't know
-97	Refused
-54	Missing by design
1	Yes
2	No

	« ts31212_v1   spPartner   ts31212 »	
Label	«	Highest general school-leaving qualification of partner
	»	highest general school-leaving qualification partner
Text	«	What is your (male) partner's highest general school-leaving qualification?
	»	What is your partner's highest general school-leaving qualification?
-98		Don't know
-97		Refused
-95	«	Implausible value
-54		Missing by design
-29	»	Value from last-mentioned sub-episode
-20	«	no school-leaving qualification
	»	No school-leaving qualification
1	«	Basic leaving certificate of the Hauptschule [school for basic secondary education], Volksschule [former name for compulsory school], 8th grade Polytechnische Oberschule (POS) [type of school in the former GDR offering intermediate secondary education]
2	«	Qualifying leaving certificate of the Hauptschule
	»	Qualifying Hauptschulabschluss
3	«	Certificate of intermediate secondary education (Realschule [intermediate secondary school], Wirtschaftsschule [type of school in Bavaria offering intermediate secondary education with a focus on commerce], entrance qualification for universities of a
	»	Certificate of intermediate secondary education (Real-/Wirtschaftsschulabschluss; Fachschul-/Fachoberschulreife; 10. grade POS)
4	«	Entrance qualification for universities of applied sciences, leaving certificate of the Fachoberschule
	»	Fachhochschulreife/completion Fachoberschule
5	«	General / subject-specific higher education entrance qualification (Abitur [higher education entrance qualification]/12th grade extended Oberschule [type of school in the former GDR leading to university entrance qualification])
	»	general/subject-specific university entrance qualification (Abitur/EOS 12. grade)
6		Leaving certificate of a special needs school
7		Other qualification

	« ts31214_v1   spPartner   ts31214 »	
Label	«	Partner: highest professional qualification
	»	highest vocational qualification partner
Text	«	What is your (male) partner's highest vocational qualification?
	»	What is your partner's highest vocational qualification?
-98		Don't know
-97		Refused
-93	«	Does not apply
-55	«	Not determinable
-54		Missing by design
-29	«	Value from the last sub-episode
	»	Value from last-mentioned sub-episode
-20	«	no vocational qualification
	»	No vocational qualification
1		Completed apprenticeship (commercial, corporate, trade-oriented, agricultural) journey person's or assistant's certificate, dual vocational education and training, GDR: skilled worker's certificate
2		Master, technician's certificate
3		Civil service vocational training (civil service examination)
4		Leaving certificate from a school for health care professionals
5		Leaving certificate of Berufsfachschule, leaving certificate of a commercial school
6	«	Leaving qualification of the Fachschule (also leaving qualification of Fachakademie [type of school in Bavaria offering advanced vocational education and the possibility to obtain the entrance qualification for universities of applied sciences])
	»	Leaving certificate of the Fachschule [school for continuing vocational training] (also leaving certificate of the Fachakademie [school for advanced vocational education and the entrance qualification for universities of applied sciences in Bavaria])
7		Leaving certificate from a Fachschule in the former GDR
8		Bachelor (e.g. B.A., B.Sc.)
9		Diplom, Master (M.A.)
10		Magister, state examination
11	«	Doctorate, habilitation [post-doctoral lecturing qualification]
	»	Doctorate, habilitation
12	«	Berufsakademie, dual university without further details
	»	Berufsakademie without further specific information
13		College of public administration without further specification
14	«	University of applied sciences
	»	University of applied sciences, former college of engineering without further details
15		University without further details
16		Higher education degree (degree course) without further specification
17		Semi-skilled vocational training with a company
19		GDR: Qualification as a semi-skilled worker
21		Other vocational qualification

	« <b>ts31219_v1</b>   <b>spPartner</b>   <b>ts31219</b> »
Label	« Institution awarding higher education qualification to partner » Institution awarding higher education degree partner
Text	« And at which educational institution did your partner acquire this leaving certificate? Was that a Berufsakademie or dual university, a college of public administration, Fachhochschule or a university? » And at which educational institution did your partner obtained this qualification? Was that a Berufsakademie or a cooperative state university, a college of public administration, a university of applied sciences or a university?
-98	Don't know
-97	Refused
-54	Missing by design
1	Berufsakademie, dual university
2	College of public administration
3	University of applied sciences
4	University (including technical university, medical university, theological college, teacher training college, veterinary college as well as colleges of music and art)
5	Other institution

	« <b>ts31221_v1</b>   <b>spPartner</b>   <b>ts31221</b> »
Label	« Doctorate partner » Doctorate Partner
Text	« Was your (male) partner awarded a doctorate or is he currently working towards his doctorate? » Has your partner completed his doctorate degree or is he currently doing a doctorate?
-98	Don't know
-54	Missing by design
1	Yes, doctorate completed
2	Yes, currently doing doctorate / did doctorate back then
3	No

	« <b>ts31223_v1</b>   <b>spPartner</b>   <b>ts31223</b> »
Label	« Employment Partner » Employment partner
Text	« Is your partner currently employed full or part-time, working 'on the side' or not employed? » Is your partner currently employed full or part-time, has a side-job or is unemployed?
-98	Don't know
-97	Refused
-54	Missing by design
1	Full-time employed
2	Part-time employed
3	Side-job
4	Unemployed

	« <b>ts31224_v1</b>   <b>spPartner</b>   <b>ts31224</b> »
Label	« Working hours, partner » Working time partner
Text	« How many hours does your (male) partner work on average per week – including any side jobs? » How many hours does your partner on average work per week – including possible side-jobs?
-98	Don't know
-97	Refused
-54	Missing by design
-21	« no fix working hours » No fixed working hours
-20	« more than 90 hours per week » More than 90 hours per week

	« <b>ts31225_v1</b>   <b>spPartner</b>   <b>ts31225</b> »
Label	« Non-employment, partner » Unemployment Partner
Text	« What does your partner currently do predominantly? » What does your partner currently mainly do?
-98	Don't know
-97	Refused
-54	Missing by design
1	Unemployed
2	Short-time working
3	One-euro-job, job creation scheme, or similar program offered by the Federal Employment Agency/Job Center or ARGE
4	Partial retirement irrespective of what phase
5	General school education
6	Vocational training
7	Vocational training for Master, technician's certificate
8	Higher education
9	Doctorate
10	Vocational retraining, advanced or further education
11	On maternity leave/parental leave
12	Housewife/househusband
13	Sick / temporarily unable to work
14	Retiree, pensioner, (preliminary) retirement
15	(Voluntary) military/community service, Federal Volunteers Service, alternative service or voluntary social/ecological year or European Voluntary Service
16	Other

	« ts31227_v1   spPartner   ts31227 »	
Label	«	Professional position, partner
	»	Professional position partner
Text	«	What is your (male) partner's current professional position?
	»	What is your partner's current occupational status?
-98		Don't know
-97		Refused
-54		Missing by design
1	«	Worker
	»	Employee
2		Employee, also employee of the public service
3	«	Civil servant, including judges
	»	Civil servant, also judge
4	«	Regular / professional soldier
	»	Regular or professional soldier
5		Self-employed person
6	«	Assisting family member
	»	assisting family member
7	«	Freelancer
	»	freelancer

	« <b>ts31228_v1</b>   <b>spPartner</b>   <b>ts31228</b> »
Label	« Exact professional position partner » Exact vocational position partner
Text	« And what is your (male) partner's exact professional position there? » And what is your partner's exact occupational status there?
-98	Don't know
-97	Refused
-54	Missing by design
10	Unskilled worker
11	Semi-skilled worker/partially skilled worker
12	Skilled worker, journeyman [trained craftsman]
13	Assistant foreman, group leader, Brigadier [former GDR: Leader of a work unit]
14	Master, construction foreman
20	Low-skill occupation, e.g. salesperson
21	Qualified occupation, e.g. office clerk, technical draftsman
22	Highly qualified occupation or leading position, e.g. engineer, research assistant, department manager
23	Occupation involving extensive management duties e.g., director, CEO, member of the executive board
24	Production or plant foreman
30	In sub-clerical class (up to and including 'Oberamtsmeister')
31	In the clerical class, from assistant to principal secretary or office inspector, inclusively
32	Executive class (from inspector to Amtsrat inclusive and/or Oberamtsrat as well as elementary, secondary or intermediate school teacher inclusive)
33	In the administrative class, including judge, e.g. teacher starting from level of Studienrat [junior position held by school teachers upon career entry], senior government official
40	Military team rank
41	Non-commissioned officer, e.g. staff sergeant, sergeant, master sergeant
42	Simple officer to captain (included)
43	Staff officers from major to general/admiral
51	Self-employed as an academic, self-employed professional, e.g. physician, lawyer, architect
52	Self-employed person in agriculture
53	Self-employed person in trade, commerce, industry, service; other self-employment or entrepreneurship



	« <b>ts31230_v1</b>   <b>spPartner</b>   <b>ts31230</b> »
Label	Management position partner
Text	« Does your partner have a leading position in his activity? » Does your partner hold a management position?
-98	Don't know
-97	Refused
-54	Missing by design
1	Yes
2	No

	« <b>ts31410_v1</b>   <b>spPartner</b>   <b>ts31410</b> »
Label	« Marriage / registered civil partnership » Marriage/ registered civil partnership
Text	« Did you marry your partner (<28109>)? » Have you married your partner or have you registered the civil partnership?
-98	Don't know
-97	Refused
-54	Missing by design
1	Yes
2	No

	« <b>ts3141m_v1</b>   <b>spPartner</b>   <b>ts3141m</b> »
Label	« Date of marriage (month)
	» Marriage date (month)
Text	« When did you marry your partner <28109>?
	» When did you marry or register your civil partnership?
-98	Don't know
-97	Refused
-93	Does not apply
-54	Missing by design
1	January
2	February
3	March
4	April
5	May
6	June
7	July
8	August
9	September
10	October
11	November
12	December
21	Beginning of the year/winter
24	Spring/Easter
27	Mid-Year/Summer
30	Fall
32	End of year

	« <b>ts3141y_v1</b>   <b>spPartner</b>   <b>ts3141y</b> »
Label	« Date of marriage (year)
	» Marriage date (year)
Text	« When did you marry your partner <28109>?
	» When did you marry or register your civil partnership?
-99	Filtered
-98	Don't know
-97	Refused
-96	Not in list
-95	Implausible value
-94	Not reached
-93	Does not apply
-92	Question erroneously not asked
-91	Survey aborted
-90	Unspecific missing
-56	Not participated
-55	Not determinable
-54	Missing by design
-53	Anonymized
-52	Implausible value removed
-51	No estimate in check module

	« <b>ts31510_v1</b>   <b>spPartner</b>   <b>ts31510</b> »
Label	End of the partnership due to separation or death of a partner
Text	Did you get divorced, did you separate or is your (male) partner deceased?
-98	Don't know
-97	Refused
-54	Missing by design
1	Divorced / civil partnership annulled
2	Separated
3	Partner deceased
4	Marital status unchanged
5	Moved back in with partner, currently living together
6	No longer living together but partnership still exists

	« <b>ts3151m_v1</b>   <b>spPartner</b>   <b>ts3151m</b> »
Label	« <b>Date of partner's death (month)</b> » <b>Date of death Partner (month)</b>
Text	When did your partner pass away?
-98	Don't know
-97	Refused
-93	Does not apply
-54	Missing by design
1	January
2	February
3	March
4	April
5	May
6	June
7	July
8	August
9	September
10	October
11	November
12	December
21	Beginning of the year/winter
24	Spring/Easter
27	Mid-Year/Summer
30	Fall
32	End of year

	« <b>ts3151y_v1</b>   <b>spPartner</b>   <b>ts3151y</b> »
Label	« Date of partner's death (year)
	» Date of death Partner (year)
Text	When did your partner pass away?
-99	Filtered
-98	Don't know
-97	Refused
-96	Not in list
-95	Implausible value
-94	Not reached
-93	Does not apply
-92	Question erroneously not asked
-91	Survey aborted
-90	Unspecific missing
-56	Not participated
-55	Not determinable
-54	Missing by design
-53	Anonymized
-52	Implausible value removed
-51	No estimate in check module

	« <b>ts3152m_v1</b>   <b>spPartner</b>   <b>ts3152m</b> »
Label	Date of moving apart (Month)
Text	« <b>When did you or your partner move out of the shared home?</b> » <b>When did you or your partner moved out of the common household?</b>
-98	Don't know
-97	Refused
-93	Does not apply
-54	Missing by design
1	January
2	February
3	March
4	April
5	May
6	June
7	July
8	August
9	September
10	October
11	November
12	December
21	Beginning of the year/winter
24	Spring/Easter
27	Mid-Year/Summer
30	Fall
32	End of year

	« <b>ts3152y_v1</b>   <b>spPartner</b>   <b>ts3152y</b> »	
Label	Date of moving apart (Year)	
Text	«	When did you or your partner move out of the shared home?
	»	When did you or your partner moved out of the common household?
-99	Filtered	
-98	Don't know	
-97	Refused	
-96	Not in list	
-95	Implausible value	
-94	Not reached	
-93	Does not apply	
-92	Question erroneously not asked	
-91	Survey aborted	
-90	Unspecific missing	
-56	Not participated	
-55	Not determinable	
-54	Missing by design	
-53	Anonymized	
-52	Implausible value removed	
-51	No estimate in check module	

### spVocExtExam

« ts15304_v1   spVocExtExam   ts15304 »		
Label		External examination qualification
Text		What leaving qualification did you obtain?
-20		no qualification
1		Completed apprenticeship (commercial, corporate, trade-oriented, agricultural) journey person's or assistant's certificate (skilled worker's certificate), dual vocational education and training
2		Leaving certificate from a school for health care professionals
3		Leaving certificate of Berufsfachschule, leaving certificate of a commercial school
4	«	Other type of leaving certificate of the Fachschule
	»	other type of leaving certificate from a Fachschule
5		Master's / foreman's certificate
6		Technician's certificate
10		Diplom from a university of applied sciences (Dipl(FH))
11		Diplom from a university
12		Bachelor's degree teaching profession
13		Bachelor (not for teaching post)
14		Master teaching post
15		Master (not for teaching post)
16	«	Magister
	»	Magister [German degree in tertiary education, pre-Bologna system, level equivalent to master]
17		First state examination for teaching post
18		First state examination (not for teaching post)
19	«	Second or third state examination
	»	Second/Third State Examination (without teaching post)
20		Doctorate
21		Habilitation
22		Medical specialist
23		Civil service examination for the subclerical class
24		Civil service examination for the clerical class
25		Civil service examination for the executive class
26		Civil service examination for the administrative class
27		IHK (Chamber of Industry and Commerce) examination
28		Other leaving qualification
29	«	Other degree from a higher education institute (e.g., ecclesiastical examination, artistic examination)
	»	Other degree from a higher education institution (e.g., ecclesiastical examination, artistic examination)
30	»	Second State Examination teaching post



### spVocTrain

« tg24146_v1   spVocTrain   tg24146 »	
Label	« Change of type of leaving qualification as against pre-episode »
	» Change of type of qualification compared with pre-episode
Text	« Will your next degree course result in the same leaving qualification as the degree course we talked about before, or is it another leaving qualification, e.g. Bachelor instead of state examination or elementary school teaching qualification instead of Gymnasium teaching qualification? »
	» Will your next degree course result in the same leaving qualification as the degree course we talked about before, or is it another leaving qualification, e.g. Master instead of Bachelor or elementary school teaching qualification instead of Gymnasium teaching qualification?
-99	« Filtered »
-98	« Don't know »
-97	« Refused »
-92	« Question erroneously not asked »
-54	Missing by design
-29	« Value from the last sub-episode »
	» Value from last-mentioned sub-episode
1	Same leaving qualification
2	Other qualification

« tg24205_v1   spVocTrain   tg24205 »	
Label	Point of time decision for master
Text	When did you make the decision for your master degree program?
-54	Missing by design
1	before starting the previous higher education program
2	During the previous higher education program
3	after ending the previous higher education program

« ts15219_v1   spVocTrain   ts15219 »	
Label	Vocational qualification
Text	« Which civil service examination did you take? »
	» Which civil service examinations did you do?
-99	« Filtered »
-98	« Don't know »
-92	« Question erroneously not asked »
-55	Not determinable
-54	Missing by design

(...)

-20	«	no qualification
	»	Without any qualification
1	«	Completed apprenticeship (commercial, corporate, trade-oriented, agricultural) journey person's or assistant's certificate (skilled worker's certificate), dual vocational education and training
	»	Completion of an apprenticeship (commercial, corporate, trade-oriented, agricultural), journeyman's or assistant's certificate (skilled worker's certificate), dual training
2		Leaving certificate from a school for health care professionals
3	«	Leaving certificate of Berufsfachschule, leaving certificate of a commercial school
	»	Leaving certificate of a Berufsfachschule, leaving certificate of a Handelsschule
4	«	Other type of leaving certificate of the Fachschule
	»	other type of leaving certificate from a Fachschule
5	«	Master's / foreman's certificate
6	«	Technician's certificate
	»	Technician's training certificate
7		Diplom
8	«	Bachelor
	»	Bachelor's degree
9	«	Master
	»	Master's degree
10	«	Diplom from a university of applied sciences (Dipl(FH))
	»	Diplom from a Fachhochschule (Dipl(FH))
11	«	Diplom from a university
	»	University Diplom
12		Bachelor's degree teaching profession
13	«	Bachelor (not for teaching post)
	»	Bachelor's degree (without teaching profession)
14	«	Master teaching post
	»	Master's degree teaching profession
15	«	Master (not for teaching post)
	»	Master's degree (without teaching profession)
16		Magister
17	«	First state examination for teaching post
	»	First state examination teaching profession
18	«	First state examination (not for teaching post)
	»	First state examination (without teaching)
19	«	Second state examination
	»	Second/Third state examination
20		Doctorate

(...)

21		Habilitation
22		Medical specialist
23		Civil service examination for the subclerical class
24		Civil service examination for the clerical class
25		Civil service examination for the executive class
26		Civil service examination for the administrative class
27		IHK (Chamber of Industry and Commerce) examination
28	«	Other leaving qualification
	»	other qualification
29		Other degree from a higher education institute (e.g., ecclesiastical examination, artistic examination)

« ts15221_v1   spVocTrain   ts15221 »		
Label	«	Aspired vocational education qualification (reconstructed)
	»	aspired vocational training qualification
Text	«	Which civil service examination [final exam for the different classes of German civil service careers] do you/did you want to do?
	»	Which civil service examinations do/did you want to do?
-98		Don't know
-97	«	Refused
-92		Question erroneously not asked
-55		Not determinable
-54		Missing by design
-20	«	no qualification
	»	No degree
1	«	Completed apprenticeship (commercial, corporate, trade-oriented, agricultural) journey person's or assistant's certificate (skilled worker's certificate), dual vocational education and training
	»	Completion of an apprenticeship (commercial, corporate, trade-oriented, agricultural), journeyman's or assistant's certificate (skilled worker's certificate), dual training
2		Leaving certificate from a school for health care professionals
3	«	Leaving certificate of Berufsfachschule, leaving certificate of a commercial school
	»	Leaving certificate of a Berufsfachschule, leaving certificate of a Handelsschule
4	«	Other type of leaving certificate of the Fachschule
	»	other type of leaving certificate from a Fachschule
5	«	Master's / foreman's certificate
6	«	Technician's certificate
	»	Technician's training certificate

(...)

7		Diplom
8	«	Bachelor
	»	Bachelor's degree
9	«	Master
	»	Master's degree
10	«	Diplom from a university of applied sciences (Dipl(FH))
	»	Diplom from a Fachhochschule (Dipl(FH))
11	«	Diplom from a university
	»	University Diplom
12		Bachelor's degree teaching profession
13	«	Bachelor (not for teaching post)
	»	Bachelor's degree (without teaching profession)
14	«	Master teaching post
	»	Master's degree teaching profession
15	«	Master (not for teaching post)
	»	Master's degree (without teaching profession)
16		Magister
17	«	First state examination for teaching post
	»	First state examination teaching profession
18	«	First state examination (not for teaching post)
	»	First state examination (without teaching)
19	«	Second state examination
	»	Second/Third state examination
20		Doctorate
21		Habilitation
22		Medical specialist
23		Civil service examination for the subclerical class
24		Civil service examination for the clerical class
25		Civil service examination for the executive class
26		Civil service examination for the administrative class
27		IHK (Chamber of Industry and Commerce) examination
28	«	Other leaving qualification
	»	other qualification
29		Other degree from a higher education institute (e.g., ecclesiastical examination, artistic examination)

	« tg2452m_v1   spVocTrain   tg2452m »	
Label	«	Start of the doctorate (month)
	»	Starting time of the doctorate (month)
Text	«	And when did you begin the content-related work on your doctorate?
	»	And when have you started with the content work for your doctorate?
-99	«	Filtered
-98		Don't know
-97		Refused
-96	«	Not in list
-95	«	Implausible value
-94	«	Not reached
-93		Does not apply
-92	«	Question erroneously not asked
-91	«	Survey aborted
-90	«	Unspecific missing
-56	«	Not participated
-55	«	Not determinable
-54		Missing by design
-53	«	Anonymized
-52	«	Implausible value removed
-51	«	No estimate in check module
1	»	January
2	»	February
3	»	March
4	»	April
5	»	May
6	»	June
7	»	July
8	»	August
9	»	September
10	»	October
11	»	November
12	»	December
21	»	Beginning of the year/winter
24	»	Spring/Easter
27	»	Mid-Year/Summer
30	»	Fall
32	»	End of year

	« <b>tg2452y_v1</b>   <b>spVocTrain</b>   <b>tg2452y</b> »
Label	« Start of the doctorate (year)
	» Starting time of the doctorate (year)
Text	« And when did you begin the content-related work on your doctorate?
	» And when have you started with the content work for your doctorate?
-99	Filtered
-98	Don't know
-97	Refused
-96	Not in list
-95	Implausible value
-94	Not reached
-93	Does not apply
-92	Question erroneously not asked
-91	Survey aborted
-90	Unspecific missing
-56	Not participated
-55	Not determinable
-54	Missing by design
-53	Anonymized
-52	Implausible value removed
-51	No estimate in check module