NEPS National Educational Panel Study

FDZ-LIfBi

Data Manual

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

Scientific Use File Version 16.0.0

LIFBI LEIBNIZ INSTITUTE FOR EDUCATIONAL TRAJECTORIES

Copyrighted Material Leibniz Institute for Educational Trajectories (LIfBi) Wilhelmsplatz 3, 96047 Bamberg Director: Prof. Dr. Cordula Artelt Administrative Director: Dr. Stefan Echinger Bamberg; May 16, 2022

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. (2022). Data Manual NEPS Starting Cohort 5– First-Year Students, From Higher Education to the Labor Market, Scientific Use File Version 16.0.0. 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 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 DZHW Hannover: Isabelle Fiedler, Stefanie Gäckle, Annika Grieb, Marie Kühn, Uta Liebeskind, Katrin Mergard, Andreas Ortenburger, Hilde Schaeper

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

Daniel Bela, Simon Dickopf, Hannes Götz, 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 Web: https://www.neps-data.de/datacenter Phone: +49 951 863 3511



1	Intro	duction		1
	1.1	About t	this manual	1
	1.2	Further	documentation	1
	1.3	Data re	lease strategy	3
	1.4	Data ac	Cess	5
	1.5	Publicat	tions with NEPS data	6
	1.6	Rules ar	nd recommendations	7
	1.7	On usin	ng the Federal State label (Bundeslandkennung)	9
	1.8	User se	rvices	9
	1.9	Contact	ting the Research Data Center	11
2	Samj	oling and	d Survey Overview	12
	2.1	From hi	igher education to the labor market	12
	2.2	Samplin	ng strategy	13
	2.3	Compet	tence 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
		2.4.15	Wave 15: 2019 (CATI)	33
		2.4.16	Wave 16: 2020 (CATI)	34
3	Gene	eral Conv	ventions	35
	3.1	File nan	nes	35
	3.2	Variable	es	37
		3.2.1	Conventions for general variable naming	37
		3.2.2	Conventions for competence variable naming	40
		3.2.3	Labels	43
	3.3	Missing	g values	44

	3.4	Generat	ted variables
4	Data	Structu	re 48
	4.1	Overvie	w
	4.2	Identifie	ers
	4.3	Panel d	ata
	4.4	Episode	or spell data
		4.4.1	Edition of the life course
		4.4.2	Revoked episodes
		4.4.3	Subspells and harmonization of episodes
	4.5	Data file	es
		4.5.1	Basics
		4.5.2	Biography
		4.5.3	CohortProfile
		4.5.4	EditionBackups
		4.5.5	Education
		4.5.6	MethodsCATI
		4.5.7	MethodsCAWI
		4.5.8	MethodsCompetencies
		4.5.9	pTargetCATI
		4.5.10	pTargetCAWI
		4.5.11	pTargetCORONA
		4.5.12	pTargetMicrom
		4.5.13	spChild
		4.5.14	spChildCohab
		4.5.15	, spCourses
		4.5.16	, spЕmp
		4.5.17	spFurtherEdu1
		4.5.18	spFurtherEdu2
		4.5.19	spGap
		4.5.20	spinternship
		4.5.21	spMilitary
		4.5.22	spParLeave
		4.5.23	spPartner 105
		4.5.24	spSchool 107
		4.5.25	spSchoolExtExam
		4.5.26	spSibling
		4.5.27	splinemp 113
		4 5 28	splocereaks 115
		4 5 29	spVocExtExam 117
		4 5 30	sprocezcezani
		4 5 31	splocinepine in
		4 5 3 2	StudyStates 121
		4.5.32 4.5.22	Weights 127
		T.J.JJ	ΨΟΒΠΟ · · · · · · · · · · · · · · · · · · ·

		4.5.34 xEcoCAPI	29
		4.5.35 xInstitution	31
		4.5.36 xPlausibleValues	33
		4.5.37 xTargetCompetencies	35
5	Spec	al Issues 1	37
	5.1	Special Types of Variables	37
		5.1.1 Service Variables	37
		5.1.2 Auxiliary variables	38
		5.1.3 Version variables	38
		5.1.4 Preload Variables	39
	5.2	Coding field of study	39
		5.2.1 Recruitment	39
		5.2.2 Panel Waves	40
	5.3	Coding of Higher Education Institutions	41
	5.4	Special features of interruption episodes in spVocTrain	42
	5.5	Teacher Education Students and Teachers 1	42
	5.6	Wave Specific Issues	45
A	Refe	rences 14	46
В	Арре	endix 14	48
	B.1	R examples	48
	B.2	Release notes	77
	B.3	Comparison of _v1 variables	87

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>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 B.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:
 - \rightarrow 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 Digitial Object Identifier.

Every release of a NEPS Scientific Use File is registered at da | ra 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:SC5:16.0.0. 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).

SUF Version	DOI	Date of release
16.0.0 (current)	doi:10.5157/NEPS:SC5:16.0.0	2022-05-16
15.0.0	doi:10.5157/NEPS:SC5:15.0.0	2021-05-19
14.1.0	doi:10.5157/NEPS:SC5:14.1.0	2020-12-02
14.0.0	doi:10.5157/NEPS:SC5:14.0.0	2020-05-27
13.0.0	doi:10.5157/NEPS:SC5:13.0.0	2020-02-14
12.0.0	doi:10.5157/NEPS:SC5:12.0.0	2019-07-26
11.0.0	doi:10.5157/NEPS:SC5:11.0.0	2018-09-06
10.0.0	doi:10.5157/NEPS:SC5:10.0.0	2018-04-19
9.0.0	doi:10.5157/NEPS:SC5:9.0.0	2017-06-23
8.0.0	doi:10.5157/NEPS:SC5:8.0.0	2016-12-23
6.0.0	doi:10.5157/NEPS:SC5:6.0.0	2016-03-31
4.0.0	doi:10.5157/NEPS:SC5:4.0.0	2014-09-30
3.1.0	doi:10.5157/NEPS:SC5:3.1.0	2014-05-16
3.0.0	doi:10.5157/NEPS:SC5:3.0.0	2013-07-05

Table 1: Release history of SUF in Starting Cohort 5

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>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:
 - \rightarrow 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 LIfBi = 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 > 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.ndicators 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>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 citing the data version (DOI) as follows:

NEPS Network. (2022). National Educational Panel Study, Scientific Use File of Starting Cohort First-Year Students. Leibniz Institute for Educational Trajectories (LIfBi), Bamberg. https://doi.org/10.5157/NEPS:SC5:16.0.0

In addition, the NEPS study is to be referred to at an appropriate place:

This paper uses data from the National Educational Panel Study (NEPS; see Blossfeld and Roßbach, 2019). The NEPS is carried out by the Leibniz Institute for Educational Trajectories (LIfBi, Germany) in cooperation with a nationwide network.

Finally, the reference article should be listed in the bibliography:

Blossfeld, H.-P., & Roßbach, H.-G. (Eds.). (2019). *Education as a lifelong process: The German National Educational Panel Study (NEPS). Edition ZfE* (2nd ed.). Springer VS. https://doi.org/10.1007/978-3-658-23162-0.

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 > 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. (2022). Data Manual NEPS Starting Cohort 5– First-Year Students, From Higher Education to the Labor Market, Scientific Use File Version 16.0.0. 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 Network instead:

NEPS Network. (2022). *Starting Cohort 5: First-Year Students (SC5), Wave 16, Questionnaires (SUF Version 16.0.0)*. 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.

Introduction

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 > 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

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:

 \rightarrow 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 ist 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:

 \rightarrow 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 co-hort 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>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 Web: → www.neps-data.de>Data Center>Contact Data Center Phone: +49 951 863 3511

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
- <u>not</u> attending higher education institutions run by Federal Ministries or Federal States for members of their public services (e.g., University of Applied Labour Studies/University of the German Federal Armed Forces Munich/Universität der Bundeswehr München)

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 German Centre of Higher Education Research and Science Studies (DZHW) in order to gain access to the students. The administration of 261 institutions declared their cooperativeness, 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 et al., 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>Data Center>Data and Documenation
>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.

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

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

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

² Reduced testing: In wave 7, the stage-specific competence test (ba) was realized in a subsample of students and graduates of business sciences only.
 ³ 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>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 16.



Figure 2: Survey progress of Starting Cohort 5 (waves 1 to 16)

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)

						20	11					
	01	02	03	04	05	06	07	08	09	10	11	12
interview target person									n=	=12,27	72	

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

Target persons

Sample Participants of the first wave willing to take part in the panel

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)

						20	12					
	01	02	03	04	05	06	07	08	09	10	11	12
interview target person										n=	=11,20	02

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)

						20	13					
	01	02	03	04	05	06	07	08	09	10	11	12
interview target person										n=	=10,18	82

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)





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 infas – Institute for Applied Social Sciences, Bonn

Mode of survey Computer-assisted telephone interview (CATI)

Target persons (Subsample B)

Current wave Selected 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 infas – Institute for Applied Social Sciences, Bonn

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

2.4.8 Wave 8: 2014 (CAWI)

						20	14					
	01	02	03	04	05	06	07	08	09	10	11	12
interview target person										n	=8,62	.8

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) 2015 01 02 03 04 05 06 07 08 09 10 11 12 n=10,096



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)

						20	16					
	01	02	03	04	05	06	07	08	09	10	11	12
interview target person				1	n=9	,089						

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)

		2016 02 03 04 05 06 07 08 09 10 11 12												
	01	02	03	04	05	06	07	08	09	10	11	12		
interview target person											(n=7)	,020		

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 infas – Institute for Applied Social Sciences, Bonn

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 (computerbased 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)

	2018											
	01	02	03	04	05	06	07	08	09	10	11	12
interview target person											n=5	,161

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)

2.4.15 Wave 15: 2019 (CATI)



Figure 17: Field times and realized case numbers in wave 15

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)

2.4.16 Wave 16: 2020 (CATI)

		2020										
	01	02	03	04	05	06	07	08	09	10	11	12
interview target person					n=6	,218						



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)

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 LIfBi 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 4.5.4 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 4.3 and section 4.4 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.

Element	Definition						
SC[1-6]	Indicator for the starting cohort						
	 1 = Newborns 2 = Kindergarten 3 = Fifth-grade students 4 = Ninth-grade students 5 = First-year university students 6 = Adults 						
[filename]	Meaning of the file name						
	<i>Prefix</i> : x = cross-sectional file; sp = spell file; p = panel file						
	<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)						
	File names of generated datasets do not have a prefix and always start with a capital letter (e.g., CohortProfile, Weights)						
[D,R,O]	Indicator for the confidentiality level						
	 D = Download version R = Remote access version O = On-site access version 						
[#]-[#]-[#](_beta)	Indicator for the release version						
	<i>First digit</i> : the main release number is incremented with every fur- ther 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						
	Second digit: 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						
	<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						
	_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						

Table 3: Naming conventions for NEPS file names

For instance, the file SC5_CohortProfile_D_16.0.0.dta refers to the *CohortProfile* data of *Starting Cohort 5* in its *Download* version of the Scientific Use File release *16.0.0*.

General Conventions

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 19. 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 19: 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	Respondent type									
	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	Topic/domain								
	Indicator to which theoretical dimension or educational stage the variable refers								
	1 = Competence development								
	2 = Learning environments								
	3 = Educational decisions								
	4 = Migration background								
	5 = Returns to education								
	6 = Interest, self-concept and motivation								
	7 = Socio-demographic information								
	a = Newborns and early childhood education								
	b = From kindergarten to elementary school								
	c = From elementary school to lower secondary school								
	d = From lower to upper secondary school								
	 e = From upper secondary school to higher ed./occ. training/labor market 								
	f = From vocational training to the labor market								
	g = From higher education to the labor market								
	h = Adult education and lifelong learning								
	s = Basic program								
	x = Generated variables								
3–7	Item number								
	Indicator for the item number which typically consists of four numeric characters								
	plus one alphanumeric character								
8–11	Suffixes (optional, see below)								
	Indicator for several types of variables; separated from the previous characters by an underscore								

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):

General Conventions

- 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. Versionized variables are listed in section B.3.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 signalizes that data in this variable is only available via on-site acces; _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).

Specific variables for (prospective) teachers

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

41

Table 5: (continued)

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)
_CD wb	Web/Internet-based test modus (inproctored)
_₩6	web/internet based test modus (unproteored)
_c	Scored item variable (s_c for partial credit-items)
_scl	Weighted likelihood estimate (WLE) ¹²
_sc2	Standard error for the WLE ²
_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)

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

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).

is surveyed in different survey waves or starting cohorts, the variable name is equiped 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 on 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 B.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 *NEPStools* 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 *NEPStools* 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.

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

Code	Meaning	Note
Item nonresp	oonse	
	n at was sheed	and a selection of few in a two sets with time west intigers (s. s.
-94	not reached	compotency test measures)
-95	implausible value	assigned by the survey agency (e.g. multiple answers to
55		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)
Not applicab	le	
-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)
-91	survey aborted	respondent quit interview, in CAWI mode
-92	question erroneously not asked	question not asked by mistake, in CAWI and CATI
-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
Edition missi	ngs (recoded into missing)	
-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

Table 6: Overview of missing codes

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 20 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 documetation report by Zielonka and Pelz, 2015.



Figure 20: 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 4.5.37). 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 areaspecific 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 21: 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

Data Structure

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

4.2 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.
- **splink** uniquely identifies episodes/spells across all datasets within each person. It is used to link data from Biography with Education or episode modules such as spVocTrain

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.

4.3 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_t) 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.

4.4 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 episode data, there is one row for each episode that was captured during the interview. 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. This means (especially for time-variant variables like ISEI or CASMIN) that the respective values indicate the status *at the given time of the episode*, and not necessarily the current status which is valid nowadays.

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 episode, this is reflected in a new data row.

In addition to such (time dependent) episode data, which we call *duration spells*, there are two other types of episode spells in our data:

- Occurring events or the transition from one state to another (e.g., change of marital status, change of educational level) are recorded in *event spells* with one row describing one state.
- the existence of children, partners, etc., is recorded in *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 or child. More detailed information on the required identifier variables can be found on the respective data file pages in in section 4.5. Please also note section 4.4.3 for a further complication of this matter.

One general remark: be aware that the number of episodes per se is independent of the survey wave. During one interview (one wave), there may be several episodes (several rows) recorded, or no episode at all. Also, the dates given in the episode relate to the time the episode was valid, whereas the wave relates to the interview date. They might not even overlap!

You should consider those two entities (spell and wave) as completely unrelated. Although there might be some situations where you have the need to know *when* the information of an episode has been collected, you are best advised to ignore the variable wave in episode data completely.

Do not try to use the variable wave to merge episode data to panel data. Although this might seem like the proper way to do this, episode data may contain multiple (or none) rows per wave and ID, while panel data contain exactly one row for every wave (of an ID). Such a merge results in the panel data obtaining an episode structure, which totally messes up the data.

A better approach seems prior to conduct such a merge, try to aggregate the episode data to *one information* for each interview date, or even just one information for the whole life course, so that in the end, you do not have more rows than waves (per respondent).

4.4.1 Edition of the life course

The life course data in the NEPS Starting Cohorts consists primarily of information on episodes of school attendance, participation in vocational preparation measures and vocational training or university education. Further it consists of information on exercise of compulsory or voluntary services (military module), employment and unemployment episodes, as well as spells of parental leave. We refer to these activities as *main activities*.

The episodes, grouped by episode type, are recorded independently in separate modules. The aim of recording these activities is to obtain chronologically complete life histories on the employment and training careers of our respondents. This requires two different edition steps of the data:

 After the episodes have been collected in the longitudinal modules, the first step in the edition process of the life courses already takes place during the interview. The episodes are summarized in the data revision module and put into their chronological order. Subsequently, they are checked for chronological gaps and overlaps. This test is carried out by a cooperative clarification of chronological gaps between interviewee and interviewer.

If chronological gaps are discovered in the data revision module, these are closed by subsequently recording additional episodes of the above-mentioned *main activities*. If there is no main activity for the examined period, the interviewee can close it with a so-called gap activity. In addition, gaps can be closed in the data revision module by adjusting the dates of the episodes between which a gap exists.

Chronological overlaps of episodes are discussed in the data revision module together with the interviewee. This may lead to a change of the dates of the episodes involved in the overlap. For inaccurate or missing dates, estimates are calculated in addition to the original dates, as far as there are reasonable indications for good estimates. For example, the imprecise specification about the starting month of an episode "Summer" is replaced by the value 7 "July" and saved in the biography file. In this way, even episodes with incomplete date specifications can be included in the chronological test and checked for gaps or overlaps with temporally adjacent episodes in the overall context of the life course (for general and specific functionality of the data revision module, see Ruland et al., 2016 and Matthes et al. 2005, 2007).

The result of this examination of the life courses during the interview are largely complete and time-consistent life courses.

2. Despite this meticulous examination during the interview, there are still minor inaccuracies in the consistency of life courses after the survey. For example, one-month overlaps of episodes are not edited in the data revision module. The same applies to gaps between successive episodes of up to two months. The test in the data revision module can also be interrupted or skipped at the request of the interviewee so that it was not carried out or not carried out completely.

For these reasons, a second, automated step of processing the time data of the life courses takes place after the end of the interview during data edition. The results of these temporal adjustments are also saved in the biography file. The automated edition is divided into several successive edition steps.

The first step is to remove one-month overlaps of episodes. A one-month overlap between two episodes is, in our definition, when the end date of a preceding episode is identical to the start date of the following episode. The procedure here is to shorten the end date of the previous episode by one month. The prerequisite is that the previous episode is longer than one month, otherwise this one-month episode would be shortened to the duration of zero. If the duration of the previous episode is only one month, the start date of the following episode is shortened by one month. If both episodes have a duration of one month, the dates are not edited.

Subsequently, one to two-month gaps between successive episodes are automatically closed. If the gap has a duration of one month, the end date of the previous episode is extended by one month. If there is a two-month gap, the start date of the following episode is additionally brought forward by one month.

Finally, chronological gaps in the life course that are larger than two months are closed by inserting new episodes for these gaps in the biography file, which close these gaps completely. These episodes are marked as *data edition gap* in the sptype variable of the biography file.

All these changes of the time specifications described are exclusively made in the biography file. The respondents' original information on the start and end dates of the episodes remain in the data files of the longitudinal modules (also see Künster 2015a, 2015b).

4.4.2 Revoked episodes

In order to reduce seam 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. At this point, approximate dates (e.g., seasons) are converted into specific months (first month of the respective season) in order to enable a check at all. 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 spms "Check module: type of event" (spms==-20 "Episode revoked in check module").

The addition of new episodes in the check module is indicated in the variable "Episode mode" (e.g., ts23550=4 in spEmp). A detailed description of the functionality of the check module for reported life courses is given in Hess et al., 2012 (in German language), which can be found on the documentation page:

> www.neps-data.de>Data Center>Data and Documentation
> Starting Cohort First-Year Students>Documentation

4.4.3 Subspells and harmonization of episodes

There is one important circumstance to consider when working with NEPS spell data. 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 still *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:

- first (right-censored) data row reported in initial wave (subspell=0 if this is the only subspell for the episode, subspell=1 if there are other subspells)
- continued episode reported in next wave(s) (subspell=2, subspell=3, etc.)

To make it easier for data users to work with these spread episode data, they are also summarized in a data line (record) according to defined rules. This data line reflects the most current information on the episode. This means that for completed episodes, the information valid at the end of the episode is selected and for episodes that were not yet completed at the last interview time, the information valid at the last interview time is selected. We call this process of summarizing information about an episode from different survey waves *episode harmonization*. It is described in detail below.

Episodes are defined by the assignment to a respondent (ID_t), by the type of episode (e.g., training episode), by an episode ID (splink, which typically consecutively numbers the episodes of the same type of episode of a case), and by the start and end date of the respective episode.

If an episode both begins and ends within the data collection period of a survey wave, then it can be assumed that this episode has been completely recorded with all the desired information (see figure 22, spell 1). In the SUF data of the corresponding longitudinal data file, there is a single data line for this episode, which contains the complete information.



Figure 22: Logic of subspells

However, there are many episodes that have not yet ended at the time of the interview of a survey wave, but are still ongoing at that time. Such persistent episodes are updated in the subsequent survey wave in which the respective person takes part. This means that further information on these episodes is recorded in the subsequent survey waves until the respondents report the episodes as finished (see figure 22, spell 2). In such cases, the information on an episode is stored separately in the SUF in one data line for each survey wave, so that the information on this episode is divided over several data lines and one data line of this episode, as well as the episode ID. There is, however, an additional variable subspell, which consecutively numbers the data lines that belong to one episode that was recorded over several survey wave, i.e., those that began and ended during the period covered by the survey wave, the variable subspell contains the value 0. The same applies to episodes that were recorded for the first time in the current survey wave and that were still ongoing at the time of the interview (see figure 22, spell 3).

The episode file for the cases shown in figure 22 corresponds to the data structure listed in table 7 before an episode harmonization. For the sake of simplification, the table contains data

from three consecutive surveys/waves, each conducted in december of the years 2009-2011. There is only one row of data for the first episode of the example case, because it was completed before the survey time of wave 2, i.e., it was completely recorded in this wave. Accordingly, the value of subspell is 0.

For the second episode, there are three data lines with the information on this episode from waves 2-4. The subspell variable for the second episode numbers the partial episodes from 1-3. The end of the second episode was reported in the fourth survey wave.

The third episode was recorded in the fourth survey wave. This episode continues, but since only a part of the episode has been reported so far, subspell is also initially given the value 0. This does not change until further information for this episode is recorded in a subsequent survey wave.

var2	var1	ongoing	end_y	end_m	start_y	start_m	subspell	wave	splink	ID_t
5	3	no	2009	april	2005	may	Θ	2	300001	1
	1	yes	2009	december	2009	june	1	2	300002	1
	•	yes	2010	december	2009	june	2	3	300002	1
8	•	no	2011	july	2009	june	3	4	300002	1
4	2	yes	2011	december	2011	august	Θ	4	300003	1

Table 7: Data lines of the example case in the SUF before spell harmonization

For episodes that last over several survey waves, the NEPS does not collect the same information in each survey wave. In the wave in which an episode is recorded for the first time, all unchangeable core information about this episode is collected. In the case of training episodes, this includes the type of training (e.g., vocational training or studies), the exact designation of the training occupation and some other parameters that distinguish this training from other training. This of course also includes the start date of the episode. This information will not be requested again when this episode is updated in later survey waves. Instead, additional characteristics of the episode, such as current pay, are recorded in these waves. As soon as the interviewee reports the episode as completed, information regarding the end is recorded. Such information is, for example, the achieved completion of a training episode and of course the end date of the episode. In this respect, the information on an episode, which was updated via different survey waves, is divided over the individual partial episodes (subspells) of this episode. The number of the partial episodes varies depending on the total duration of the episode.

In order to make it easier for data users to work with the data of updated episodes, the information from the partial spells of episodes is summarized in an additional data line. Therefore, besides the data lines for the partial episodes, there is also a data line that gives an overall overview of the updated episode and is referred to as the *harmonized episode*.

Thus, episode harmonization is only used if there are several partial spells for an updated episode from different survey waves.

The data line for the harmonized episode is added to the already available records in the longitudinal file. The variable subspell always has the value 0 for harmonized episodes. In our example case shown above, an additional data line would be added for the second episode as a summary of the three partial episodes of this episode in the longitudinal file (see table 8), since only the second episode has several partial spells in different survey waves.

ID_t	splink	wave	subspell	start_m	start_y	end_m	end_y	ongoing	var1	var2
1	300001	2	Θ	may	2005	april	2009	no	3	5
1	300002	2	1	june	2009	december	2009	yes	1	
1	300002	3	2	june	2009	december	2010	yes		•
1	300002	4	3	june	2009	july	2011	no		8
1	300002	4	Θ	june	2009	july	2011	no	1	8
1	300003	4	0	august	2011	december	2011	yes	2	4

Table 8: Data lines of the example case in the SUF before spell harmonization

Since a harmonized spell is a summary of all partial spells of an updated episode, exactly one piece of information must be selected from the partial spells for each variable, which is then transferred to the harmonized spell. In most cases the rule for selecting the relevant information that is transmitted is obvious. But if it is not, the following rules are applied:

- first For all questions that are only asked when a new episode is entered, i.e., when the episode is reported for the first time, the information for the harmonized spell is taken from the first partial episode because it can only be found there and is valid for the complete duration of the episode (see var1 in table 8).
- **last** For information that is either updated in every survey wave or that can only be found in the last partial spell after the end of the episode, the information for the harmonized spell is taken from the last partial episode (see var2 in table 8).

There is an exception concerning the application of the harmonization rule *last*. If an already established question in the longitudinal modules is generally not asked in a certain survey wave, then the undetermined value of the associated variable is replaced with the value -54 *missing by design* during data edition. The reasons for not asking the question can be manifold. If this question follows the harmonization rule *last*, the value -54 is not stored into the harmonized episode. Instead, the existing partial episodes of the episode concerned are searched for a value that deviates from the value -54 and this value is stored in the harmonized episode. The same procedure is used with the value -55 *not applicable*. The idea is that the value determined in this way is a good estimate of the missing last information on this item of this episode.

- **first nonmissing** The harmonization of most of the variables follows either the selection rule *first* or *last*. However, there are exceptions to this rule. An exception occurs, for example, if a new variable is introduced when recording episodes, which basically follows the *first* rule, but which should also be collected for episodes updated in the current survey wave. In such cases, the information on this variable is then also contained in the data for updated episodes, but is not in the first partial spell, but in a later partial episode. In these cases, the first valid value to be found in any partial spell of an episode is selected.
- **last nonmissing** There is a similar exception for variables that measure a changing state until a target state is reached. In the case of employment episodes, this can be, for example, changing from a fixed-term position of a specific job to a permanent one. In cases in which

an employment is temporary when it is first recorded, the question about the time limitation of the position is asked each time the episode extends over several survey waves. This continues until the employment either ends or the state of employment changes to *permanent*. Once this change from fixed-term to permanent job has been completed, the question of a time limitation is no longer asked when the episode is updated, since the reverse change from a permanent to a fixed-term job within the same job is hardly considered realistic. The information about the delimitation of the episode is therefore not necessarily in the first or last part of the spell. Here the last valid value of a partial spell of this episode is relevant. Therefore, in this case, the *last nonmissing* rule (last valid value to be found in the partial spells of an episode) is used for harmonization.

There is another exception in cases in which the continuation of an episode in the current survey wave is contradicted by the respondent during the life course assessment in the data revision module (see Ruland et al., 2016 for more information). This exception only affects episode types that are included in the life course assessment in the data revision module (episodes from the data files spSchool, spVocPrep, spVocTrain, spMilitary, spEmp, spUnemp, spParLeave, spGap). In such cases, we assume that the partial spells recorded in previous waves of the survey contain correct information on this episode up to the part of the episode that was contradicted, because they were subjected to a life course revision carried out together with the respondent in the previous waves of the survey. According to this logic only the part of the episode recorded in the current survey wave is contradicted by the respondent and not the complete episode. The information already collected and stored in a data line on the current partial spell (which was contradicted in the data revision module) can still be found in the longitudinal file, but is marked in the variable spms with the code -20 as episode canceled in the data revision module. During the harmonization, this cancellation has been considered by only filling the harmonized episode with values from the partial spells that are not marked as canceled, i.e., all partial spells except for the contradicted partial spell of the current survey wave. The end date of this episode is set to the interview time of the survey wave in which the last, uncontradicted information on this episode was recorded.

Coded occupational information is recoded in the harmonized episodes based on the information available there. Therefore, there may be differences between the values of the partial episodes and the harmonized episodes for these generated variables.

The Research Data Center keeps track on which harmonization rule was applied to variables of the longitudinal data for which episodes were updated across survey waves. Those harmonization tables are currently not publicized, but you can obtain the rules for specific variables upon request.

Data users can decide whether they want to use the harmonized spells for data analysis or whether the information from the subspells that reflects the changes in characteristics of these episodes over time is important to them. Both information are available in the longitudinal data files.

If the harmonized episodes are to be used, including the episodes that only consist of a single partial spell and therefore did not have to be harmonized, then it is sufficient to select all records

for which the value of the variable subspell is 0.

keep if subspell==0

Thereafter, all episodes should be excluded that were contradicted in the data revision module (variable spms = -20) and which at the same time do not belong to the harmonized episodes (variable spext = 0)¹. As described above, this step already has been included in the harmonization process for the harmonized episodes.

If, on the other hand, you do not want to use the harmonized episodes but the original partial spells of the episodes, then all records should be dropped where the variable subspell has the value 0 and simultaneously the variable spext has the value 1. Subsequently, it is also necessary to exclude all partial episodes that were contradicted in the data revision module (variable spms = -20).

4.5 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 23 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 16-0-0
** path where the data can be found on your local machine
global datapath Z:/Data/${cohort}/${version}
```

1 Also the variable spgen indicates whether an episode was originally reported as finished (spgen=0) or whether it is a harmonized (generated) episode (spgen=1).



Figure 23: 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.

4.5.1 Basics

					« go back t	o overview	
Description				Exemplary	variables		
Simplified in	nformation al	bout respondents	s in a	ID_t	ID target		
plain format	t			tx29000	Age at interview mor	nth (years)	
				t70000m	Date of birth: month		
				t70000y	Date of birth: year		
wide format	t: 1 row = 1 re	espondent		t700001	Gender		
ID variables need	ded to identify a s	ingle row		tx29003	Mother tongue: Gerr	man	
		ingle low		tx29004	Citizenship: German		
ID_t				tx29005	Born in Germany		
Other ID variable	es useful for linkag	26		t741001	Size of household (pe	ersons)	
none		<u>,</u>		tx29060	currently employed		
none				- tx29904	Main spells of type 'Emp'		
Number of varia	bles / number of r	rows in file			(number)		
84 / 17,909	Э						
Contains data fro	om waves						
Exemplary data	snanshot						
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	30.58	[w] female	yes	2	yes	9	
7011856	30.00	[w] temale	yes	3	yes	2	

This file contains the latest reported basic information on each respondent, e.g., sociodemographic variables like age in years (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. This dataset only contains data from CATI interviews, information from CAWIs is not integrated. **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!

4.5.2 Biography

				« go back to overview			
Description			Exempla	Exemplary variables			
Integrated and File structure spell format: 1 ID variables needed ID_t splink Other ID variables to wave sptype Number of variable 10 / 231,136 Contains data from 1 2 3 12 13 14	d edited life cours L row = 1 episode d to identify a single row useful for linkage es / number of rows in fi waves 4 5 6 7 15 16	e data of 1 respondent / le 8 9 10 11	ID_t splink wave sptype startm starty endm endy spms splast	ID target Link for spell merging Wave Spell type Episode start (month) Episode end (month) Episode end (year) Type of event Episode is ongoing			
Exemplary data sna	apshot						
ID_t	splink	wave	sptype	starty	endy		
7004828	360001	1	36	2001	2001		
7006993	220001	1	22	1997	2001		
7014365	270001	1	27	2010	2010		
7015940	240002	13	24	2014	2018		
7027895	240002	7	24	2010	2013		

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

63

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 4.4.2) 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
```

4.5.3 CohortProfile

					« g	b back to overview	
Description			Exemplary variables				
Paradata o	n the coh	nort's panel sample		ID_t wave	ID target Wave		
Iong 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				cohort tx80220 tx80521	NEPS Starting Cohort Participation/drop-out status Data available: interview target		
				tx80522	person Data available target person	competence data	
ID_i ID_tg	D_i ID_tg			tx80524	e: institution participation in		
19 / 286,5 Contains data f	544 from waves	DEF OT FOWS IN THE					
1 2 3 12 13 1 Exemplary data	3 4 5 4 15 1 a snapshot	6 7 8 9 10 6					
ID_t	wave	tx80220	tx80521		tx80522	tx80524	
7011366	1	Participation	yes		yes	yes	
7011366	2	Temporary drop-out	no	missing	by design	not determinable	
7011366	3	Temporary drop-out	no	missing	by design	not determinable	
7011379	1	Participation	yes		yes	yes	
7011379 7011379	3	Participation Participation	yes yes	missing	by design by design	yes 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)
```

Data Structure

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
4.5.4 EditionBackups

			« go ba	ck to overview
Description	Exe	emplary variables	5	
Backup of original data that were modified dur- ing the data edition process	ID_ wa	_t ve	ID targe Wave	t
File structure	ua	aset	Dataset	name
long format: 1 row = 1 changed value of a vari- able in a datafile	me	name rgevars	Variable ID-Varia merging	bles for
ID variables needed to identify a single row	SOU	rcevalue_num	Original	value (if
dataset varname ID_t wave splink subspell part- ner child	ed	tvalue_num urcevalue_str	numerio New val Original	c) ue (if numeric) value (if
Other ID variables useful for linkage			string)	
mergevars	ed	tvalue_str	New val	ue (if string)
Number of variables / number of rows in file				
14 / 21,545				
Contains data from waves				
1 2 3 4 5 6 7 8 9 10 11 12 13 14 15 16				
Exemplary data snapshot				
ID_t wave dataset varname	mergevar	s sourceval …	ue_num	editvalue_num
70084014 16 SpVoCEXTEXAM 1515304 1D	_t wave exa t wave exa	m	00	30.00
7009210 1 pTargetCATI t731306	ID_t wav	e	5.00	2.00
7014551 1 pTargetCATI t731306	ID_t wav	e	5.00	2.00
7025970 . spVocTrain ts15221 ID_t spl	ink subspel	1	28.00	27.00

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

Data Structure

- 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 orginal 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
```

4.5.5 Education

						" Bo buck to	overview
				Exemplary	variab	les	
upward tr	ransitions in	n educational		ID_t number datem	ID tar Sort r valid	get number since (month)	
: 1 row = 1	event (epi	sode) of 1 re-		— datey — tx28101 — tx28102 — tx28103	valid Recer years Recer	since (year) nt CASMIN of education = f(nt ISCED-97	CASMIN)
				tx28109	Chan	ge in educational	
es useful for lin	ıkage			– splink	classi Link for Exam	fication or spell merging number	
tx28100				- tx28100	Source of information of		
bles / number	of rows in file				educa	itional qualificati	on
3							
om waves							
4 5 15 16	6 7 8	9 10 11					
snapshot							
number	datey	tx28101	tx28102	tx281	93	splink	tx28100
1	2003	Θ	-20		0	220001	22
2	2007	3	10		2	220002	22
3	2010	5	13 10		3 0	220003	22
4	1999	0	-20		9	220001	24 22
1 2	2005	3	-20		2	220001	22
2	2005	5	10		∠ 2	220002	22
3	2006	5			<u> </u>	220000	,,
	upward tr : 1 row = 1 ded to identify es useful for lin bles / number bles / number 3 om waves 4 5 15 16 snapshot number 1 2 3 4 1 2	upward transitions in : 1 row = 1 event (epi ded to identify a single row es useful for linkage bles / number of rows in file bles / number of rows in file 3 om waves 4 5 6 7 8 15 16 snapshot number datey 1 2003 2 2007 3 2010 4 2015 1 1999 2 2005	upward transitions in educational : 1 row = 1 event (episode) of 1 re- ded to identify a single row es useful for linkage bles / number of rows in file 3 om waves 4 5 6 7 8 9 10 11 15 16 7 8 9 10 11 snapshot	upward transitions in educational: 1 row = 1 event (episode) of 1 re-ded to identify a single rowes useful for linkagebles / number of rows in file3om waves456789101516snapshot120030-2007310320105119990-2005310	Exemplary 1upward transitions in educational ID_t numberdatemdateytx28101tx28102tx28103tx28103tx28109es useful for linkagebles / number of rows in filesmapshotnumberdateytx28101tx28102tx28101tx28100tx28101tx28101tx28101tx28101tx28101tx28102tx28101tx28101tx28101tx28102tx28101tx28101tx28102tx28101tx28101tx28102tx28101tx28101tx28101tx28101tx28101tx28101tx28101tx28101tx28101tx28102tx28101tx28103tx28101tx28104tx28105tx28105tx28105tx28104tx28105 </td <td>upward transitions in educational ID_t ID tar ID_t ID tar number Sort r datey valid : datey valid : datey valid : tx28101 Recer tx28102 years tx28103 Recer tx28109 Chang classif splink suseful for linkage splink waves 5 4 5 6 7 snapshot 1 2003 0 -20 1 2003 0 -20 0 2 2007 3 10 2 3 2010 5 13 3 4 2015 8 18 9 1 1999 0 -20 0 2 2005 3 10 2</td> <td>upward transitions in educational ID_t ID target number Sort number datem valid since (year) tx28101 Recent CASMIN datey valid since (year) tx28102 years of education af (tx28102 ded to identify a single row tx28103 es useful for linkage splink bles / number of rows in file splink 3 m waves 4 5 6 7 8 9 10 11 15 16 11 12 2007 3 10 2 220001 1 2007 3 10 2 220002 3 10 2 220001 1 1999 0 -20 0 220001 11 1999 0 220001 10 2 220001 1 1999 0 -20 0 220001 10 220001 1 1999 0 10 2 220002 220001</td>	upward transitions in educational ID_t ID tar ID_t ID tar number Sort r datey valid : datey valid : datey valid : tx28101 Recer tx28102 years tx28103 Recer tx28109 Chang classif splink suseful for linkage splink waves 5 4 5 6 7 snapshot 1 2003 0 -20 1 2003 0 -20 0 2 2007 3 10 2 3 2010 5 13 3 4 2015 8 18 9 1 1999 0 -20 0 2 2005 3 10 2	upward transitions in educational ID_t ID target number Sort number datem valid since (year) tx28101 Recent CASMIN datey valid since (year) tx28102 years of education af (tx28102 ded to identify a single row tx28103 es useful for linkage splink bles / number of rows in file splink 3 m waves 4 5 6 7 8 9 10 11 15 16 11 12 2007 3 10 2 220001 1 2007 3 10 2 220002 3 10 2 220001 1 1999 0 -20 0 220001 11 1999 0 220001 10 2 220001 1 1999 0 -20 0 220001 10 220001 1 1999 0 10 2 220002 220001

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 <mark>en</mark>
** 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
```

4.5.6 MethodsCATI



This dataset offers a variety of information on the data collection, e.g., gender ($t \times 80301$) and age ($t \times 80302$) of the interviewer; interview date (intm, inty); interview duration ($t \times 80209$); incentives ($t \times 80210$); and individual survey participation ($t \times 80220$).

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
```

Data Structure

tab wave tx80220

 $\star\star$ how many different interviewers did CATI surveys? distinct ID_int

4.5.7 MethodsCAWI

				« go b	ack to overview	
Description			Exemplary	variables		
Paradata from	the targets CAWI int	terview	ID_t wave	ID target Wave		
File structure			tx80208	Interview: length c	of	
long format: 1	row = 1 target in 1 v	vave		questionnaire (min	nutes)	
ID variables needed	to identify a single row		tx80225	Interview: last deli	very status	
			tx80250	Interview: winners	s of the lottery	
ID_t wave			tx80200	Interview: number	of all contact	
Other ID variables us	seful for linkage			attempts		
none	none		tx80206	Interview: number of		
none				interruptions		
Number of variables	/ number of rows in file					
21 / 21,972						
Contains data from v	vaves					
1 2 3 4 12 13 14 1		9 10 11				
Exemplary data snap	ishot	+~20209		+~20250	+++20206	
10_1 7001982	11	14.11		0	0	
7001982	14	-54.00		-54	-54	
7002030	11	25.96		0	Θ	
7002030	14	18.14		-54	Θ	
7002115	11	26.18		0	0	
7002115	14	21.80		-54	1	
7002191	11	19.14		0	Θ	

This dataset offers a variety of information on the data collection, e.g., teacher over-sample (tx80122); interview duration (tx80208); winners of the prize draw (tx80250); and the number of interruptions during the interview (tx80206).

Importantly, MethodsCAWI contains all contacted respondents whether an interview was realized or not. Thus, MethodsCAWI includes more cases than pTargetCAWI. MethodsCAWI provides data from wave 11 onwards because paradata on CAWI data prior to wave 11 was collected in a different way. Perhaps paradata for earlier waves will be included in future releases.

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)

Data Structure

```
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
```

4.5.8 MethodsCompetencies

					« go back to	overview
Description				Exemplary	y variables	
Paradata fro	m the targets	s competency t	ests	ID_t	ID target	
				ID_i	Institution ID	
				wave	Wave	
long format:	1 row = 1 ta	rget in 1 wave		ID_tg	Test group ID	
ID variables need	led to identify a si	ngle row		testm	Test: Survey day (month	ר)
	-			testy	Test: Survey day (year)	
ID_t wave				tx80422	Survey mode (realized)	
Other ID variable	s useful for linkag	e		ID_int	Interviewer: ID	
	int			tx80301	Interviewer: gender	
				tx80302	Interviewer: age group	
Number of variab	oles / number of r	ows in file		tx80303	Interviewer: highest	
73 / 29 /17	,				school-leaving qualifica	tion
75 / 25,417				tx80661	Number participants	
Contains data fro	m waves			tx80628	Questions about test ta	sks
		7 0 0	10 11			
	4 5 0					
12 13 14						
Fueren la muelata e						
ID t	wave	ID int	tx80301		tx80302	tx80303
7002115	12	-54	-54	mis	sing by design	-54
7002163	1	1385	1		30-49 years	7
7002163	12	-54	-54	mis	sing by design	-54
7002189	1	1385	1		30-49 years	7
7002189	12	-54	-54	min s:	sing by design	-54
7002204	1 5	1318	2		50-49 years	10
7002204	12	-54	-54	mis	sing 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

4.5.9 pTargetCATI

					« go back to	o overview	
Description				Exemplary	variables		
Data from r	espondents	CATI question	inaires	ID_t	ID target		
File structure				ID_i	Institution ID		
long format	: 1 row = 1	target in 1 way	ve	wave t431000	Wave Migration sentiment		
ID variables nee	ded to identify	a single row		t531214	Tuition loan		
ID_t wave				t531250	Post-recording final gr Source of finance: fam	ade nily	
Other ID variabl	es useful for lin	kage		tg24503	Employment context o	loctorate	
ID_i				t712001	Kindergarten Gender		
Number of varia	ables / number o	of rows in file		t70000y	Date of birth: year		
1,050 / 10	1,038			t514001	Satisfaction with life		
Contains data fr	om waves			t514008	Satisfaction with cours Size of household	se of study	
1 2 3 12 13 14	4 5 15 16	6 7 8 9	9 10 11				
Exemplary data	snapshot	170.1.000		1700001	170000		
ID_T 7003570	wave	t/244⊍3 -54 0	tg245⊍3		T/0000y	T514008	
7014070	13	- 54.0	5	[m] male	1991	0 10	
7014178	9	-54.0	4	[w] female	1989	8	
7018526	10	2.6		[m] male	1990	8	
7026228	16	0	3	[w] female	1991	8	

The data in file pTargetCATI are from computer assisted telefone 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.²

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)).

² includes all students of the panel sample regardless of their questionnaire participation.

Data Structure

The file also includes information on the study program the respondents had started in winter term 2010/2011. The data were collected in the initial questionnaire and are stored in variable tg01003_g1 (type of higher education institution), variables beginning with tg0400 (different classifications of up to three subjects, information on majors or minors), variables tg02001, tg02001_g1 (intended degree), variables tg03001_g1, tg03001_g2 (type of intended teaching degree), and variables tg15207_g1R, tg15207_g2R (location of the higher education institution).

The initial questionnaire was mainly administered as a self-administered written survey. Partly, it was integrated into the first telephone interview. In this case, only the basic questions were asked. To avoid overstraining the participants, questions of minor relevance were omitted. Among these questions are those on admission restrictions (tg10001, tg11001, tg11002, tg11003), availability, use and quality of measures to facilitate integration into higher education (variables beginning with the string tg0800), and attitudes of parents and peers towards the study decision (variables beginning with tg1500).

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

4.5.10 pTargetCAWI

				« go back	to overview	
Description			Exemplary	y variables		
Data from re	spondents CAW	questionnaires	ID_t	ID target		
File structure			wave	Wave		
	1		ID_i	Institution ID		
long format:	1 row = 1 target	In 1 wave	t242020	Quality equipment: I	iterature	
ID variables need	led to identify a single	row	t242107	Higher education inst	titution	
	,,			activities: sport		
ID_t wave			t289902	Living in shared living	ł	
Other ID variable	s useful for linkage		t514001	Satisfaction with life		
			t272061	Motivation for courses/trainings		
ו_טו			t30300b	Amount of rent		
Number of variab	oles / number of rows i	n file	_ tg51004	Course of study cance	eled/inter-	
1 400 / 54 /				rupted/completed		
1,408 / 54,4	465		tg74011	Time budget: doctor	ate work	
Contains data fro	m waves		t241011	Time budget semeste	er: courses	
				Time budget semeste		
1 2 3	4 5 6 7	8 9 10 11				
12 13 14	15 16					
Exemplary data si	napshot					
ID_t	wave	t289902	t514001	t30300b	tg51004	
7004341	8	1	7	400	2	
7008707	8	1	5	345	2	
1010/58	8	1	(460	2	
1012056 7012410	8 0	1	8	270	2	
1013419	0	Ţ	4	200	T	

Apart from computer assisted telefone 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 (\pm 30300b), satisfaction with life (\pm 514001), having a roommate (\pm 289902), and there are also variables to help you to identify if a target is currently studying (\pm 51000, \pm 51001, \pm 51004). In contrast to CATIs, CAWIs are self-administered. Furthermore, biographical data such as episodes of employment or episode of vocational training were not collected.

Since wave 11, data on the device used during the online survey (tg5910*), the screen size (tg5911*), and also the survey setting (tg5920*) are collected and published in the SUF. These data enable new possibilities of method research. 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
```

4.5.11 pTargetCORONA

					« go back t	o overview
Description				Exemplary	variables	
Data regardiı demic on res	ng the impac pondents life	t of the corona p	an-	ID_t wave	ID target Wave	
File structure				t514001	Satisfaction with life	
long format:	1 row = 1 tar	get in 1 wave		t514010	Satisfaction with study/training/school	
ID variables neede	ed to identify a sir	ngle row		tm00001	Impact: Coronavirus in	nfection -
ID_t wave				tm00007	no Impact: Quarantine - I	no
Other ID variables	s useful for linkage	2		- tm00013	Employment status be	efore
none				_ tm00015	coronavirus pandemic Systemically importan	t
Number of variab	les / number of ro	ows in file			profession	
185 / 2,859						
Contains data from	n waves					
1 2 3 12 13 14						
Exemplary data sr	napshot					
ID_t	wave	t514001	t514010		tm00013	tm00015
7004851	•	7	7	I	was employed	no
7007964	•	8	6	I	was employed	yes
7018040	•	ŏ	10	ц т	was employed	yes
7029737		6	7	I	was employed	no

This data has been established to investigate the medium and long-term effects of the corona pandemic on skills development and educational pathways over the life course. The following questions aire in particular:

- How do learning environments change and which potentials and risks become clear through the beginning digitalization of learning?
- Are there effects on upcoming educational decisions and are there medium and long-term effects on social educational inequality
- What are the effects on educational outcomes, such as income, but also non-monetary returns, e.g., health and labor market participation

Data is collected by means of a cross-cohort questionnaire program adapted to the current situation of the respective participants. In order to collect this data in a timely manner, the first

questions were administered via online survey in Starting Cohorts 2-6 in May 2020. As this time span did not overlap with regular waves, data from this survey is marked with a missing wave (wave==.). The integration of the corresponding questions is planned in an additional module on the corona pandemic for the forthcoming main surveys in all Starting Cohorts.

Example 11 (Stata): Working with pTargetCORONA

```
** open the file
use ${datapath}/${cohort}_pTargetCORONA_D_${version}.dta, clear
label language en

** note that the wave is missing for some cases,
** as this reflects the pre-wave survey in may 2020
tab wave

** but rows can be uniquely identified by ID_t and wave
isid ID_t wave
```

4.5.12 pTargetMicrom

						« go bac	k to overview
Description					Exemplary varial	oles	
Small-scale place of re	e regional in sidence	dicators	on responden	ts'	ID_t wave	ID target Wave	
File structure					regio	Indicator	for
panel form of 1 respor	nat: 1 row = ndent	1 region	al level in 1 wa	ve	ID_regio	enrichme System-fi enrichme	nt level ree ID of nt level
ID variables ne	eded to identify	/ a single ro	W		mso_k_ausland	Share for	eigners
ID_t wave	regio				mso_k_familiembe_k_haustyp	Family st Type of h	ructure ouse
Other ID varial	bles useful for li	nkage		_	mgm_k_dom	Dominan milieu®	t microm geo
ID_regio					mgs_k_dom	Dominan	t
Number of var	iables / number	of rows in	file			geo-subn	nilieu
188 / 197	,552				mmo_k_volume	n Move vol	ume
-					mpi_k_dichte	Car densi	ty
Contains data	from waves				mas_k_berufsuv	Occupati	onal disability
1 2 3 12 13 1	345.41516	6 7			mas_k_krankzuv	insurance Additiona insurance	e al health
Exemplary data	a snapshot						
ID_t	wave	regio	ID_regio	mso_k_aus	land mbe_k_	haustyp	<pre>mpi_k_dichte</pre>
7009879	7	1	145167		8	6	1
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 in wave 5 (data for waves 1, 3 and 7 have been enriched at a basic level).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 12 (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 detailled 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
```

4.5.13 spChild



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 13 (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
```

4.5.14 spChildCohab

					« go back t	o overview			
Description				Exemplary	variables				
file listing co	habitation sp	ells with child	Iren	ID_t child	ID target Child number Spell number cohabitation with child pell Number of subspell				
spell format respondent	: 1 row = 1 c	cohabitation t	ime of 1	spell subspell					
ID variables need ID_t spell su	ded to identify a si bspell	ngle row		wave Wave ts3331m Start date Living together Child (month)			wave Wave ts3331m Start date Livin (month)		ther Child
Other ID variable	es useful for linkag	е		— ts3331y	Start date Living together Child (year)				
child wave				ts3332m	End date Living with child				
Number of varial 20 / 4,382	bles / number of r	ows in file		ts3332y	End date Living with child Currently living together with child				
Contains data fro 1 2 3 12 13 14	4 5 6 15 16	789	10 11						
Exemplary data s	snapshot								
ID_t	child	spell	subspell	wave	ts3331y	ts3332y			
7003340	1	101	1	3	2012	2012			
7004215	1	101	1	3	2010	2012			
1004184 7012623	⊥ 2	707 T0T	1	3	2011	2012			
7014667	1	101	2	5	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 14 (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:
\star 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)
\star\star 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
```

4.5.15 spCourses



This module comprises courses and trainings attended during episodes of employment (spEmp), unemployment (spUnemp), parental leave (spParLeave), military, or civilian service (spMil-itary), as well as episodes from the spGap module. It comprises all spells from the past 12 months prior to the first interview that were recorded in the modules mentioned above. For follow-up interviews data on courses is collected up to three years retrospectively in case of temporary drop-outs between waves. The starting and end dates of the spells in this module represent the original starting and end dates of episodes (in which a course was taken) but not the start and end of the courses themselves.

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. For each of these episodes, information on up to three courses is included in wide format and therefore the course enumerators is stored in wide format (course_w1, course_w2, and course_w3), whereas in the other course modules (spFurtherEdu1 andspFurtherEdu2) 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 15 (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

4.5.16 spEmp



This extensive module covers all spells of regular employment, including traineeships, preparatory service (e.g, for the teaching and legal profession), and internships (only in case that the target persons are not studying). Information on internships while studying is included in spInternship. 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), type of student employment (tg2608b), quality of student jobs (t265401-t265423) and internships (tg26300-tg2630i). Have a look at pTargetCATI and pTargetCAWI for more fine-grained information on teacher training and the situation of teachers.

```
Example 16 (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
```

4.5.17 spFurtherEdu1

					« go back t	o overview
Description				Exemplary vari	iables	
information a	bout additiona	courses		ID_t wave	ID target Wave	
entity format:	1 row = 1 cour	se of 1 respondent		– course – t271048	Course number Course is ongoing	
ID variables needed	d to identify a single	row		t271049	Termination Cour	se
ID_t course				t272000_0 	Content other cou Other courses 2	urse
Other ID variables	useful for linkage			t271051	Other course	
wave				t272000_g13	Content other cou (course ID)	urse
Number of variable	es / number of rows	in file				
18 / 8,029						
Contains data from	n waves					
1 2 3 12 13 14	4 5 6 7 15 16	8 9 10 11				
Exemplary data sna	apshot					
ID_t 7013825	wave 13	course	t2710	148	t271050	t271051
7013885	16	1608		no	no	ves
7015410	15	1502		no	no	no
7017051	9	902		no	no	yes
7018082	15	1506		no	no	no

This module contains information on further courses (also private courses) since the last interview 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).

Example 17 (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)
```

Data Structure

```
** 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
```

4.5.18 spFurtherEdu2

Description		Exemplary	variables	
information about courses		ID_t	ID target	
File structure		wave	Wave	
		course	Course number	
entity format: 1 row = 1 course of 1 res	pondent	t279040	Professional/private reas	sons
ID variables needed to identify a single row		t279046	Course costs Employer	
		t272040	Provider	
ID_t course		t279041	Motivation for course	
Other ID variables useful for linkage			attendance	
		t272043	Certificate	
wave		t272003	Course assessment: lear	ned
Number of variables / number of rows in file			new things	
28 / 16 910				
Contains data from waves				
1 2 3 4 5 6 7 8 9 12 13 14 15 16	10 11			
Exemplary data snapshot				
ID_t wave course	t279046		t279041	t272043
7003501 7 701	fully		some effort	1
(00/151 15 1501 7000200 12 1201	tully		some effort	3
7010365 10 1002		2	some ettort	1
7013432 15 1503	fully	a	lot of effort	1

« go back to overview

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.

Example 18 (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)
```

4.5.19 spGap



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 19 (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

4.5.20 spInternship

					« go back to				
Description				Exemplary	variables				
reported int	ternship epis	odes		ID_t	ID target				
File structure				splink	Link for spell merging Spell number Number of subspell Wave Start month Internship episode Start vear Internship episode				
spell format	t: 1 row = 1 i	nternship episoo	de of 1	spell					
ID variables nee	ded to identify a s	single row		tg3607m					
ID_t spell su	ıbspell			tg3607y tg3608m	y Start year Internship episode m End month Internship episode				
Other ID variable	es useful for linka	ge		tg3608y	End year Internship episode Ongoing of internship episode Type of internship Average working hours				
wave splink				tg36109 tg36110					
Number of varia	bles / number of	rows in file		tg36111					
40 / 39.37	1				Internship				
				tg36119	Placement as an inter	'n			
Contains data fro	om waves								
1 2 3 12 13 14	4 5 6 15 16	789	10 11						
Exemplary data	snapshot	I	1			I			
ID_t	spell	subspell	wave	tg3607y	tg3608y	tg36111			
7003809	2	1	5	2013	2013	40			
7014121	∠ 4	2	12	2015	2013	25			
7014498	5	2	9	2012	2014	25			
7017996	10	1	10	2016	2016	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, renumeration, learning content, and other key aspects have been surveyed.

Example 20 (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

4.5.21 spMilitary

					« go back to overview		
Description				Exemplary variables			
military / civilian service and voluntary gap years				ID_t ID target splink Link for spell merging			
File structure				subspell	Number of subspell		
spell format: $1 row = 1 episode of 1 respondent$			lent	spell Spell number			
speniormat		souc of i respond		wave	Wave		
ID variables needed to identify a single row				ts21201	Type of military service episode		
ID_t spell subspell				ts2111m	Start Military service episode - month		
Other ID variables useful for linkage				ts2111y	Start Military service episode -		
wave calial					year		
wave splink				ts2112m	End Military service e	pisode -	
Number of variables / number of rows in file					month		
21 / E 1EE				ts2112y	End Military service e	y service episode -	
21 / 5,155					year		
Contains data from waves							
1 2 3 4 5 6 7 8 9 10 11 12 13 14 15 16							
Exemplary data	snapshot						
ID_t	splink	subspell	spell	wave	ts2111y	ts2112y	
7002206	250001	1	1	1	2008	2010	
7005236	250001	2	1	12	2016	2016	
/006146	250001	3	1	9	2012	2013	
1013361 7031640	250001	3 1	1	5	2008	2013	
1031040	230001	Ŧ	Ŧ	9	2013	2013	

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 21 (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
4.5.22 spParLeave

					« go back t	o overview	
Description				Exemplary variables			
episodes of	parental leav	e		ID_t child	ID target Child number		
spell format of 1 respond	t: 1 row = 1 p dent	oarental leave	episode	spell subspell wave	Spell number Number of subspell Wave		
ID variables nee	ded to identify a si	ingle row		ts2711m	Start Parental leave (r	month)	
ID_t spell su	ubspell			ts2711y ts2712m	Start Parental leave (/ear) ionth)	
Other ID variable	es useful for linkag	je		ts2712y	End Parental leave (year)		
wave child s	splink						
Number of varia	ables / number of r	ows in file					
29 / 4,041							
Contains data fr	om waves						
1 2 3 12 13 14	4 5 6 15 16	789	10 11				
Exemplary data	snapshot						
1D_t	child 2	spell	subspell	wave	ts2711y	ts2712y	
7008949	2	101	1	12	2010	2017	
7014571	1	101	2	9	2014	2014	
7019173	- 4	404	2	9	2013	2014	
	-	101	-	10	2015		

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. The employment spell is not neccessarily interrupted if the mother is on parental leave as part-time work is legal.

Example 22 (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

4.5.23 spPartner

« go back to overview Description Exemplary variables history of partners ID target ID_t partner Partner number File structure Number of subspell subspell entity format: 1 row = 1 partner of 1 respondent ts31204 Partner: born Germany/abroad ts31211 Partner German ID variables needed to identify a single row ts31203 Gender of partner ID_t partner subspell ts3141m Marriage date (month) ts3141y Marriage date (year) Other ID variables useful for linkage Year of birth Partner ts3120y wave Start date Partnership - month tg2811m Number of variables / number of rows in file tg2811y Start date Partnership - year tg2804m End date Partnership episode 109 / 78,849 (month) Contains data from waves End date Partnership episode tg2804y (year) 3 5 9 10 ts31206 Age at immigration Partner 13 15 Exemplary data snapshot ID_t subspell ts31203 ts3120y tg2811m tg2804y partner tg2811y tg2804m 7014135 1 [m] male 1967 1992 12 2010 1 2 7011662 [m] male 1971 7 1993 4 2011 1 1 0 2 2008 3 2011 7016992 1 [w] female 1991 7018280 1 0 [m] male 1994 5 2014 8 2014 7002045 [m] male 1988 2013 2016 1 0 12 1

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 partners since the beginning of winter term 2010. For earlier partners, only information on the year of birth and education is available. 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 23 (Stata): Working with spPartner (find R example here)

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

Data Structure

```
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)
\star\star 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)
```

4.5.24 spSchool

					« go back ti	JOVEIVIEW	
Description				Exemplary	variables		
general scho	oling history			ID_t	ID target		
File structure				splink	Link for spell merging		
spell format	$\cdot 1 row = 1 sc$	hool enisode of	1 re-	subspell	Number of subspell		
spondent	. 1100 - 130			spell	Spell number		
spondent				wave	Wave		
ID variables need	led to identify a sin	gle row		ts11204	Type of school		
ID_t spell su	bspell			ts1111m	Start date month Scho episode	loc	
Other ID variable	s useful for linkage			— ts1111y	Start date year School	episode	
				ts1112m	End month School epi	sode	
wave splink				ts1112y	— ts1112y End year School episode		
Number of varial	oles / number of ro	ws in file		ts11209	School-leaving qualified	cation	
				ts11214	ng		
/3 / 4/,058					qualification		
Contains data fro	m waves			ts11218	Final grade school-lea	ving	
					certificate		
1 2 3	4 5 6	7 8 9 10					
12 13 14	15 16						
Exemplary data s	napshot						
ID_t	splink	subspell	spell	wave	ts1111y	ts1112y	
1005320 7012883	220002	U 1	2	1	2000	2009	
7012883	220003	2	3	1	2008	2011	
7014127	220005	2	5	5	2000	2013	
7015073	220004	2	4	3	2007	2011	

This module covers each respondent's general education history from school entry until the date of completion, or in case of enduring episodes, date of interview. This includes

- 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 24 (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
```

4.5.25 spSchoolExtExam

					« go bao	ck to overview	
Description				Exemplary var	iables		
school certif external stuc	ficates acqu dents' exam	uired by rec ination	cognition or	ID_t wave	ID target Wave		
File structure				exam	Exam numbe	r	
entity forma	t: 1 row = 1	exam of 1 r	respondent	ts11300	Awarded qua Germany?	lification in	
ID variables need	led to identify a	single row		ts1130m	Date month o	qualification	
ID_t exam				ts1130y	Date year qu	alification was	
Other ID variable	s useful for link	age			awarded		
wave				ts11302	Awarded sch qualification	ool-leaving	
Number of variab	oles / number of	f rows in file		ts11300g1	Awarded qualification in		
28 / 812					Germany? (edited) Country of awarded		
Contains data fro	m waves						
1 2 3 12 13 14	4 5 6 15 16	6 7 8					
Exemplary data s	napshot						
ID_t	wave	exam	ts11300	ts1130y	ts11302	ts11300_g1	
7006792	1	1	1	2007	•	1	
7012231	1	1	1	2004	• 4	1	
7013713	1	1	1	2007	4	1	
7018475	- 7	- 1	1	2013	5	1	

spSchoolExtExam comprises information about school exam certifications that have not been acquired through "regular" schooling. 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 classes)
- certificates that are automatically awarded by advancing through grades in upper secondary education or by completing vocational training

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

```
\star\star aim of this example is to evaluate the age of the respondent \star\star 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)
nepsmiss ts11302
** show some deviation
tabulate ts11302, summarize(age)
```

4.5.26 spSibling

						« go back to c	verview		
Description				E	Exemplary variables				
siblings of re	espondent	t			D_t /ave	ID target Wave			
entity forma	at: 1 row =	= 1 sibling of :	1 respondent	ti t	ibling x80211 g3270m	Sibling number Survey/Test instrument Month of birth Sibling			
ID_t sibling		,		tį tį	g3270y g32706	Year of birth Sibling Is sibling still alive?	20		
Wave	blos (pumbo	r of rows in file			g32708 g32709 g32711	Unemployment Siblings Highest school-leaving	IB		
11 / 26,932	2	or tows in me		tį	g32724	Sibling lives with parents			
1 2 3 12 13 14	4 5 15 16	6 7 8	9 10 11						
Exemplary data s	snapshot								
ID_t	wave	sibling	tg3270m	tg3270y		tg32708	tg32711		
7007331	1	1	5	1990		unemployed	5		
7014616	1	1	1	1988	fu	ll-time employed	5		
7014829	1	1	12	1979	fu	ll-time employed	5		
7015440 7018659	1	1 2	12 6	1983 1986	fu	ll-time employed unemployed	1 5		

spSibling contains all siblings of the respondent reported in wave 1. Each sibling is stored in one row, containing information about the date of birth (tg3270m/y), employment status (tg32708), and highest degree (tg32711).

Example 26 (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 // keep wave-1-data only as data on sibling was only collected once
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</pre>
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 identificator
keep ID_t total*
duplicates drop
```

4.5.27 spUnemp



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 27 (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
```

Data Structure

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

4.5.28 spVocBreaks

				« go back to o	verview
Description			Exemplary	variables	
breaks during vo	ocational training		ID_t splink	ID target Link for spell merging	
spell format: 1 r respondent	row = 1 break of 1 ep	isode of 1	break ts1531y	Training interruption numbe Start year Interruption of training episode	er
ID variables needed to	identify a single row		ts1531m	Start month Interruption of	
ID_t splink break	K		ts1532y	End year Interruption of	
Other ID variables user	ful for linkage		ts1532m	End month Interruption of	
Number of variables /	number of rows in file		tg2419a	training episode Status during interruption:	
17 / 1,986				semester off Status during interruption:	
1 2 3 4	5 6 7 8 9	10 11			
Exemplary data snapsh	hot				
ID_t	splink	break		ts1531y	ts1532y
7015794	240004	1		2016	2016
7012784	240005	1		2019	2020
7007363	240001	1		2012	2013
7015501	240001	2		2016	2016

This module covers all breaks of further trainings, vocational and/or academic, that a respondent ever attended – with a special focus on academic education. Information on vocational breaks were part of spVocTrain in prior data releases. Since release 16-0-0 break episodes are being extracted and edited to spVocBreaks. The data structure of breaks has been transformed from wide format to long format. In this dataset ...

- different types of breaks (break semesters, de-registrations, non-formal breaks) are included.
- includes several breaks within a single episode of a person.
- closes gaps between succeeding breaks (< 3 months) and combines overlapping breaks to continuous breaks.
- breaks within breaks were deleted.
- dates of beginnings and endings were corrected and stored as variables with _g1-suffixes.

- every break is in a separate row.
- splink helps you to merge data to spVocTrain and Biography as well.

Example 28 (Stata): Working with spVocBreaks

```
** example 1: merge spVocBreaks and spVocTrain
** open the vocational breaks
use ${datapath}/SC5_spVocBreaks_D_${version}.dta , clear
** reshape study breaks to wide format to match data with spVocTrain; first add _w-
 suffix to variables
foreach var of varlist ts15310_g1 ts15310_g2 ts1531y ts1531m ts1531m_g1 ts1531y_g1
 ts1532c ts1532y ts1532m ts1532y_g1 ts1532m_g1 tg2419a tg2419b tg2419c {
       rename `var' `var'_w
}
reshape wide *_w, i(ID_t splink) j(break)
** save this file temporarily
tempfile tmp
save `tmp'
** open the data file
use ${datapath}/SC5_spVocTrain_D_${version}.dta, clear
** only keep full or harmonized episodes
keep if subspell==0
** merge the previously created temporary data file to this
merge 1:1 ID_t splink using "`tmp'" , keep(using match)
** you now have put together information of breaks and the vocational track to
 analyze the students with breaks. The number of total episodes with breaks reduces
 the amount of rows of the combined dataset.
** example 2: merge spVocBreaks and Biography (further data preparation to analyze
 data is recommended)
** open the vocational breaks
use ${datapath}/SC5_spVocBreaks_D_${version}.dta , clear , clear
*merge breaks with biography data
merge m:1 ID_t splink using ${datapath}/SC5_Biography_D_${version}.dta, clear
\star\star now you could cut those vocational episodes using dates of episodes and breaks to
 re-define vocational episodes
*****
```

4.5.29 spVocExtExam

				« go back	to overview	
Description		Exemplary va	Exemplary variables			
vocational ed side of the reg	ucation certifica gular German e	ates acquired out- ducational system	ID_t wave	ID target Wave		
Eilo structuro			exam	Exam number		
entity format:	: 1 row = 1 exan	n of 1 respondent	ts15301_g1	Professional/spec title (KldB 1988)	ialization	
ID variables neede	d to identify a single i	row	ts15301_g4	Professional/spec	ialization	
ID_t exam			ts15301_g6	title (ISCO-08) Professional/spec	ialization	
wave	useful for linkage		ts1530m	End month Extern	nal	
Number of variable	es / number of rows i	n file	ts1530y	End year External		
30 / 3,187				examination		
Contains data from) waves		- ts15304	External examina	tion	
1 2 3 4 5 6 7 8 9 10 11				External examination in Germany/abroad		
12 13 14	15 16					
Exemplary data spa	anshot					
ID_t	wave	exam	ts1530m	ts1530y	ts15304	
7004603	10	1	4	2016	19	
7005077	12	1	4	2017	30	
7010062	13	1	6	2018	30	
7016720	12	1	2	2017	30	
7035175	10	1	9	2015	17	

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 trainig 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 29 (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 <mark>en</mark>
** 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)
nepsmiss ts15304
** show some deviation
tabulate ts15304, summarize(age)
```

4.5.30 spVocPrep



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 30 (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
```

4.5.31 spVocTrain

					« go back to	overview	
Description				Exemplary	variables		
vocational e	ducation his	story		ID_t spell	ID target Spell number		
spell format	: 1 row = 1 e	pisode of 1 resp	oondent	subspell ts15201	Number of subspell Type of vocational tra	ining	
ID variables need	ded to identify a	single row		ts1511m	Starting month Trainir	ng episode	
ID_t spell subspell				—— ts1511y — ts1512m	Starting year Training episode End month Training episode		
Other ID variable	es useful for linka	age		ts1512y	End year Training epis	ode	
wave splink				ts15215	Company size of train company	ing	
Number of varia	bles / number of	f rows in file					
164 / 133,5	546						
Contains data fro	om waves						
1 2 3 12 13 14	4 5 6 15 16	5 7 8 9	10 11				
Exemplary data	snapshot						
ID_t	spell	subspell	ts1511m	ts1511y	ts1512m	ts1512y	
7028397	2	3	9	2010	5	2013	
7018057	1	3 3	10 10	2010	5	2013	
7013700	2	э 1	10	2010	2	2013	
1010100	-	-	10	2010	Ζ.	2011	

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

- tertiary education at universities (including colleges of education, theology, and art and music), universities of applied sciences, Berufsakademien/cooperative state universities, colleges of public administration). Note: Up to three subjects (majors and minors) are recorded.
- doctoral or postdoctoral studies
- 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

Data Structure

accredited training courses to receive licenses

In case of higher education study, new episodes are generated if

- a subject changes over the course of studies, or
- the intended degree changes over the course of studies (e. g., from master's degree to state examination), or
- the higher education institution changes.

If a higher education episode follows immediately a preceding higher education episode, interviewees are asked whether the type of degree was changed (tg24146), whether different subjects were chosen (tg24159) or whether the respondent moved to another higher education institution (tg24121). New information on the intended degree (ts15221), subjects (variables beginning with tg2416, tg2417), and higher education institution were only collected when the aforementioned questions were answered with *yes*. In case of a negative answer, the variable ts15221 (intended degree) takes the value of the preceding episode while the variables containing information on the subjects take the value -29 (value from the last sub-episode). The information of the preceding episode was integrated into the service variables tg24162_g1, tg24165_g1, tg24168_g1, tg24170_g1-tg24170_g5, tg24173_g1tg24173_g5 and tg24176_g1-tg24176_g5 (see section section 5.1.1).

Information on the subjects, intended degree, and the higher education institution of the first study episode in winter term 2010/2011 was collected in the initial questionnaire, which was mainly administered as a written survey and partly integrated into the first telephone interview. This data can be found in pTargetCATI in the variables tg0400* (information on the subjects), tg01003_g1 (type of higher education institution), tg15207_g1R, tg15207_g2R (location of the higher education institution), and tg02001* (intended degree). The information was not newly collected in the first telephone interview but was integrated into the service variables tg2417* and tg01003_ha (see section section 5.1.1). The variable h_aktstu in spVocTrain indicates which episode refers to the first study episode in winter term 2010/2011 (h_aktstu==1 "Episode is 1st study episode WT 2010 (start of study)").

In the telephone interview following wave 1 (i. e., in wave 3, 5 or 7, depending on panel participation), the vocational education history from winter term 2010/2011 onwards was newly collected with an improved survey instrument. This has led to duplicate and/or right-censored episodes in the dataset spVocTrain. In order to deal with those episodes, the variable $t \times 20100$ 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.

Example 31 (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
```

4.5.32 StudyStates

Description Exemplary variables Data on state of studies derived from spVoc-ID_t ID target Train wave Wave tx24000 Recommendation: use File structure person (less than three long format: 1 row = 1 respondent in 1 wave episode per wave) tx24001 chronological order of ID variables needed to identify a single row interviews ID_t wave or ID_t tx24001 tx24021 Change of episode number tx24022 Episode number (splink) Other ID variables useful for linkage tx24100 Status of studies ID_t tx24022 (completed,on-going) tx15318 Successful completion of Number of variables / number of rows in file training 45 / 218,517 tx15317 Vocational qualification Contains data from waves tx15310 Status of study interruption tx24190 Type of study interruption 5 6 7 8 9 10 11 1 2 3 4 12 13 14 15 Exemplary data snapshot ID_t tx24022 tx24100 tx15318 tx15317 tx15310 wave 7018735 1 240001 ongoing, no compl. ongoing not reached no interrupt. 7018735 2 240001 ongoing, no compl. ongoing not reached no interrupt. 7018735 240001 ongoing, no compl. no interrupt. 3 ongoing not reached 7018735 4 240001 ongoing, no compl. ongoing not reached no interrupt. 7018735 5 240001 ongoing, no compl. ongoing not reached no interrupt. 7018735 6 240003 completed + ongoing Bachelor no interrupt. yes 7018735 7 240003 completed + ongoing Bachelor no interrupt. yes 8 240003 completed + ongoing Bachelor 7018735 yes no interrupt. completed + ongoing 7018735 9 240004 yes Bachelor no interrupt. 7018735 10 240004 all completed Master yes no interrupt. 11 no further episodes 7018735 -21 no interrupt. . . 7018735 12 -21 no further episodes no interrupt. . 7018735 13 -21 no further episodes no interrupt. . 7018735 -21 14 no further episodes no interrupt. 7018735 15 -21 no further episodes no interrupt.

The file StudyStates contains all target persons for each wave as long as they have not dropped out. As soon as a person drops out, it will not be part of the dataset since that certain wave.

« go back to overview

For each respondent in each wave, StudyStates contains information on status of studies/tertial education (tx24100), which vocational qualification is achieved (tx15317), which subjects were choosen or switched. There is also data on wether a person continued its studies at a different educational institution (tx24011) as well as information on study-breaks (tx15310, tx24190).

This dataset boils down information from spVocTrain to meet the data structure of longformat files such as pTargetCATI – one line for each person per wave. This procedure inevitable goes along with information loss. As some target persons' data show multiple educational episodes at the same time, defining a proper status of studies, for instance, is virtually impossible to achieve without sound assumptions. Therefore we introduced an indicator that recommends the usage of a person within the dataset (tx24000). The data is recommended to use as long as the person has less than three episodes in any wave.

This dataset is still a kind of beta version of the desired outcome, therefore we suggest to use this dataset to get a quick insight on data concering study states. Data will be polished and will be more usuable in the next release!

Example 32 (Stata): Working with StudyStates

```
*** 1. enriching StudyStates with episode data from spVocTrain ***
** open spVocTrain file
use "${datapath}/SC5_spVocTrain_D_${version}.dta" , clear
** only keep full or harmonized episodes and save file temporarily
keep if subspell == 0
tempfile spvoc
save "`spvoc'", replace
** open StudyStates file
use "${datapath}/SC5_StudyStates_D_${version}.dta" , clear
** rename tx24022 to splink and keep only valid episodes
rename tx24022 splink
keep if splink > 0
** merging StudyStates with spVocTrain, only keeping desired variables
merge m:1 ID_t splink using "`spvoc'", keep(master matched) nogenerate keepusing(
 tg24203 tg24205 tg24162_g1 tg24165_g1 tg24168_g1)
*** 2. merging StudyStates with pTarget-data ***
** open StudyStates file
use "${datapath}/SC5_StudyStates_D_${version}.dta" , clear
** merging StudyStates with pTargetCATI, only keeping desired variables
merge 1:1 ID_t wave using "${datapath}/SC5_pTargetCATI_D_${version}.dta", keep(master
  matched) nogenerate keepusing(t406000 t731351_g1 t31300a t34006k tg2411a
 t66003a_g1 t531260)
```

```
** merging data with pTargetCAWI, only keeping desired variables
merge 1:1 ID_t wave using "${datapath}/SC5_pTargetCAWI_D_${version}", keep(master
    matched) nogenerate keepusing(t291501 t291502 t321104 t66007a t513051 tg54112
    t30300b)
```

4.5.33 Weights



Weighting variables (starting with w_) are included in the Weights dataset. Also, you find cluster (ID_cl) and stratification (stratum) identifiers here. ID_i resembles university at sampling, which took part prior to wave 1. 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.

« go back to overview

Example 33 (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
```

4.5.34 xEcoCAPI

Description **Exemplary** variables additional competencies for students of eco-ID target ID_t nomics and business administration wave Wave tx80921 Participation status: File structure Economics-subsample wide format: 1 row = 1 student testm Test: Survey day (month) testy Test: Survey day (year) ID variables needed to identify a single row bas7mar1_c Economic competence: ID_t marketing 1 bas7_sc1 Economic competence: WLE Other ID variables useful for linkage bas7_sc2 Economic competence: wave ID_int SE(WLE) Number of variables / number of rows in file ID_int Interviewer: ID tg90308 Number semesters 136 / 600 economics Contains data from waves tg24160_g2 Subject group subject 1 (destatis 2010/11) 5 6 7 8 9 10 11 tx80200 Interview: number of all contact attempts Exemplary data snapshot ID_t tx80921 bas7_sc1 bas7_sc2 tg90308 tg24160_g2 tx80200 wave 7006150 7 6 0.12031 0.38248 7 3 2 7009130 7 6 0.37047 0.39715 7 3 3 7009996 6 0.48922 0.48893 7 3 2 7 0.61990 8 3 7012529 7 6 0.77208 3 7014554 7 6 0.02677 0.40591 7 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 specificially 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 $t \times 80921$. Additional methods data – like number of contact tries ($t \times 80200$) and reasons for item-nonresponse in testing (e.g., $t \times 80411$) – are available as well.

« go back to overview

CAPI data are basically focussing on the student's area of studies (e. g., $tg24160_g2$). For more information see Lauterbach (2015).

Example 34 (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 in now available only in waves which have
** surveyed the topic
tab wave bas7_sc1
```

4.5.35 xInstitution

				« go b	ack to overview
Description			Exemplary var	riables	
context informa	ation about the institu	ition	ID_i	Institution ID	
File structure			tg04001_g7	Subject group	WT 2010 (for
wide format: 1	row = 1 area of studie	s in 1 insti-		merging with c	ontext data)
tution			tg91102_R	HEI region: BIK	-region type
	a tile attractions		lg92104_0	excellence 200	6 or 2007
ID variables needed t	to identity a single row		— tg92301_0	HEI: Funding b	ody
ID_I tg04001_g	/		_tg92601_R	HEI: Students 2	010 total
Other ID variables us	eful for linkage			(aggr. Tercent.	universities)
none			tg93204_0	SG: Students 2	010: male
Number of variables	/ number of rows in file		tg93205_0	SG: Students 2	er professor
127 / 3 510	-				
12, , 3,310					
Contains data from w	vaves				
Exemplary data snap	shot	+-02104 0		00001 0	+-02001 D
1002112	tg04001_g7	tg92104_0	τg	92301_0	tg92601_R
1002095	1	1		1	3
1002222	10	1		1	3
1002222	3	1		1	3
1002042	5	1		1	3

Data file xInstitution contains context data (e.g., size of the institution, regional unemployment rate) for all 413 higher education institutions which were listed in the codebook of the Federal Statistical Office in 2010/2011. However, higher education institutions with different locations are only considered once with combined information for all locations (for detailed information, see Weber, 2014).

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 also note that the context information up to now has not been updated and refers to most recent information available in 2010.

Example 35 (Stata): Working with xInstitution (find R example here)

** open datafile

```
use ${datapath}/SC5_pTargetCATI_0_${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 iformation
keep if seq == max
\star\star 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)
\star\star assuming that the less students at university the more intensive the support by
the
\star\star 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
```

4.5.36 xPlausibleValues

				U			
Description				Exemplary variables			
Plausible Va	lues of competer	nce data	ID_t	ID target			
File structure			wave_w1	Row contains da	ta from		
				wave 1 (2010/20	011		
wide format	: 1 row = 1 respo	ondent		(CATI+competer	icies))		
ID variables need	ded to identify a single	row	wave_w5	Row contains da	ta from		
				wave 5			
ID_t			wave_w1	2 Row contains da	ta from		
Other ID variable	es useful for linkage			wave 12			
			mas1_pv	1 Math: cross-sec	tional		
wave_w*				plausible value 1	L		
Number of varial	bles / number of rows	in file	mas1_pv2	2 Math: cross-sec	tional		
114 / 11 72	20			plausible value 2	2		
114 / 11,/5	99		mas1_pv	10 Math: cross-sec	tional		
Contains data fro	om waves			plausible value 1	10		
			mas1_pv:	1u Math: longitudii	nal		
1 2 3	4 5 6 /			plausible value 1	L		
12 13 14							
Exemplary data s	snapshot						
ID_t	wave_w1	mas1_pv1	mas1_pv2	mas1_pv10	mas1_pv1u		
7013717	1	1.10363	0.55658	1.13098	1.28266		
7009412	1	0.47345	0.42787	1.14043	1.96270		
7010495	1	0.96598	1.93949	0.37123	1.59425		
7005346	1	1.21510	0.82586	1.39770	2.41945		
7010340	1	1.03923	0.99005	0.71114	2.68605		

« go back to overview

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 et al., 2020) and on our website:

 \rightarrow www.neps-data.de > Data Center > Overview and Assistance > Plausible Values

Example 36 (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!
```

4.5.37 xTargetCompetencies

Exemplary var	riables	
ID_t	ID target	
wave_w1	Row contains data	from
	wave 1 (2010/2011	
	(CATI+competencie	s))
wave_w5	Row contains data	from
	wave 5	
mas1r092_c	Mathematical competence:	
	Item 2	
mas1_sc1	Mathematical com WI F	petence:
mas1_sc2	Mathematical com	petence:
	SE(WLE)	
res1_sc1	Reading competen	ce: WLE
res1_sc2	Reading competen	ce:
	SE(WLE)	
rsci0051_c	Reading speed: Ite	m 51
rss1_sc3	Reading speed: Sur	n
5 ma	s1_sc1	mas1_sc2
1 1	.56820	0.65016
1 0	88910	0.55859
1 0	.21208	0.77751
1 0	.36226	0.55137
	wave_w1 wave_w5 mas1r092_c mas1_sc1 mas1_sc2 res1_sc1 res1_sc2 rsci0051_c rss1_sc3 ics3_sc1 ics3_sc2	wave_w1 Row contains data is wave 1 (2010/2011 (CATI+competencie) wave_w5 Row contains data is wave 5 mas1r092_c Mathematical completencie) mas1r092_c Mathematical completencie) mas1_sc1 Mathematical completencie) mas1_sc2 Mathematical completencie) res1_sc1 Reading competencie) res1_sc2 Reading competencie) res1_sc2 Reading speed: liter rss1_sc3 Reading speed: Sur ics3_sc1 ICT-Literacy WLE ics3_sc2 ICT-Literacy SE of W 1 0.88910 1 0.36226

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 in wave 1 competence tests were conducted in paper-and-pencil mode in groups at the participating higher education institutions. The wave 5 test included a mode experiment with an individual online test on the one hand and three different modes applied in a group setting (conventional paper-based assessment, paper-based assessment with digital pens, and computer-based assessment). Please also note that data from the web-based assessment are not available yet.

Due to overlaps in survey periods in wave 12, data was also collected from target persons who did not participate in the last three CATI waves (including CATI in wave 12) and who are usually treated as final drop-outs.

```
Example 37 (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 Special Types of Variables

5.1.1 Service Variables

field 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 tg2416*_g1 and tg2417* contain information on the respective field of study, thus the variables tg24162_g1, tg24165_g1, tg24168_g1, tg24170_g1-tg24170_g5, tg24173_g1-tg24173_g5, and tg24176_g1-tg24176_g5 provide complete subject information for all study episodes. Working with the service variables is recommended.

- type of higher education institution The variable tg01003_g1 (type of higher education institution, two levels) is originally a part of the first wave recruitment information contained in dataset pTargetCATI. The variable ts15201 (type of vocational training program, seventeen 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 higher education institution) provides an aggregated version of ts15201 in spVocTrain partly using information from tg01003_g1 for first wave spells, as seen in table 9.
- intended vocational qualification Because of a programming error, variable ts15221 (Intended vocational qualification) misleadingly contains information on the achieved vocational qualification in waves 9 and 10 for some respondents. To correct this mistake, the variable ts15221_g1 (Intended vocational qualification, revised) was introduced. This variable contains the correct information on the intended vocational qualification for all target persons in the data file and for all subspells of a vocational episode. The information on the intended vocational qualification for wave 1 was collected in the initial questionnaire and is stored in variable tg02001 in pTargetCATI. This variable was used to provide the information in ts15221_g1 in spVocTrain for wave 1.

	tg01003_ha/tg01003_g1		ts15201
1	University of applied sciences (incl. Beruf- sakademie/cooperative state university)	7 8 9	Degree course at a Berufsakademie/cooperative state university Degree course at a college of public administration Degree course at a university of applied sciences (not a college of public administration)
2	University	10	Degree course at a university, including college of ed- ucation, art college, music college

Table 9: Harmonization of type of higher education institution

vocational education history In the telephone interview following wave 1 (i. e., in wave 3, 5, or 7, depending on panel participation), the vocational education history from winter term 2010/2011 onwards was newly collected with an improved survey instrument. 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.

5.1.2 Auxiliary variables

Additionally to the reported information, auxiliary variables are generated automatically during the course of an interview. They are used to manage the interview process and to ensure that the questions are addressed adequately for each interviewee. Part of this supplementary stored information is useful when analysing the data. Therefore some of these automatically generated variables are released in different data sets. Auxiliary variables can be identified by their variable label, which begins with "Auxiliary variable:" and indicates, that it is not a recorded but an automatically generated information about the target person.

5.1.3 Version variables

It rarely happens that errors in the programming of the questionnaire appear during the field time, which could jeopardize the correct execution of the interview. In these cases, error corrections in the field are required. This information is stored in so-called version variables, so that it is documented later on in the data which target persons have received which application version of the instrument. One can identify such variables by their variable name starting
with Version_. A description of the error correction can be called up using the command infoquery var in Stata.

5.1.4 Preload Variables

In order to disburden the process of the online survey, information from prior waves is used to guide through the survey. It is important to note that only information from CATI waves is used to generate this preload information. Since some information in the online waves is only updated when a target person reports changes since the last CATI wave, some of these preloaded information is released in the data. You can identify the preloaded information by its label starting with "Preload:", indicating that this variable contains the information known at the time of the last CATI survey.

5.2 Coding field of study

5.2.1 Recruitment

- **data collection** Information on field of study of the first study program in winter semester 2010/11 was collected mainly in PAPI and sometimes also in CATI mode (for information on sampling in SC5, see Aßmann et al., 2011, and Zinn et al., 2017). The data of the PAPI questionnaires were entered by the data collecting institute (infas) and delivered to NEPS. Information on the field of study was delivered to NEPS as original string variable.
- **coding** Coding of field of study was done 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 because the classification scheme used changed between recruitment and first wave data collection: sampling was based on the classification of 2009/10 while the coding of recruitment information was based on the classification of 2010/11 (see Statistisches Bundesamt, 2011).

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 field of study is based on the Federal Statistical Office (Destatis) for the winter semester 2010/11 (Statistisches Bundesamt, 2011). Coding decisions can differ from Destatis recommendations for coding degree programs into fields 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 field of study has been recorded using lists – in CATI as well as in online surveys. In cases where interviewers were unable to fit a respondents answer into the respective list, the field of study has 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 respondents' 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, same subjects of study were assigned to *different* codes within the list. In other cases, multiple subjects were listed under the same code. The idea behind this was for the added subjects within the same code to serve as covariates, so interviewers could classify the respondents' answer into the *right* code in the list. 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 for collecting the field of study 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 fields of study in CATI and online mode.

- **coding** Coding of open responses regarding field 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 field of study largely follows the Destatis 2010/11 classification of fields of study.
- **derivation of SUF-variables** In the Scientific Use File, several alternative variables containing information on the field of study are offered. Variables with the suffix _g1R and _g2 contain aggregations of subjects according to the 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 (Statistisches Bundesamt, 2011). All derivations are based on a transcoding table provided by the Federal Statistical Office.

5.3 Coding of Higher Education Institutions

- **data collection** In the initial questionnaire the information on Higher Education Institutions (HEI) was collected as an open string variable. In the following CATI and CAWI panel waves, new information on HEI was collected using lists. These lists are based on the destatis classification valid at the time of data collection (see further description below) and were extended to include coded open answers from further waves.
- **coding** The open answers from the initial questionnaire were coded by the NEPS department *From Higher Education to the Labour market* using the Destatis classification of winter term 2010/2011. In CATI and CAWI panel waves, new open answers are coded using the extended list for data collection.
- **classification used** Data collection and coding of HEI are based on the Destatis classifications of HEI since winter term 2010/2011. To match the current higher education area, the list used for data collection was updated annually according to the changes as presented in the Destatis classification. These changes include:
 - Adding new HEI
 - Deleting existing HEI
 - Integration of one HEI into another
 - Division in different locations/sites of HEI
 - Renaming of HEI
 - Renaming and changing of type of HEI
 - Merging of locations
 - Fusion of HEI

If the HEI codes stay the same, these changes have no consequences for data collection and coding. But there are modifications that result in new codes for existing institutions. Hence, it is possible that the institution ID (ID_i) in the data differs between respondents or between spells of one respondent, even though the respondents attend the same institution. Since winter term 2010/2011, these changes are listed in table table 10 (source: Destatis).

ID_i before	ID_i after	changed in
1003035	1002062	winter term 2013/14
1003089	1002366	winter term 2015/16
1003079	1002260	winter term 2016/17
1003142	1002987	winter term 2017/18
1003105	1002207	winter term 2018/19
1003137	1002165	winter term 2018/19

Table 10: Changes in HEI codes during survey

5.4 Special features of interruption episodes in spVocTrain

There are no sensible harmonization rules for interruption episodes. In the data, interruption episodes are filed in wide format (first interruption _w1, second interruption _w2, and third interruption _w3). The variables for a first / second / third interruption are being harmonized. However, these do not necessarily correspond between subspells. A second (persistent) break in the first wave may correspond to a first break in the second wave. For example:

- **subspell 1** two interruptions; the first interruption episode is completed; the second interruption episode is right-censored because it continues at the point of the survey
- subspell 2 the continuing interruption is completed and stored in the variables for a first interruption

In addition, there was no range definition for the first interruption. The stimulus in the question was designed to report only interruptions since the last interview date, but if one is asked about the time of the first interruption in a course of study, then a wide range of information is possible. The target person could even think that the start time of a first of several interruption episodes is meant (beginning of the first interruption in the first subspell). This is why there is a small percentage in the data that reports a start date of the interruption before the interview date.

5.5 Teacher Education Students and Teachers

The sample of Starting Cohort 5 includes an oversample of teacher education students (see section 2.2 for further information). Since wave 8, the survey program is supplemented with specific questions for all (prospective) teachers – whether or not they belong to the oversample or basis sample. The participants who belong to the oversample can be identified by the variable $t \times 80121$ (*Sample: oversample of teacher education students*) stored in CohortProfile. This information is important for analyzing selected waves. Because of funding issues the

oversample could not be invited to take part in wave 7 but remained in the sample as temporary dropouts. In wave 14 the oversample did not answer the questions regarding *job requirements* (pTargetCAWI: tg781*; tg782*; tg783*; tg784*) due to survey methodological reasons.

In the first panel wave, information whether a participant belongs to the group of teacher education students was collected in two different surveys – the initial paper and pencil questionnaire and the first telephone interview. Because of different question wording between the two surveys, this information is stored in different variables in pTargetCATI. Generally speaking there are two types of variables, described below, that help identify if a participants course of study leads to a teaching degree. First, there are those variables containing information about whether someone's degree is a teaching degree (intended teaching degree) and second, there are variables containing information about the type of teaching degree that is intended (type of teaching degree). Some of these variables contain only information about the study course of the first panel wave and others are applicable for all information about teaching degrees, irrespective of the panel wave it was collected. That helps to identify survey participants who entered the study with a non-teaching degree course of study and changed into a teaching degree course of study later on.

Identify intended teaching degree Variable tg02001 (Intended degree WT 2010 (PAPI questionnaire)) contains information collected in the initial questionnaire. The variable distinguishes several degree types, including those that lead to a teaching degree in German teacher education (e.g., bachelor's degree, state examination). So-called polyvalent Bachelor courses with the option of specializing in teacher education are subsumed under the category "Bachelor (not in teaching)". In variable tg02001_g1 (Intended degree WT 2010 (incl. polyvalent Bachelor; PAPI questionnaire)), this specific potential pathway to a teaching degree is recorded separately in the category "Polyvalent Bachelor with a teaching option".

In the data of the telephone interviews in spVocTrain, teacher education students can be identified by the variable ts15221_g1 (Intended vocational qualification, revised) in combination with tg24201 (Intended teaching degree). In contrast to the initial PAPI questionnaire, participants who report a higher education episode and want to obtain a non-teaching degree are asked whether they are studying with the aim of becoming a teacher. This new question was introduced because it is possible to be enrolled in a non-teaching degree course and to decide later for a teaching degree. To identify teacher education students in the first study episode of the winter term 2010/2011, additional information on the type of episode (variable h_aktstu) has to be used. For ease of use, a new service variable tg02001_ha (Intended degree WT 2010 (start of study; CATI and PAPI questionnaire)) was generated and introduced in pTargetCATI (see below).

tg24201_g1 (Intended teaching degree WT 2010 (start of study; CATI)) contains information on whether the first study program in winter term 2010 was started with the aim of becoming a teacher. The information comes from the first telephone interview (variable tg24201 if h_aktstu=1 and wave=1). Because of different question wording, the variable differs from tg02001 (*Intended degree WT 2010 (PAPI questionnaire*)). tg24201_g1 is part of pTargetCATI.

tg02001_ha (Intended degree WT 2010 (start of study; CATI and PAPI questionnaire)) combines information on the intended degree collected in the first PAPI questionnaire and the first telephone interview. It refers to the first study program in winter term 2010 and updates the variable tg02001 (Intended degree WT 2010 (PAPI questionnaire)) with information whether a teaching degree is intended, collected during the first telephone interview (tg24201). tg02001_ha is part of pTargetCATI.

Coding of type of (intended) teaching degree Data on the type of intended teaching degree were collected using an open-ended question. The coding of the answers is based on the classification of teaching careers proposed by the Standing Conference of the Ministers of Education and Cultural Affairs of the Länder in the Federal Republic of Germany (Kultusministerkonferenz/KMK) but does not distinguish teacher education programs encompassing several levels. In case that more than one type of teaching degree or a program that span several levels was mentioned without specifying a main focus, the answer was coded under the highest level.

The corresponding variables are tg51420_g1 (Type of intended teaching degree (differentiated; CAWI)) in pTargetCAWI and tg24202_g1 (Type of intended teaching degree (differentiated; CATI)), tg24202_ha (Type of intended teaching degree WT 2010 (start of study; CATI)), and tg03001_g2 (Type of intended teaching degree WT 2010 (differentiated; PAPI questionnaire)) – all included in pTargetCATI. tg24202_g2 (Type of intended teaching degree WT 2010 (start of study; CATI)) contains information on the type of the intended teaching degree and refers to the first study program in winter term 2010. The information comes from the first telephone interview (variable tg24202_g1, if h_aktstu=1 and wave=1). tg24202_g2 is part of pTargetCATI.

Important notes on auxiliary variables The auxiliary variables tg60011 (wave 8), tg60014 (wave 9), tg60015 (wave 10), tg60016 (wave 13), have been generated to navigate through the survey. Please use these variables only to get a first overview of the data. Use the original episode files for analyses!

tg60012 and tg60013 do not only include information on the phase of teacher education a participant is pursuing or has completed but also on being employed as teacher. In addition, information on intentions are taken into account: Participants who have completed the second phase of teacher training (preparatory service) but are not yet employed as a teacher are asked whether they intend to work as a teacher. Respondents who have completed the first phase of teacher education at universities or equivalent institutions but have not yet started the second phase are asked whether they want to complete preparatory service. These two auxiliary variables are available from wave 11 onwards and are part of either pTargetCATI (tg60013) or pTargetCAWI (tg60012). Since variable tg60013 is used for guidance during the interview, there have been adjustments in this variable in wave 13 against wave 12 to ensure a stricter filtering into certain teacher-related questions. It is now achieved that target persons who only have a qualification in a teacher-related bachelor's degree, but no longer intend to study a teacher-related master's degree, will no longer receive the teaching-related questions (as was still the case in wave 12). From wave 16 on there was a new answering category added so that variable $\pm g60013$ contains now information about actual interrupted teaching employment episodes ("6 = interrupted employment as a teacher (e.g. due to parental leave)"). Survey participants with this status category are presented a reduced teacher context questionnaire. All these successive adjustments result in the existence of three versions of this variable ($\pm g60013_v1$ in wave 12 vs. $\pm g60013_v2$ in wave 13 and wave 15 vs. $\pm g60013$ from wave 16 onwards) in the dataset pTargetCAWI. Since wave 14 the dataset pTargetCAWI contains an additional auxiliary variable to state the actual phase of teacher education of each participant ($\pm g60017$). The only difference to variable $\pm g60012$ is that those participants who denied or revoked a teacher (education) context later in the questionnaire are now transferred to category "0 = no teaching reference or status unknown".

5.6 Wave Specific Issues

The 2018 online survey (wave 14) experimented with various incentives. 50% of the sample was still offered to take part in a lottery. 25% of the sample was offered a cash incentive of EUR 10 and 25% of the sample had the choice between participating in the lottery or a cash incentive of EUR 10. The distribution across the three incentive groups was random. The assignment to the different incentive groups is documented in the variable tg59000 (*Auxiliary Variable: Assignment B139 incentive group*) in pTargetCAWI. For further information see the respective field report at

 \rightarrow www.neps-data.de > Data Center > Data and Documentation > Starting Cohort First-Year Students > Documentation



- Aßmann, C., Steinhauer, H. W., Kiesl, H., Koch, S., Schönberger, B., Müller-Kuller, A., Rohwer, G., Rässler, S., & Blossfeld, H.-P. (2011). Sampling designs of the National Educational Panel Study: Challenges and Solutions (H.-P. Blossfeld, H.-G. Roßbach, & J. von Maurice, Eds.). Education as a Lifelong Process: The German National Educational Panel Study (NEPS), 14, 51–65. https://doi.org/10.1007/s11618-011-0181-8
- Blossfeld, H.-P., & Roßbach, H.-G. (Eds.). (2019). Education as a lifelong process: The German National Educational Panel Study (NEPS). Edition ZfE (2nd ed.). Springer VS. https://doi. org/10.1007/978-3-658-23162-0
- 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" (Nichttraditionelle Studierende) (DZHW: Data Manual). DZHW - Deutsches Zentrum für Hochschul- und Wissenschaftsforschung GmbH.
- FDZ-LIfBi. (2022). Data Manual NEPS Starting Cohort 5– First-Year Students, From Higher Education to the Labor Market, Scientific Use File Version 16.0.0. Bamberg, Germany, Leibniz Institute for Educational Trajectories, National Educational Panel Study.
- Hess, D., Steinwede, A., & Schneider, B. (2012). *Erhebung von retrospektiven Längsschnittdaten* - *Prüfmodul*. Bonn, infas Institut für angewandte Sozialwissenschaft GmbH.
- Künster, R. (2015a). Startkohorte 6: Erwachsene (SC6) Datenversion 5.0.0. Technical Report 1: Edition und Korrektur der Lebensverlaufsdaten. Leibniz-Institut für Bildungsverläufe, Nationales Bildungspanel. Bamberg, Germany.
- Künster, R. (2015b). Startkohorte 6: Erwachsene (SC6) Datenversion 5.1.0. Technical Report: Korrektur der Lebensverlaufsdaten. Leibniz-Institut für Bildungsverläufe, Nationales Bildungspanel. Bamberg, Germany.
- 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.
- Matthes, B., Reimer, M., & Künster, R. (2005). TrueTales ein neues Instrument zur Erhebung von Längsschnittdaten. In Arbeitsbericht 2 des Projektes "Frühe Karrieren und Familiengründung: Lebensverläufe der Geburtskohorte 1971 in Ost- und Westdeutschland".
- Matthes, B., Reimer, M., & Künster, R. (2007). Techniken und Werkzeuge zur Unterstützung der Erinnerungsarbeit bei der computergestützten Erhebung retrospektiver Längsschnittdaten. Methoden, Daten, Analysen – Zeitschrift für Empirische Sozialforschung, 1(1), 69–92.

- NEPS Network. (2022-a). National Educational Panel Study, Scientific Use File of Starting Cohort First-Year Students. Leibniz Institute for Educational Trajectories (LIfBi), Bamberg. https: //doi.org/10.5157/NEPS:SC5:16.0.0
- NEPS Network. (2022-b). Starting Cohort 5: First-Year Students (SC5), Wave 16, Questionnaires (SUF Version 16.0.0). 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.
- Ruland, M., Drasch, K., Künster, R., Matthes, B., & Steinwede, A. (2016). Data-Revision Module
 A Beneficial Tool to Support Autobiographical Memory in Life-Course Studies. In H.-P.
 Blossfeld, J. von Maurice, M. Bayer, & J. Skopek (Eds.), *Methodological Issues of Longitudinal Surveys. The Example of the National Educational Panel Study* (pp. 367–384).
 Springer VS.
- 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. 10).
 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.
- Statistisches Bundesamt (Ed.). (2011). Bildung und Kultur. Fachserie 11 Reihe 4.1 Studierende an Hochschulen. Wintersemester 2010/2011.
- Steinwede, J., & Aust, F. (2012). *Methodenbericht, NEPS Startkohorte 5 CATI-Haupterhebung Herbst 2010, B52.* Bonn, Germany, infas.
- Weber, A. (2014). Data Manual: Starting Cohort 5 Context Data. DZHW. Hannover.
- 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.
- 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.



B.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 38 (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 39 (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 40 (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 = "Teilepisodennummer", ... 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 41 (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 42 (R): Working with Biography

'** import the data file'
Biography =

Example 43 (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
```

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

Example 44 (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",
 # 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
)
'** see that this only added information to the rows corresponding to spSchool'
cbind(addmargins(table(Education$tx28100, Education$source)))
```

Example 45 (R): Working with MethodsCATI

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



```
'** 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 summirized in groups
for (i in 1:length(MethodsCompetencies$ID_tg)) {
 if(!is.na(MethodsCompetencies$ID_tg[i]) & MethodsCompetencies$ID_tg[i] < 0){</pre>
   MethodsCompetencies$groupsize[i] = NA
   #sets all observations to NA for which ID_tg < 0 (here -55 and -54)</pre>
 }
}
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:%S")))
    #assigns the values from tx80603 and tx80604 in time format to them
}
MethodsCompetencies$duration =
        \label{eq:methodsCompetencies} MethodsCompetencies \\ \$tx80604\_t - MethodsCompetencies \\ \$tx80603\_t \\ \end{cases}
#creates a new variable "duration", subtracting start time from end time
summary(MethodsCompetencies$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 47 (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))
'\star\star 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 48 (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.
' \star\star 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 49 (R): Working with pTargetMicrom

```
'** open pTargetMicrom datafile. Note that this data file is only available OnSite!'
pTargetMicrom = read.dta13("SC6_pTargetMicrom_0_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_0_version.dta", convert.factors = T)
pTargetMicrom = merge(CohortProfile, pTargetMicrom, by = c("ID_t", "wave"), all =
TRUE)
```

Example 50 (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 51 (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)))})
```

Example 52 (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 53 (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$subspell == 0)
'** open the Biography data file'
Biography = read.dta13("SC5_Biography_D_version.dta", convert.factors = T)
'** merge the spEmp to Biography'
#After merging, Stata merge has one variable more than R, because in Stata
#a merge indicator is produced during the merging process and in R isn't.
#Since we need a merge indicator here, the merge command has to be extended:
Biography = transform(merge(
 x = cbind(Biography, source = "master"),
 #x contains the Biography data set plus one extra column "source",
 #where source = "master"
 y = cbind(spEmp,source = "using"),
 #y contains the spEmp data set plus one extra column "source",
 #where source = "using"
 all.x = TRUE, by = c("ID_t", "splink")),
 #merges x and y by ID_t and splink
 source = ifelse(!is.na(source.x) & !is.na(source.y), "both",
 #in the merged dataset, source = "both" if the observations is in x AND in y
              ifelse(!is.na(source.x), "master", "using")),
              #otherwise, source = "master" if the obs. is only in x
              #and source = "using" if the obs. is only in y
 source.x = NULL.
 source.y = NULL
 #the columns "source" in x and y are deleted
)
#Since the variables wave and spms are in both data sets, Biography AND spEmp
#check before merging by: intersect(names(Biography), names(spEmp))
#and since the variables are not part of the merging process,
#both versions are contained in the new merged data set as
#wave.x/wave.y and spms.x/spms.y
#For each variable, one of the versions can be dropped by:
Biography$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 54 (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 55 (R): Working with spFurtherEdu2

```
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"</pre>
#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.
<code>'**To</code> 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 56 (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 \boldsymbol{x} and \boldsymbol{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 57 (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 58 (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 59 (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 60 (R): Working with spPartner

Example 61 (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 62 (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 63 (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$target_bdate[i]) {</pre>
     spSibling$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 identificator'
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 64 (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 65 (R): Working with spVocExtExam

```
'** 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 66 (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.
```

```
Appendix
```

```
#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 67 (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",
 #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 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 68 (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, 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 69 (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$t514008)),]
#remove observations with other missings in t514008
pTargetCATI = subset(pTargetCATI, !(t514008 == "Don't know"
                                  | t514008 == "Refused"
                                  | t514008 == "Does not apply"
```
```
t514008 == "Missing by design"))
'** some respondents reported satisfaction with studies in 7th and in 9th waves
** to keep the latest information, create a seq and a max variables'
pTargetCATI = within(pTargetCATI, {seq = ave(ID_t, ID_t, FUN = seq_along)})
pTargetCATI = within(pTargetCATI, {max = ave(ID_t, ID_t, FUN = length)})
'** only keep the latest reported iformation'
pTargetCATI =
        subset(pTargetCATI, pTargetCATI$seq == pTargetCATI$max)
'** only keep the variables relevant for the merge and the analysis'
pTargetCATI =
        subset(pTargetCATI, select = c("ID_t", "ID_i", "tg04001_g7", "t514008"))
'** merge two variables from xInstitution'
#open datafile xInstitution
xInstitution = read.dta13("SC5_xInstitution_0_version.dta", convert.factors = T)
#merge xInstitution to pTargetCATI
pTargetCATI =
 merge(x = pTargetCATI,
             y = xInstitution[,c("ID_i", "g04001_g7", "tg92601_R", "tg92104_0")],
              by = c("ID_i", "g04001_g7"), all.x = TRUE)
<code>'**</code> assuming that the less students at university the more intensive the support by
** the university staff per student and the more satisfied are students with their
** studies tabulate Satisfaction with studies by Students 2010 total
** note that the following analysis is feasible in both, RemoteNEPS and Onsite'
cbind(addmargins(table(pTargetCATI$t514008, pTargetCATI$tg92601_R)))
cbind(addmargins(prop.table(table(pTargetCATI$t514008, pTargetCATI$tg92601_R))))
' \star \star assuming that students at excellence universities are more satisfied with
** their studies, tabulate the distribution of satisfaction by tg92104_0
** note that the following analysis is only feasible in the Onsite version of SUF,
** since the variable tg92104_0 is anonymized in RemoteNEPS'
cbind(addmargins(table(pTargetCATI$t514008, pTargetCATI$tg92104_0)))
cbind(addmargins(prop.table(table(pTargetCATI$t514008, pTargetCATI$tg92104_0))))
```

Example 70 (R): Working with xTargetCompetencies

```
#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$wave_w7)
'** to work with competence data, you might want to merge it to CohortProfile.
** if you want to keep the panel logic (and not only add all competencies
** to every wave), you need a mergeable wave variable in xTargetCompetencies.
\star\star here, we focus on math competencies, that have been tested in wave 1.'
xTargetCompetencies$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(xTargeCompetencies, 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)
```

B.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 16.0.0 ** (doi:10.5157/NEPS:SC5:16.0.0) ** known issues: StudyStatesAddon: - this dataset will be possibly provided in another major release to analyze target persons with multiple episodes within waves Basics: - use this dataset just to get an idea of the sample, please! Do not analyze data using this dataset! pTargetCAWI: - the variable name tg69225 of the item "T-Structure (class management): monitoring 5" indicates a wrong construct affiliation. The name should be tg69725 as it is the fifth item belonging to the construct "monitoring" of which all other associated items start with "tg6972 *". The name will be corrected with the next release. -----* Changes introduced to NEPS:SC5 by version 16.0.0 * in general: - disagint has been extended for revoked episode in check module - removed generated date variables (*12?_g1, *11?_g1) in spell datasets - new variable names for intm inty intd corresponding to those in SC3 and SC4 - information taken from CATI/spVocTrain and CAWI can differ and show inconsistencies, this is due to the nature of survey data New dataset added: spVocBreaks - data on vocational breaks were extracted from spVocTrain - break episodes were reshaped to long format - break episodes were cleaned (closing small gaps, drop breaks within breaks) StudyStates: added new variables (including CAWI-information) >> type of higher education institution (CATI/spVocTrain): tx92401 tx92402 >> type of higher education institution (CAWI): tx92403 >> Status of the course of study (CAWI): tx24101>> Higher education institution ID (CAWI): tx24013 >> Highest vocational qualification (CAWI): tx15316

Multiple episodes in waves cause problems to generate states, change of subjects, study breaks etc.

pTargetCATI:

- versionized variable t514008_v1 for wave 7 had to be added. In this wave all target persons were asked about their satisfaction with the course of studies - no matter if they were actually studying at this time. Because of this, interviewers were instructed to tell non-students to choose the " does not apply"-button. From wave 9 on, non-students were filtered and interviewer instructions were removed.
- job related information on social capital was implemented in wave 15 by mistake. After this problem has been recognized, the items were removed from the instrument. As only a part of the respondents answered these questions, this variables were removed for wave 15. These variables are: t324110 t32411k t32411l t32411m t32411o t32411p t32411q t32411n t32411r t32411s t32411b t32411d t32411e t32411c t325110 t32511k t32511l t32511m t32511o t32511p t32511q t32511n t32511r t32511s t32511b t32511d t32511e t32511c
- variable t515039_g1 "Job characteristics: episode number of employment episode" was generated for wave 15 for merging information from spEmp to information on job characteristics in pTargetCATI

pTargetCORONA :

- information on satisfaction with course of study, school or apprenticeship (variable t514010) is also available for those respondents who previously indicated that they were employed, although the response option "does not apply" was available; in these cases, it is unclear what respondents were referring to with their answer to the satisfaction question, so this variable should be treated with caution, depending on the research question
- information on satisfaction with work (variable t514009) is also available for those respondents who previously indicated that they studied, were in vocational educational training or did nothing, although the response option "does not apply" was available; in these cases, it is unclear what respondents were referring to with their answer to the satisfaction question, so this variable should be treated with caution, depending on the research question
- information on health-related limitations are also available in pTargetCORONA (variable t521055). However, some respondents indicated very strong or strong limitations; they stated that they are in good or very good physical and mental health, though. In these cases, it is unclear, how respondents interpreted the question regarding limitations in daily activities.

```
spPartner:
       - variables ts31226_g1 - ts31226_g16 "Partner: Profession" were edited. If the
          partner's profession did not change since the last interview (th32369 ==
          1), the code -29 "Value from the last sub-episode" was implemented.
-----
* Changes introduced to NEPS:SC5 by version 14.1.0 *
New dataset added: pTargetCORONA
pTargetCATI:
       - added values to t751004_g* - variables for wave 12 and wave 13
_____
* Changes introduced to NEPS:SC5 by version 14.0.0 *
General remarks:
       - some variable labels were corrected
       - supplemental meta information on several variables was added
CohortProfile:
       - further smoothing on ID_i has been done (should be stable ... by now)
MethodsCAWI:
       - variables added: tx80102, tx80103, tx80210, tx80310
       - erroneously added waves 6 and 8 were dropped
pTargetCATI:
       - versionized variables had to be added due to change of item-battery-
          composition:
             tg51101_v1, tg51102_v1, tg51103_v1, tg51104_v1, tg51108_v1, tg51109_v1,
                  tg51110_v1, tg51111_v1, tg51112_v1, tg51113_v1, tg51114_v1,
                 tg51115_v1, tg51116_v1, tg51117_v1, tg51118_v1
      - polarity of categories in variable t428050 have changed since wave 11 as the
          field instrument was edited in this way:
             - in releases prior to version 11 the coding was: 1="not at all"; 2="
                 hardly"; 3="average"; 4="strongly"; 5="very strongly";
             - since version 11 the coding is: 1="very stringly"; 2="strongly"; 3="
average"; 4="hardly"; 5="not at all"
spPartner:
      - variable ts31410 corrected
xEcoCAPI:
       - plausible values for competency data were added
xPlausibleValues:
       - new dataset since release 13-0-0: provides plausible values for competency
          data stored in xTargetCompetencies
* Changes introduced to NEPS:SC5 by version 13.0.0 *
******
Known issues:
```

MethodsCA comp	AWI: waves 6 and 8 erroneously added to MethodsCAWI, data in those lines is oletely missing, please drop these waves
******	**********************
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
CohortPr	ofile: - some checks on plausibility and smoothing on ID_i has been done
EditionsB	Backup: – new dataset since release 12–0–0: provides original data prior to coding and smoothing during the process of data preparation
pTargetC/	ATI: - 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
spVocTra	in: - 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
======================================	s 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 droppped, except for intro-variables in spChild and spPartner.
pTargetC/	 ATI: 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.
pTargetCA	\WI: - for several variables open answers were (belatedly) coded.
spChild :	 variable ts33204_g1 was generated to provide information on the status of th child. Therefore category "other child in household" was added.
spEmp :	

	match the corresponding variable's name in other starting cohorts.
spSchoolE	xtExam :
-	additional information on external examinations from wave 1 and 3 was gathered from file spSchool.
spVocExtE	xam :
-	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.
spVocTrai	n:
-	variables tg24162_g1, tg24165_g1 and tg24168_g1 were generated to provide information on major or minor subjects for each subspell of an episode.Fo 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
-	variable tg01003_ha was edited and now excludes administration and business
-	servicevariables with information on subject of studies (tg2417*) were
a	
General r	
_	emarks:
-	emarks: several variables surveyed have been renamed to *_v1 and *_v2 in prior releases:
-	emarks: 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
-	emarks: 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).
-	emarks: 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).
- CohortPro	emarks: 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). ofile: testy testm testd erroneously coded to -56 for testing data in wave 7 have
- CohortPro -	emarks: 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). ofile: testy testm testd erroneously coded to -56 for testing data in wave 7 have now been coded with correct dates
- CohortPro - pTargetCA\	emarks: 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). ofile: testy testm testd erroneously coded to -56 for testing data in wave 7 have now been coded with correct dates NI:
- CohortPro - pTargetCA\ -	emarks: 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). ofile: testy testm testd erroneously coded to -56 for testing data in wave 7 have now been coded with correct dates NI: there have been changes during the field phase regarding interviewer
- CohortPro - pTargetCA\ -	<pre>emarks: 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). ofile: testy testm testd erroneously coded to -56 for testing data in wave 7 have now been coded with correct dates NI: there have been changes during the field phase regarding interviewer instructions in variable "tg51001"; the now indicator variable "Version ts51001" contains information about the set of the se</pre>
– CohortPro – pTargetCA\ –	emarks: 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). offile: testy testm testd erroneously coded to -56 for testing data in wave 7 have now been coded with correct dates NI: there have been changes during the field phase regarding interviewer instructions in variable "tg51001"; the new indicator variable "Version_tg51001" contains information abo the version of the survey instrument
- CohortPro - pTargetCAN - MethodsCA	emarks: 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). ofile: testy testm testd erroneously coded to -56 for testing data in wave 7 have now been coded with correct dates MI: there have been changes during the field phase regarding interviewer instructions in variable "tg51001"; the new indicator variable "Version_tg51001" contains information abo the version of the survey instrument MI:
- CohortPro - pTargetCA\ - MethodsCA -	<pre>emarks: 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). ofile: testy testm testd erroneously coded to -56 for testing data in wave 7 have now been coded with correct dates NI: there have been changes during the field phase regarding interviewer instructions in variable "tg51001"; the new indicator variable "Version_tg51001" contains information abo the version of the survey instrument WI: a new data file including para data from the CAWI interviews has been added</pre>
- CohortPro - pTargetCA) - MethodsCA -	<pre>emarks: 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).</pre> offile: testy testm testd erroneously coded to -56 for testing data in wave 7 have now been coded with correct dates NI: there have been changes during the field phase regarding interviewer instructions in variable "tg51001"; the new indicator variable "Version_tg51001" contains information abo the version of the survey instrument WI: a new data file including para data from the CAWI interviews has been added
– CohortPro – pTargetCA\ – MethodsCA –	<pre>emarks: 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).</pre>

Appendix

General remarks: - several variables surveyed prior to wave 10 have been renamed to st _v1 and st_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-shortquestionaire, 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 concering 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

 ${\tt CohortProfile:}$

- 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 questionaire 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 questionaire 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 chidl 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 nepsmiss 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;

<u>185</u>

use for orientation, not for analyses!

xInstitution:

 dataset containing detailed information on the targets' institutions has been added for onsite access in Bamberg

B.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

Label Generation status Text -54 missing by design	
Text -54 missing by design	
-54 missing by design	
 no migrant background 	
1 1st generation	
2 1.5th generation	
3 2nd generation	
4 2.25th generation	
5 2.5th generation	
6 2.75th generation	
7 3rd generation	
8 3.25th generation	
9 3.5th generation	
10>>3.75th generation	

	« t400500_g2v1 pTargetCATI t400500_g2 »
Label	Missing/contradicting imformation about country of birth for generation
	status
Text	
-54	missing by design
1	Unique assignment possible
2	Information for target person unknown
3	Information for one parent unknown
4	Information for both parents unknown
5	Information for one grandparent unknown
6	Information for two grandparents unknown
7	Information for three grandparents unknown
8	Information for four grandparents unknown
9	no assignment to a generation status possible

Appendix

	« t400500_g3v1 pTargetCATI t400500_g3 »
Label	Group of origin
Text	
-54	missing by design
1	Germany
2	Italy
3	Poland
4	Romania
5	Turkey
6	Former Yugoslavia
7	Former Soviet Union
8	Central and South Amerika, Carribean
9	Northern and Western eurospe
10	North America
11	Oceania/Polynesia
12	other country of Middle East and North Africa
13	other country of Africa
14	other country of Asia
15	other Central and Eastern Europe
16	other Southern Europe
17	abroad, but cannot be assigned to a specific group of origin

	<pre>« t514008_v1</pre>	pTargetCATI t514008 »
Label	« Sat	sfaction with higher education
	» Sat	sfaction with course of study
Text	Но	v satisfied are you with your higher education?
-98	doi	i't know
-97	ref	ised
-93	doe	es not apply
-54	mis	sing by design
Θ	COR	npletely dissatisfied
1	1	
2	2	
3	3	
4	4	
5	5	
6	6	
7	7	
8	8	
9	9	
10	COr	npletely satisfied

	« t	:516201_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

Appendix

	~~	t516202_g1v1 pTargetCATI t516202_g1 »
Label		Party election (another party)
Text		Which other party is this?
-98		don't know
-97		refused
-55		not determinable
-54		missing by design
-52		implausible value removed
1		Citizens' initiatives
2		Voter participation with invalid vote
3		Indifferent
4		Already disbanded parties (graue Panther 2008, Deutsche Bierunion after 1990, SDP after 1990)
5		APPD - Anarchistic pogo party Germany
6	«	AUF-Party, for work, environment, family (Christain)
	>>	AUF-party, for work, environment, family (Christian)
7		Bayernpartei - Bavarian party
8		PBC - Party of Bible-abiding Christians
9		BIG - Alliance for innovation and justice (party of Germans of Turkish origin)
10		Die Frauen - Feminist party The Women
11		Die Freien Wähler - The Free Voters
12		Die Freiheit - Civil rights' party for more freedom and democracy
13		Die PARTEI - Party for Labour, Rule of Law, Animal Protection, Promotion of Elites and Grassroots Democratic Initiative
14		Die Tierschutz-Partei - Party Humans Environment Animal Rights
15		Die Violetten - The Purples
16		DKP - German Communist Party
17		FAMILIE - Family Party of Germany
18	~~	MLPD - Marxist-Leninist Party of Germany
	>>	MLPD - Marxist–Leninist Party of Germany
19		ÖDP - Ecological Democratic Party
20		PDV - Party of Reason
21	«	Pro NRW - Pro North-Rhine Westphalia
	>>	Pro NRW
22		SSW - South Schleswig Voters' Association
23	«	AfD - Alternative for Germany
24	~~	Bündnis 21/RRP - Pensioners' Party
	>>	Bündnis 21/RRP - Rentnerinnen- und Rentner-Partei [Pensioners' Party]
25		KPD - Communist Party of Germany
26		UDP - Union of German Patriots
27	>>	Alfa (Allianz für Fortschritt und Aufbruch [alliance for progress and awakening]) since 2016 LKR (Liberal-Konservative Reformer [liberal-conservative reformers])
28	>>	Deutsche Mitte - German Middle
29	>>	V-Partei (Party for Change, Vegetarians and Vegans)
30	>>	Liberale (New Liberals - The Social Liberal)
31	>>	Die Republikaner - The Republicans

	« t525	008_v1 pTargetCATI t525008 »
Label		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	«	have 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

	« t525 2	09_v1 pTargetCATI t525209 »
Label	~~	Alcohol consumption
	>>	Alcohol consumption frequency last 12 months
Text	~	How often do you consume alcoholic drinks?
	>>	How often do you consume alcoholic drinks? Think about the average over 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

	« tg2450a	a_v1 pTargetCATI tg2450a »
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 your doctoral program. 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 free doctoral student without institutional integration. In which institutional context are you currently doing your doctorate?
-98		don't know
-97	«	refused
-92	«	question erroneously not asked
-54		missing by design
-52	«	implausible value removed
-20	>>	none of it
Θ		not specified
1		specified

	« tg245	0b_v1 pTargetCATI tg2450b »
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 your doctoral program. 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 free doctoral student without institutional integration. In which institutional context are you currently doing your doctorate?
-98		don't know
-97	«	refused
-92	«	question erroneously not asked
-54		missing by design
-52	«	implausible value removed
-20	>>	none of it
Θ		not specified
1		specified

	« tg2	450c_v1 pTargetCATI tg2450c »
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 your doctoral program. 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 free doctoral student without institutional integration. In which institutional context are you currently doing your doctorate?
-98		don't know
-97	«	refused
-92	«	question erroneously not asked
-54		missing by design
-52	«	implausible value removed
-20	>>	none of it
Θ		not specified
1		specified

	« tg245	<mark>0d_v1</mark> pTargetCATI tg2450d »
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 your doctoral program. 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 free doctoral student without institutional integration. In which institutional context are you currently doing your doctorate?
-98		don't know
-97	~~	refused
-92	~~	question erroneously not asked
-54		missing by design
-52	~	implausible value removed
-20	>>	none of it
Θ		not specified
1		specified

	« tg2450	e_v1 pTargetCATI tg2450e »
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 your doctoral program. 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 free doctoral student without institutional integration. In which institutional context are you currently doing your doctorate?
-98		don't know
-97	«	refused
-92	«	question erroneously not asked
-54		missing by design
-52	«	implausible value removed
-20	>>	none of it
Θ		not specified
1		specified

	<pre>« tg2450</pre>	of_v1 pTargetCATI tg2450f »
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 your doctoral program. 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 free doctoral student without institutional integration. In which institutional context are you currently doing your doctorate?
-98		don't know
-97	«	refused
-92	«	question erroneously not asked
-54		missing by design
-52	«	implausible value removed
-20	»>	none of it
Θ		not specified
1		specified

	« tg24	.50g_v1 pTargetCATI tg2450g »
Label	«	Doctorate context - alongside studies
	>>	Doctorate context - while studying
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 your doctoral program. 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 free doctoral student without institutional integration. In which institutional context are you currently doing your doctorate?
-98		don't know
-97	«	refused
-92	«	question erroneously not asked
-54		missing by design
-52	«	implausible value removed
-20	>>	none of it
Θ		not specified
1		specified
	« tg24	50h_v1 pTargetCATI tg2450h »
Label	« tg24 «	50h_v1 pTargetCATI tg2450h » Doctorate context - without institutional integration
Label	« tg24 « »	50h_v1 pTargetCATI tg2450h » Doctorate context - without institutional integration Doctorate context - without institutional integration, free doctorate student
Label Text	« tg24 « »	 boh_v1 pTargetCATI tg2450h » Doctorate context - without institutional integration Doctorate context - without institutional integration, free doctorate student [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?
Label Text	« tg24 « » «	 bbh_v1 pTargetCATI tg2450h » Doctorate context - without institutional integration Doctorate context - without institutional integration, free doctorate student [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 your doctoral program. 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 free doctoral student without institutional integration. In which institutional context are you currently doing your doctorate?
Label Text	« tg24 « » «	 by the second second
Label Text -98 -97	<pre></pre>	 bb-v1 pTargetCATI tg2450h » Doctorate context - without institutional integration Doctorate context - without institutional integration, free doctorate student [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 your doctoral program. 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 free doctoral student without institutional integration. In which institutional context are you currently doing your doctorate? don't know refused
Label Text -98 -97 -92	<pre></pre>	 bb-v1 pTargetCATI tg2450h » Doctorate context - without institutional integration Doctorate context - without institutional integration, free doctorate student [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 your doctoral program. 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 free doctoral student without institutional integration. In which institutional context are you currently doing your doctorate? don't know refused question erroneously not asked
Label Text -98 -97 -92 -54	<pre></pre>	 botorate context - without institutional integration Doctorate context - without institutional integration, free doctorate student [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 your doctoral program. 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 free doctoral student without institutional integration. In which institutional context are you currently doing your doctorate? don't know refused question erroneously not asked missing by design
Label Text -98 -97 -92 -54 -52	<pre></pre>	 botorate context - without institutional integration Doctorate context - without institutional integration, free doctorate student [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 your doctoral program. 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 free doctoral student without institutional integration. In which institutional context are you currently doing your doctorate? don't know refused question erroneously not asked missing by design implausible value removed
Label Text -98 -97 -92 -54 -52 -20	<pre>« tg24 « » « « » « « « « « « « « « « « « » »</pre>	 b50h_v1 pTargetCATI tg2450h » Doctorate context - without institutional integration Doctorate context - without institutional integration, free doctorate student [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 your doctoral program. 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 free doctoral student without institutional integration. In which institutional context are you currently doing your doctorate? don't know refused question erroneously not asked missing by design implausible value removed none of it
Label Text -98 -97 -92 -54 -52 -20 0	<pre></pre>	 b50h_v1 pTargetCATI tg2450h » Doctorate context - without institutional integration Doctorate context - without institutional integration, free doctorate student [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 your doctoral program. 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 free doctoral student without institutional integration. In which institutional context are you currently doing your doctorate? don't know refused question erroneously not asked missing by design implausible value removed none of it not specified

	« tg2	450i_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 your doctoral program. In the following we would like to ask you some questions. A doctoral thesis ca 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 free doctoral student without institutional integration. In which institutional context are you currently doing your doctorate?
-98		don't know
-97	«	refused
-92	«	question erroneously not asked
-54		missing by design
-52	«	implausible value removed
-20	>>	none of it
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
	>>	[AUX] Auxiliary variable: Teaching groups, current status
-92	>>	question erroneously not asked
-54		missing by design
Θ		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
	>>	completed teaching degree course and intended Referendariat [period as a trainee teacher] or completed teaching degree course and intended employment as a teacher
3		ongoing teaching Referendariat
4		completed Referendariat [period as a trainee teacher] and employment as a teacher is intended
5	~~	employment as teacher
	>>	employment as a teacher
6	>>	interrupted employment as a teacher (e.g. due to parental leave)

	« tg60	022_v1 pTargetCATI tg60022 »
Label	<<	Screening employed teachers
	>>	Screening Employed teachers
Text	«	As announced at the beginning of the interview, there is a survey section for student teachers, trainee teachers and teaching staff Just to be on the safe side and to simplify the interview process, would you please tell me briefly whether you are currently employed as a teacher?
	*	As announced at the beginning of the interview, the interview includes a survey section for student teachers, Referendare [trainee teachers] and teachers. Just to make sure and to simplify the interview process, please tell me briefly if you are currently employed as a teacher.
-98	«	don't know
-97	~	refused
-93	~	does not apply
-92	«	question erroneously not asked
-54		missing by design
-52	<<	implausible value removed
1		yes
2		no
3	>>	yes, but I have currently interrupted my employment as a teacher

	« tg60031_v1 pTargetCATI tg60031 »
Label	Preload Completed teaching degree course
	» Preload completed teacher education program (CATI), as of 13th wave
Text	[AUTO] Preload Completed teaching degree course
-54	missing by design
Θ	no teaching degree course completed
1	teaching degree course completed

	« th21	1300_v1 pTargetCATI th21300 »
Label		Number of selected courses
Text	«	[AUTO]: Via random selection, select two courses that were completed between <intmpre intmjpre=""> and <20102(intm/intj)></intmpre>
	>>	[AUTO]: Randomly select one of the courses that was completed from <intmpre intmjpre=""> to <20102(intm/intj)></intmpre>
-54		missing by design
Θ		no course selected
1		1 course selected
2	~~	2 courses selected

Appendix

	«	ts15911_v1 pTargetCATI ts15911 »
Label		Graduate
	>>	Auxiliary variable: highest degree
Text		[AUX]
	>>	[AUX] Highest degree
-54		missing by design
0		No higher education qualification
1		BA, MA, Diploma, state examination
	>>	BA
2	~~	Doctorate
	>>	MA, Diploma, state examination
3	>>	Doctorate

pTargetCAWI

	« t24240	0_g2v1 pTargetCAWI t242400_g2 »
Label	«	AUX: subject group ref subject questions about learn. env. (destatis 2010/11)
	>>	Subject group reference subject learning environment (destatis 2010/11)
Text	«	[AUTO] Auxiliary variable: Field of study referenced for questions about learning env.
	>>	The following is about your experiences in your current course of study. If you are studying several subjects, these can be very different, e.g. in terms of content and/or organization of teaching. For this reason, we ask you to select the main or teaching subject to which you would like to refer to in the next questions.
-99		filtered
-97		refused
-96		not in list
-92	«	question erroneously not asked
-91		survey aborted
-55	>>	not determinable
-54		missing by design
-29		Value from the last sub-episode
-28		Value from recruitment pTargetCATI
-20		no further subject
1		Linguistic and cultural studies
2		Sport
3		Law, economics and social science
4		Mathematics, sciences
5		Human medicine/health sciences
6		Veterinary medicine
7		Agricultural-, forest- and nutrition sciences
8		engineering
9		Arts, art science
10		Outside the study area structure

	« t24240	0_g5v1 pTargetCAWI t242400_g5 »
Label	~	AUX: ISCED-97 ref subject questions about learning environment (1-digit level)
	>>	ISCED-97 reference subject learning environment (1-digit level)
Text	«	[AUTO] Auxiliary variable: Field of study referenced for questions about learning env.
	>>	The following is about your experiences in your current course of study. If you are studying several subjects, these can be very different, e.g. in terms of content and/or organization of teaching. For this reason, we ask you to select the main or teaching subject to which you would like to refer to in the next questions.
-99		filtered
-98		don't know
-97		refused
-96		not in list
-92		question erroneously not asked
-91		survey aborted
-55		not determinable
-54		missing by design
Θ		general educational programs
1		Education
2		Humanities and Arts
3		Social sciences, Business and Law
4		Natural sciences, mathematics and computer science
5		engineering, manufacturing and construction
6		Agriculture and Veterinary
7		Health and welfare
8		Services
9		not known or unspecified

	« t289	900_v1 pTargetCAWI t289900 »
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?

	« tg511	01_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
Θ		not specified
1		specified

	<pre>« tg51102</pre>	2_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
Θ		not specified
1		specified

	« tg5110 :	3_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
Θ		not specified
1		specified

	« t	g51104_v1 pTargetCAWI tg51104 »
Label	~~	Curr. activity: vocational training
	>>	Voc. train./further educ.: vocational training
Text	**	[MF] Which of the following activities are your 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 both
Θ		not specified
1		specified

	« tg5110	8_v1 pTargetCAWI tg51108 »
Label	«	Curr. activity: retraining or further education
	>>	Voc. train./further educ.: retraining, further education
Text	«	[MF] Which of the following activities are your 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 both
Θ		not specified
1		specified

	« tg51109	9_v1 pTargetCAWI tg51109 »
Label	~~	Curr. activity: (voluntary) services,
		(military/alternative/community/social)
	>>	other activities: volunt. military service/social year/fed. volunt. service
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
Θ		not specified
1		specified

	« tg5111(9_v1 pTargetCAWI tg51110 »
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
Θ		not specified
1		specified

	« tg5111	1_v1 pTargetCAWI tg51111 »
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

	« tg51112	2_v1 pTargetCAWI tg51112 »
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		nitered
-97		refused
-92		question erroneously not asked
-91		survey aborted
-54		missing by design
Θ		not specified
1		specified

	« tg5111	3_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
Θ		not specified
1		specified

	« tg51114	4_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
Θ		not specified
1		specified

	« tg511	<mark>15_v1 pTargetCAWI tg51115</mark> »
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
Θ		not specified
1		specified

	« tg5111(6_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
Θ		not specified
1		specified

	« tg5111	7_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
Θ		not specified
1		specified

	« tg5111	8_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
Θ		not specified
1		specified

	« tg5130	0_v1 pTargetCAWI tg51300 »
Label	~	Change of subject since starting university
	>>	Change of subject 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)>?</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

	« tg5140(0_v1 pTargetCAWI tg51400 »
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 state examination)?</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

	« tg5150	0_v1 pTargetCAWI tg51500 »
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)>?</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

	« ts232	228_v1 spEmp ts23228 »
Label	«	Type of education required
	>>	Type of required training
Text		What kind of training is usually required to do this job?
-98		don't know
-97		refused
-92	«	question erroneously not asked
-55		not determinable
-54		missing by design
1		no qualification
2		a training on the job
3		a completed vocational training
4		a completed training at a Fachschule
5		a master craftsman's/craftswoman's certificate or technician certificate
6	~	a completed higher education qualification (university of applied
		sciences or university)
7		a doctorate or habilitation
8	>>	a Bachelor's degree (university of applied sciences or university)
9	>>	a Master's degree or state examination, a diploma or a Magister
		(degrees from a university of applied sciences or university)

« ts2390	1_v1 spEmp ts23901 »
	Auxiliary variable: current employment
	[AUX] Auxiliary variable Current employment
~	implausible value
	not determinable
	missing by design
~	currently employed
>>	Current employment
~	employed within the last year, but not currently
>>	Completed employment
~~	not employed within the last year / end not assignable
	« ts2390 « « » « » « » «

	~	ts23911_v1 spEmp ts23911 »
Label		Auxiliary variable: type of employee
Text		[AUX] Employee type
-55		not determinable
-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	~~	Freelancer
7	~~	Self-employed
8	~~	Positions in an assisting capacity
9	~~	Vocational training positions
13	>>	semi-skilled or unskilled work/student assistant
14	>>	private student tuition/homework supervision
spInternship

	« tg3611	1_v1 spInternship tg36111 »
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
-55		not determinable
-54		missing by design
-21		no fixed working hours
-20	~~	more than 50 hours per week
	>>	more than 90 hours per week

spPartner

	« ts3	1223_v1 spPartner ts31223 »
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
-55		not determinable
-54		missing by design
1	«	primarily working
	>>	full-time employed
2	>>	part-time employed
3	«	part-time employed
	>>	in a side job
4		unemployed

	« ts31510	∂_v1 spPartner ts31510 »
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
-55	«	not determinable
-54		missing by design
1		divorced/civil partnership annulled
2		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

	« ts15	304_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 apprenticeship (administrative, company-based, industrial, agricultural), journeyman's/journeywoman's certificate or apprenticeship certificate (craft certificate), dual vocational training
	"	agricultural), journeyman's/journeywoman's certificate or apprenticeship certificate (craft certificate); dual vocational training
2		graduation from a school of public health
3		certificate from a Berufsfachschule [vocational school] or a Handelsschule [type of vocational school for commercial professions]
4		certificate from another Fachschule
5	«	master craftsman's/craftswoman's diploma
	>>	Master craftsman's/craftswoman's diploma
6		technician certificate
7	~~	diploma
8	~~	Bachelor
9	~~	Master
10	~~	diploma from university of applied sciences (Dipl(FH))
	>>	diploma from a university of applied sciences (Dipl(FH))
11	~~	Diploma from a university
	>>	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 state examination (in teaching)
18		first state examination (not in teaching)
19	**	second/third state examination
	>>	second/third state examination (not in teaching)
20	~~	Doctorate
	>>	doctorate
21	~~	Habilitation
	>>	habilitation
22		medical specialist
23		civil service examination for the subcierical class
24		civil service examination for the cierical class
25		civil service examination for the executive class
26		civil service examination for the administrative class
27		other qualification
28 29		other qualification other higher education qualification (e.g. ecclesiastical examination, artistic examination)
30	>>	second state examination (in teaching)

spVocTrain

	« tg24205_v1 spVocTrain tg24205 »
Label	Point of time decision for Master
Text	When did you make the decision for your Master's degree program?
-55	not determinable
-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

	« ts3122	1_v1 spPartner ts31221 »
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

	« ts3122	3_v1 spPartner ts31223 »
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

	« ts3122	24_v1 spPartner ts31224 »
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

	« ts31225_v1 spPartner ts31225 »
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

	<pre>« ts31227</pre>	7_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

	« ts312	28_v1 spPartner ts31228 »
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, journeyperson [trained craftsperson]
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

	« ts312	30_v1 spPartner ts31230 »
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

	«	ts31410_v1 spPartner ts31410 »	
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	

	« ts3141m	n_v1 spPartner ts3141m »	
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	

	« ts3141y_v1 spPartner ts3141y »	
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	

	« ts31510_v1 spPartner ts31510 »
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

	« ts3151m_v1 spPartner ts3151m »	
Label	 Date of partner's death (month) 	
	» Date of death Partner (month)	
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	

	« ts3151y_v1 spPartner ts3151y »	
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	

	« ts3152	m_v1 spPartner ts3152m »
Label		Date of moving apart (Month)
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?
-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

	« ts3152	<mark>'y_v1 spPartner ts3152y</mark> »	
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

	« tg	24146_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 o	Point of time decision for master		
Text	When o	When did you make the decision for your master degree program?		
-54	Missing	Missing by design		
1	before	before starting the previous higher education program		
2	During	During the previous higher education program		
3	after er	nding the previous h	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	L_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)

	« tg2452	2m_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

	~~	tg2452y_v1 spVocTrain tg2452y »
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