

# 1649<sup>2026</sup>

**SOEP** Survey Papers  
Series D - Variable Descriptions and Coding

## SOEP-Core v41 – Activity Biography in the Datasets "pbiospe" and "artkalen"

Jascha Dräger

Running since 1984, the German Socio-Economic Panel (SOEP) is a wide-ranging representative longitudinal study of private households, located at the German Institute for Economic Research, DIW Berlin.

The aim of the SOEP Survey Papers Series is to thoroughly document the survey's data collection and data processing.

The SOEP Survey Papers is comprised of the following series:

- Series A** – Survey Instruments (Erhebungsinstrumente)
- Series B** – Survey Reports (Methodenberichte)
- Series C** – Data Documentation (Datendokumentationen)
- Series D** – Variable Descriptions and Coding
- Series E** – SOEPmonitors
- Series F** – SOEP Newsletters
- Series G** – General Issues and Teaching Materials

The SOEP Survey Papers are available at <http://www.diw.de/soepsurveyspapers>

**Editors:**

Dr. Jan Goebel, DIW Berlin

Dr. Christian Hunkler, DIW Berlin

Prof. Dr. Philipp Lersch, DIW Berlin and Humboldt-Universität zu Berlin

Dr. Levent Neyse, DIW Berlin and Berlin Social Science Center (WZB)

Prof. Dr. Carsten Schröder, DIW Berlin and Freie Universität Berlin

Prof. Dr. Sabine Zinn, DIW Berlin and Humboldt-Universität zu Berlin

Please cite this paper as follows:

Jascha Dräger. 2026. SOEP-Core v41 – Activity Biography in the Datasets "pbiospe" and "artkalen". SOEP Survey Papers 1649: Series D. Berlin: DIW Berlin / SOEP.



This work is licensed under a Creative Commons Attribution-ShareAlike 4.0 International License.  
© 2026 by SOEP

ISSN: 2193-5580 (online)

DIW Berlin  
German Socio-Economic Panel (SOEP)  
Anton-Wilhelm-Amo-Straße 58  
10117 Berlin, Germany

Contact: [soeppapers@diw.de](mailto:soeppapers@diw.de)

# SOEP-Core v41 – Activity Biography in the Datasets "pbiospe" and "artkalen"

Jascha Dräger

# Activity Biography in the Datasets “*pbiospe*” and “*artkalen*”

by Jascha Dräger

(based on earlier work by Rainer Pischner, Henning Lohmann, Marco Giesselmann, Mila Staneva, Paul Schmelzer, Maik Hamjediers, and Tim Eder)

*Pbiospe* and *artkalen* encompass activities over the life course and distinguish between spells in education, employment (full- and part-time employment as well as minor employment and registered unemployment), retirement, housekeeping, parental leave and others (see Table 6 for an overview). Please note that these spells do not capture transitions between educational institutions or jobs and employers; spells reflect a continuous status for instance in full-time employment regardless of potential job changes. In the yearly Individual Questionnaire, the respondents are also asked to report job changes between the previous and current years and this information can be added to *artkalen* using the *pgen* data files.<sup>1</sup>

The spell file *artkalen* is collected from Individual Questionnaires as a calendar-matrix of months of the previous year and respective statuses of 15 categories<sup>2</sup> (for example Question 141 in 2023.<sup>3</sup> The information from each annual Individual Questionnaire is then attached to the information of previous surveys in the *pkal* Data File. Then, Spells are generated for all consecutive months with the same category reported. The generated spell file starts with the year before the entrance into the sample and ends with a respondent’s last observation.

The spell file *pbiospe* is based on the information on activity status over the life course, which is collected as a matrix from every respondent answering the Biography Questionnaire (for example Question 202 in 2023). The observations start at the age of 15 and end at the current age (up to age 65). Information on activity status covers only the period up to the time the biography is collected. To update the ongoing occupational career in *pbiospe*, information from the yearly Individual Questionnaire is also used by aggregating the recorded spells from *artkalen* into yearly values.

## Content of the Datasets

Table 1 contains a list of all the variables in the datasets. The variables *begin* and *end* indicate the beginning and the end of a spell. These variables are age entries in *pbiospe* and month entries in *artkalen* starting with January 1983 as the first month. *PBIOSPE* does also include reference to calendar years, e.g. *beginy* and *endy*. The *spellnr* is a serial identifier of spells of each activity status of a given person. The variable *spelltyp* contains information on the activity

---

<sup>1</sup> To add this information, it is necessary to split each spell at the time point of an interview. Afterwards, reports of job changes can be merged at each respective time point.

<sup>2</sup> For more detailed information on the categories in the Biography and Calendar Questionnaire can be found in Table 6.

<sup>3</sup> For persons who were temporarily unavailable for interviewing, it is sometimes possible to fill in the gaps in their occupational status. If these persons fill out the additional questionnaire for temporary dropouts later on, we can use the information collected there (see files \$PLUECKE).

status during the spell, e.g., employed full-time or unemployed (For detailed information see Table 6)

Variables	Information in:	
	<i>pbiospe</i>	<i>artkalen</i>
cid	Original Household Number	Original Household Number
pid	Never Changing Person ID	Never Changing Person ID
spellnr	Serial Number of the Spell per Person	Serial Number of the Spell per Person
spelltyp	Type of Spell	Type of Spell
begin	Age Spell Begins	Month Spell Begins
end	Age Spell Ends	Month Spell End
beginy	Year Spell Begins	
endy	Year Spell Ends	
zensor	Censor Variable	Censor Variable
source		Originating Dataset
other_flag		Flags spells generated from open-ended responses for "other" spells
duplicate_flag		Flags copied spells
spellinf	Spell Construction Information	
erhebj	Survey Year Biography Data	
kalyear	First Observation in Calendar Data	
beginb*	Age Spell Begins, Xth Initial Biography Spell	
endb*	Age Spell Ends, Xth Initial Biography Spell	
beginyb*	Year Spell Begins, Xth Initial Biography Spell	
endyb*	Year Spell Ends, Xth Initial Biography Spell	
begink*	Age Spell Begins, Xth Initial Calendar Spell	
endk*	Age Spell Ends, Xth Initial Calendar Spell	
beginyk*	Year Spell Begins, Xth Initial Calendar Spell	
endyk*	Year Spell Ends, Xth Initial Calendar Spell	

Table 1: All Variables in the ARTKALEN and PBIOSPE Datasets; \*These variables list all the Spells that were combined into this Spell. As they are generated dynamically the number of variables varies.

Both the Biography and the Individual Questionnaire allow for multiple activity statuses for a given year or month. No concept of main activity is used. A common combination is, for instance, “housewife/-husband” and “working part-time”. There are several other plausible combinations, but also combinations that are less plausible. However, a list of valid combinations of activity statuses defined according to legal or similar constructs would need to be based on very strong assumptions. In addition—in case of the yearly matrix in the

Biography Questionnaire—activities are reported that took place in a calendar year in consecutive months, which makes it impossible to exclude combinations of activities. Therefore, no data cleaning is performed at this stage. Consequently, the data may contain information on more than one activity for a given point in time.

### Dealing with multiple (inconsistent) biography responses

Most respondents (97.79%) have answered the biography questionnaire only once. However, 2.0% of respondents have answered two biography questionnaires, 0.18% have answered three biography questionnaires, and 0.02% have answered four biography questionnaires. Some respondents report inconsistent spells across their repeated responses to the biography.

We always use the spells provided in the first biography questionnaire. We update this information with spells from more recent biography questionnaires if the initial biography contained gaps or more recent biographies contain spells happening after the first biography questionnaire. This updating is done separately by the year of the activity. Information from more recent biographies that are inconsistent with earlier biographies are discarded. This decision is based on the assumption that respondents remember more recent activities better than activities that are longer ago.

	Age 15	Age 16	Age 17	Age 18
School (1 <sup>st</sup> Biography; Age 17)				
Apprenticeship (1 <sup>st</sup> Biography; Age 17)		X	X	
School (2 <sup>nd</sup> Biography; Age 18)	X	X		
Apprenticeship (2 <sup>nd</sup> Biography; Age 18)		X	X	X
School (combined)	X			
Apprenticeship (combined)		X	X	X

Table 2: Combining multiple (partially inconsistent) biography responses of an imaginary respondent. Green cells are activities that are consistently reported in both biographies. Yellow cells are activities that update missing information from previous biographies. Red indicates inconsistently reported activities which are discarded.

Table 2 gives an example of how inconsistent cases are handled. The imaginary respondent answered the biography questionnaire a first time at age 17, and then again at age 18. At age 17, the respondent only reported that they did an apprenticeship at age 16 and 17. At age 18 the respondent reported to have attended school at ages 15 and 16, and an apprenticeship for ages 16 to 18. For this imaginary case, we would use the information on the activities at age 16 and 17 which were reported in the first biography and update it with the new information on activities at ages 15 and 18 provided in the second biography. However, we would ignore that the respondent reported to also have attended school at age 16 in the

second biography because this is inconsistent with the reported activity at age 16 in the first biography.

### Differentiation between Apprenticeship/Retraining and Minijob/Part-time spells in *artkalen*

For some activities, more detailed information was collected in later survey years.

Until 1998, participants were asked whether they did any kind of vocational training (*spelltyp* == 4), which includes both First Job Training, Apprenticeship and Continuing Education, Retraining. Since 1999, there are separate questions on whether participants did a First Job Training or Apprenticeship (*spelltyp* == 41) or whether they did Continuing Education or Retraining (*spelltyp* == 42). Thus, information on the more detailed activity is only available since January 1999 (month 193 since January 1983). To enable continuity, Spells with type 41 and 42 January 1999 are copied as type 4 Spells and expand the existing Spells if applicable. The spells that contain copies of more detailed spells are marked by the variable *duplicate\_flag*.

Likewise, until 2003, respondents were asked whether they are part-time employed or whether they did a 'Mini job', but these two activities were not differentiated (*spelltyp* == 3). Since 2004, there are separate questions on whether participants were part-time employed (*spelltyp* == 31) or whether they did a 'Mini job' (*spelltyp* == 32). Thus, information on the more detailed activity is only available since January 2004 (month 252 since January 1983). Again, to enable continuity, spells of type 31 or 32 are copied as type 3 spells and expand the existing spells if applicable. Table 3 gives an example of how this was done.

Before Duplication									
cid	pid	spelltyp	spellnr	begin	end	zensor	source	other_flag	duplicate_flag
868	8605	4	3	165	193	1	1	0	0
868	8605	41	4	193	206	1	1	0	0

After Duplication									
cid	pid	spelltyp	spellnr	begin	end	zensor	source	other_flag	duplicate_flag
868	8605	4	3	165	206	1	1	0	1
868	8605	41	4	193	206	1	1	0	0

Table 3: Process of Spell duplication for continuity

If you want to revert this change in *artkalen* you simply need to set the end variable for all Spells with *spelltyp* == 4 and *duplicate\_flag* == 1 to month 192 if they exceed this month. For Spells with *spelltyp* == 3 and *duplicate\_flag* == 1 you need to set the end variable to month 252 if they exceed it.

## Aggregation of 'Other' Spells into usable spell types

Open Field Response	Spelltyp in <i>artkalen</i>
"Short Work Hours"	2: Short Work Hours
"Maternity Leave"	7: Maternity Leave
"Not Working"	5: Unemployed
"Pre-Retirement"	6: Retirement
"1-Euro Job"	15: Minijob
"Social Aid"	5: Unemployed

Table 4: Aggregation of Spelltypes from Open-Field Responses

The calendar Questionnaire also includes an open field response option for the years from 1986 to 2020. Some responses from this field are used to generate useful Spell-Information. Table 4 depicts which information is used for which *spelltypes*. Spells that are created in this way are flagged using the variable *other flag*. There are 2923 Spells with aggregated 'other'-information currently in *artkalen*, this adds up to 0.6 percent of Spells.

## Censoring; Coding of the *zensor* variable

Right: Left:	Not censored	Censored missing	Censored before gap
Not censored	1	2	3
Censored missing	4	5	6
Censored after gap	7	8	9

Table 5: Overview of *zensor*-Variable; Note: (99) Gap Spells are marked with -2

Missing information on the beginning or end of a spell causes what is known as censoring problems. There are two types of missing data. First, data can be missing on periods outside the observation window (before the age of 15 and after the age of 65 in the case of *pbiospe* or before and after the panel participation in *artkalen*). Second, data can be missing on years within the observation window due to item non-response in particular years or due to temporary drop-outs (the latter applies to calendar information only). In this case, we speak of "gaps." There are nine different patterns (see Table 5).

## Aggregation of *artkalen* information into *pbiospe*

To enhance the data of the Biography Information *artkalen* data is aggregated into yearly values. *Spelltyp* in *pbiospe* is less granulated than the types in *artkalen*, so *spelltyp* is aggregated as depicted in Table 6. The category 3 "military / civilian service" denotes several activities, such as being imprisoned, serving in the military, doing former community service or any voluntary service (FSJ or Bundesfreiwilligendienst). If you are more interested in distinguishing these, it is necessary to trace down the information in the person data \$P (e.g. variable *pkal09a* in *bhp*) and to merge them onto *artkalen* or *pbiospe*.

The *artkalen* Spells are first converted into the new spelltypes (excluding gaps) and transformed to yearly begin and end variables<sup>4</sup>, before extending Spells if necessary to achieve a congruent set of Spells. A question that arises when merging the data is how to handle overlapping pieces of information. The basic principle is to assign a value of a given status in a given year if the status is recorded in the calendar or in the biography information or both. An example might help to illustrate this: the calendar records full-time employment for the years 2005 and 2007 while the biography records full-time employment for the period from 2000 up to 2006. The merged data from *pbiospe* contains a spell that begins in 2000 and ends in 2007.

	<i>pbiospe</i>	<i>artkalen</i>
1	School/University	School, University (8)
2	Apprenticeship/Training	Vocational Training (4), First Job Training, Apprenticeship (41), Continuing Education, Retraining (42)
3	Military/Civilian service	Military, Community Service (FSJ, BufDi, Zivildienst) (9)
4	Full-time employed	Full-Time Employment (1)
5	Part-time employed	Part-Time Employment/marginal employment (3), Second Job (11), Part-Time Employment (31), Mini-Job (up to 556€ (in 2025)) (32)
6	Unemployed	Unemployed (5)
7	House-Husband/Wife	House-Husband/Wife (10)
8	Retired	Retired (6)
9	Other	Other (12) Maternity Leave (7)
10	Short Work Hours	Short Work Hours (2)
99	Gap	Gap (99)

Table 6: Aggregation of Spelltypes from *artkalen* to *pbiospe*

However, the initial information is restored by including additional variables, which allows for alternative ways of merging the data (see below). The variables *spellinf*, *erhebj*, and *kalyear* contain general information on the sources of the information captured in each spell.

In total, *pbiospe* contains 684,422 spells from 134,185 respondents. For 16,802 respondents (12.5% of respondents in *pbiospe*) the spells are solely based on biography data. For 23,301 respondents (17.4%) the spells are solely based on calendar data. For 94,082 respondents (70.1%) the spells are based on both biography and calendar data.

Table 7 shows that the majority of spells are based on biography information only (52.55 percent). One-third of all spells (32.40 percent) are not observed in the Biography

<sup>4</sup> Note that this means that also when only 1 month of activity is reported in ARTKALEN the whole year will be recorded with the activity in PBIOSPE. Also, if this happens in consecutive years those Spells are expanded, so if a person has 1 month Retraining (14) Spells in March 2000, 2001, and 2002 respectively, this person is recorded with a Training (2) Spell in PBIOSPE that begins in 2000 and ends in 2002.

Questionnaire but only in the calendar data. The remainder of spells contain information from biography as well as calendar data. Usually, these spells combine one period observed in the Biography Questionnaire with a period observed in the calendar. Only 0.44 percent of the spells combine more than one period in any of the two sources (*spellinf*=4, 5 or 6).

<i>spellinf</i>	N	%	% cum.
Biography only	347,089	52.55	52.55
Calendar only	214,002	32.4	84.95
1 biography, 1 calendar spell	96,458	14.6	99.56
2+ biography, 1 calendar spell(s)	770	0.12	99.68
1 biography, 2+ calendar spell(s)	2,064	0.31	99.99
2+ biography, 2+ calendar spell(s)	75	0.01	100

Table 7: Construction of Spells in PBIOSPE

The variables *beginb\**-*endyk\** document the initial information from the two different sources and are probably not of interest to most users. However, based on these variables, users can fully separate the Biography data from the aggregated *artkalen* data. This is advisable if you want to use the more detailed *artkalen* information and combine it with the yearly information from *pbiospe* for earlier years only. The variable names indicate the “source” of the original information utilized (B: Biography -Questionnaire or K: calendar information from the yearly survey). As an example, we discuss one of the spells that combines information on more than one period from any of the two sources. The spell number 3 of person 9205 starts in 1983 and ends in 1994 (*spelltyp*=4: full-time employment). As the variable *spellinf* (=5) shows, this a spell that combines one period from the biography data with two periods from the calendar data. According to the biography data, the person worked full-time from 1983 (*beginyb1*) until 1992 (*endyb1*). There is overlapping information from the calendar data available from 1986 onwards (*kalyear*). According to these data, the person worked full-time from 1986 (*beginyk1*) to 1990 (*endyk1*) and from 1993 (*beginyk2*) to 1994 (*endyk2*). During the years 1991 and 1992, no full-time employment is recorded in the calendar data, which contradicts the information from the biography data.

pid	spellnr	spelltyp	beginy	endy	spellinf	erhebj	kalyear	beginyb1	endyb1	beginyk1	endyk1	beginyk2	endyk2
9205	3	4	1983	1994	5	1998	1986	1983	1992	1986	1990	1993	1994

Table 8: Example Observation in *pbiospe*

In *pbiospe*, no attempt is made to “resolve” such contradictions, as this would require rather strong assumptions. More important, such assumptions would differ according to the research question, which makes it even more difficult to provide a standard solution. Therefore, in such cases, we generate spells in the same manner as in less difficult cases, namely by combining the information from the calendar and the biography data. In the given example, this results in a full-time employment spell that starts in 1983 and ends in 1994. As mentioned above, there are very few spells that combine information on two or more periods

(*spellinf*=4, 5, 6). There are even fewer such spells where the period of overlap is as long as in this example, where the biography data was collected many years after the persons joined the survey (*erhebj*=1998, *kalyear*=1986). However, users who are interested in combining biography and calendar data in a different manner can use the variables *beginb*\*-*endyk*\* to fully separate the two types of data and to recombine the data on the basis of different rules of aggregation.

### Overlap between PBIOSPE and ARTKALEN

As stated above, the calendar information is used to update the biography information. However, there is also a certain overlap of the periods covered by the two types of data. This is shown in Table 9 It indicates, for persons included in *pbiospe*, the year in which the biography information was collected (variable *erhebj*). This year is usually also the last year for which biography information is available. The table also shows the first year recorded in the calendar data (variable *kalyear*). In most cases (61.1 percent), the earliest calendar information is available for the year before the biography interview. This is the case for persons who answered the Biography Questionnaire in their first year as survey respondents. The calendar in the Individual Questionnaire refers to the year before the survey. There is no overlap for 14.6 percent of the respondents, and an overlap of at least 2 years for 22.8 percent of the respondents. However, for 1.5 percent of respondents the earliest calendar information is only available at least one year after the biography. For these respondents, there is a gap between the information from the biography and the calendar. This pattern emerges almost exclusively for respondents from the samples M1, M3, M4, M5, and M6.

Year of Biography data collection	1+ year later	same year	1 year earlier	2+ years earlier	N
1984	0.0	0.1	100.0	0.0	10,993
1987	0.0	0.0	36.4	63.6	505
1988	0.0	0.0	100.0	0.0	164
1989	0.0	0.5	99.5	0.0	193
1990	0.0	0.0	100.0	0.0	180
1991	0.0	0.0	100.0	0.0	157
1992	0.0	0.0	8.4	91.6	3,930
1993	0.0	0.0	76.6	23.4	304
1994	0.1	0.1	97.8	2.0	922
1995	0.1	0.1	99.0	0.8	1,037
1996	0.0	0.2	97.3	2.5	483
1997	0.0	0.0	98.5	1.5	477
1998	0.5	0.2	98.1	1.2	415
1999	0.1	0.1	26.6	73.3	1,821
2000	0.0	0.0	90.2	9.8	235
2001	0.0	0.0	6.3	93.7	7,530
2002	0.0	0.2	48.1	51.7	526

2003	0.0	0.1	16.9	83.0	2,193
2004	0.0	0.0	68.6	31.4	433
2005	0.0	0.0	89.0	11.0	292
2006	0.0	0.0	92.2	7.8	217
2007	0.0	0.0	16.2	83.9	1,858
2008	0.0	0.0	68.9	31.1	309
2009	0.0	0.0	89.5	10.5	190
2010	0.0	4.1	79.2	16.8	7,887
2011	0.0	2.5	91.7	5.8	5,670
2012	0.0	3.9	95.6	0.6	2,205
2013	1.3	91.4	7.0	0.3	3,891
2014	1.3	40.0	57.7	0.9	447
2015	0.5	81.0	14.7	3.8	1,440
2016	17.0	76.2	6.1	0.7	3,398
2017	8.5	27.7	63.4	0.4	5,784
2018	5.1	14.1	79.4	1.4	1,647
2019	1.6	3.9	92.9	1.6	3,085
2020	4.3	13.0	81.8	1.0	3,411
2021	1.5	6.7	12.7	79.1	1,513
2022	0.1	0.5	95.8	3.6	7,762
2023	0.0	58.0	27.5	14.5	5,369
2024	0.0	0.0	82.2	17.8	5,209
Total	1.5	14.6	61.1	22.8	94,082

*Table 9: Overlap between the data from Biography vs Calendar Questionnaire; \*Year of Calendar data Collection.*